



**HAL**  
open science

# Image restoration with generalized Gaussian mixture model patch priors

Charles-Alban Deledalle, Shibin Parameswaran, Truong Q. Nguyen

► **To cite this version:**

Charles-Alban Deledalle, Shibin Parameswaran, Truong Q. Nguyen. Image restoration with generalized Gaussian mixture model patch priors. 2018. hal-01700082v1

**HAL Id: hal-01700082**

**<https://hal.science/hal-01700082v1>**

Preprint submitted on 3 Feb 2018 (v1), last revised 8 Jun 2018 (v2)

**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.

# Image restoration with generalized Gaussian mixture model patch priors

Charles-Alban Deledalle<sup>\*†</sup>, Shubin Parameswaran<sup>†</sup>, and Truong Q. Nguyen<sup>†</sup>

**Abstract.** Patch priors have become an important component of image restoration. A powerful approach in this category of restoration algorithms is the popular Expected Patch Log-likelihood (EPLL) algorithm. EPLL uses a Gaussian mixture model (GMM) prior learned on clean image patches as a way to regularize degraded patches. In this paper, we show that a generalized Gaussian mixture model (GGMM) captures the underlying distribution of patches better than a GMM. Even though GGMM is a powerful prior to combine with EPLL, the non-Gaussianity of its components presents major challenges to be applied to a computationally intensive process of image restoration. Specifically, each patch has to undergo a patch classification step and a shrinkage step. These two steps can be efficiently solved with a GMM prior but are computationally impractical when using a GGMM prior. In this paper, we provide approximations and computational recipes for fast evaluation of these two steps, so that EPLL can embed a GGMM prior on an image with more than tens of thousands of patches. Our main contribution is to analyze the accuracy of our approximations based on thorough theoretical analysis. Our evaluations indicate that the GGMM prior is consistently a better fit for modeling image patch distribution and performs better on average in image denoising task.

**Key words.** Generalized Gaussian distribution, Mixture models, Image denoising, Patch priors.

**AMS subject classifications.** 68U10, 62H35, 94A08

**1. Introduction.** Image restoration is the process of recovering the underlying clean image from its degraded or corrupted observation(s). The images captured by common imaging systems often contain corruptions such as noise, optical or motion blur due to sensor limitations and/or environmental conditions. For this reason, image restoration algorithms have widespread applications in medical imaging, satellite imaging, surveillance, and general consumer imaging applications. Priors on natural images play an important role in image restoration algorithms. Image priors are used to denoise or regularize ill-posed restoration problems such as deblurring and super-resolution, to name just a few. Classical image priors include Gibbs distributions [23], total variation which imposes Laplacian [53, 60] or hyper-Laplacian [30] prior on image gradients, and generalized Gaussian [33, 39, 6, 14] or scaled mixture of Gaussian [47] priors for wavelet or curvelet coefficients [3]. Alternatively, modeling the distribution of patches of an image (*i.e.*, small windows usually of size  $8 \times 8$ ) has proven to be a powerful solution. In particular, popular patch techniques rely on non-local self-similarity [5], Fields of experts [51], learned patch dictionaries [1, 19], sparse or low-rank properties of stacks of similar patches [9, 11, 32, 29], patch re-occurrence priors [36], Gaussian mixture model prior on image gradients [21] or image patches [62, 61, 26].

Of these approaches, a successful approach introduced by Zoran and Weiss [62] is to model patch priors of clean natural images using Gaussian Mixture Models (GMM). The agility of this model lies in the fact that a prior learned on clean image patches can be effectively

---

<sup>\*</sup>Institut de Mathématiques de Bordeaux, CNRS, Université de Bordeaux, Bordeaux INP, Talence, France ([charles-alban.deledalle@math.u-bordeaux.fr](mailto:charles-alban.deledalle@math.u-bordeaux.fr)).

<sup>†</sup>Department of Electrical and Computer Engineering, University of California, San Diego, La Jolla, USA ([cdeledalle@eng.ucsd.edu](mailto:cdeledalle@eng.ucsd.edu), [sparames@ucsd.edu](mailto:sparames@ucsd.edu), [tqn001@eng.ucsd.edu](mailto:tqn001@eng.ucsd.edu)).

employed to restore a wide range of inverse problems. It is also easily extendable to include other constraints such as sparsity or multi-resolution patches [57, 43]. The use of GMMs for patch priors make these methods computationally tractable and flexible. Although GMM patch prior is effective and popular, in this article, we argue that a mixture of Generalized Gaussian distributions (GGMM) is a better fit for image patch prior modeling. Compared to Gaussian models, Generalized Gaussian distributions (GGD) have an extra degree of freedom controlling the shape of the distribution and they encompass Gaussian and Laplacian models.

The superior patch prior modeling capability of a GGMM over a GMM is illustrated in Figure 5. The figure shows histograms of six orthogonal 1-D projections of a clustered subset of a patch database onto the eigenvectors of its covariance matrix. To illustrate the difference in the shapes ( $\nu$ ) and scales ( $\lambda$ ) of the distributions, we have chosen a few projections corresponding to both the most and the least significant eigenvalues. It can be seen that GGD is a better fit on the obtained histograms than a Gaussian model. Additionally, different dimensions of the patch follow a different GGD. Hence, it does not suffice to model all the feature dimensions of a patch database as Laplacian or Gaussian. Therefore, we propose to model patch priors as Generalized Gaussian Mixture models (GGMM) with a separate shape and scale parameters for each feature dimension. This differs from the recent related approach in [42] that considered GGMM where each component has a fix shape parameter for all directions.

**Organization.** After explaining the considered patch prior based restoration framework in Subsection 2.1 and our motivations in Subsection 2.2, we derive our GGMM based restoration scheme in Section 3. Unlike [42], that incorporates a GGMM prior in a posterior mean estimator based on importance sampling, we directly extend the maximum *a posteriori* formulation of Zoran and Weiss [62] for the case of GGMM priors. While such a GGMM prior has the ability to capture the underlying distribution of clean patches more closely, we will show that it introduces two major computational challenges in this case. The first one can be thought of as a classification task in which a noisy patch is assigned to one of the components of the mixture. The second one corresponds to an estimation task where a noisy patch is denoised given that it belongs to one of the components of the mixture. Due to the interaction of the noise distribution with the GGD prior, we first show in Section 3 that these two tasks lead to a group of one-dimensional integration and optimization problems, respectively. Specifically, for  $x \in \mathbb{R}$ , these problems are of the following forms

$$(1) \quad \int_{\mathbb{R}} \exp\left(-\frac{(t-x)^2}{2\sigma^2} - \frac{|t|^\nu}{\lambda_\nu^\nu}\right) dt \quad \text{and} \quad \operatorname{argmin}_{t \in \mathbb{R}} \frac{(t-x)^2}{2\sigma^2} + \frac{|t|^\nu}{\lambda_\nu^\nu} .$$

for some  $\nu > 0$ ,  $\sigma > 0$  and  $\lambda_\nu > 0$ . These two problems are studied in Section 4 and Section 5, respectively. In general, they do not admit closed-form solutions but some particular solutions or approximations have been derived for the estimation/optimization problem [39, 7]. By contrast, up to our knowledge, little is known for approximating the classification/integration one (only crude approximations were proposed in [56]).

**Contributions.** A major contribution of this paper is to develop an accurate approximation for the classification/integration problem. In particular, we show that our approximation error vanishes for  $x \rightarrow 0$  and  $x \rightarrow \pm\infty$  when  $\nu = 1$ , see Theorem 1 and Theorem 2. We next generalize this result for  $\frac{2}{3} < \nu < 2$  in Theorem 3 and Theorem 4. On top of that,

we prove that the two problems enjoy some important desired properties in [Proposition 2](#) and [Proposition 3](#). These theoretical results allow the two quantities to be approximated by functions that can be quickly evaluated in order to be incorporated in fast algorithms. Our last contribution is experimental and concerns the performance evaluation of the proposed model in image denoising scenario, see [Section 6](#). Together with this paper, we have released our implementation at <https://bitbucket.org/cdeledalle/ggmm-epll>.

**2. Background.** In this section we provide a brief background on the use of patch-based priors in image restoration and a quick overview of popular patch and image priors.

**2.1. Image restoration with patch based priors.** We consider the problem of estimating an image  $\mathbf{u} \in \mathbb{R}^N$  ( $N$  is the number of pixels) from noisy linear observations  $\mathbf{v} = \mathcal{A}\mathbf{u} + \mathbf{w}$ , where  $\mathcal{A} : \mathbb{R}^N \rightarrow \mathbb{R}^M$  is a linear operator and  $\mathbf{w} \in \mathbb{R}^M$  is a noise component assumed to be white and Gaussian with variance  $\sigma^2$ . In this paper, we will focus on standard denoising problems where  $\mathcal{A}$  is the identity matrix, but in more general settings, it can account for loss of information such as blurring. Typical examples for operator  $\mathcal{A}$  are: a low pass filter (for *deconvolution*), a masking operator (for *inpainting*), or a projection on a random subspace (for *compressive sensing*). To reduce noise and stabilize the inversion of  $\mathcal{A}$ , some *prior* information is used for the estimation of  $\mathbf{u}$ . Recent techniques [[19](#), [62](#), [57](#)] include this *prior* information as a model for the distribution of patches found in natural clean images. The Expected Patch Log-Likelihood (EPLL) algorithm [[62](#)] restores an image by maximum *a posteriori* estimation, corresponding to the following minimization problem:

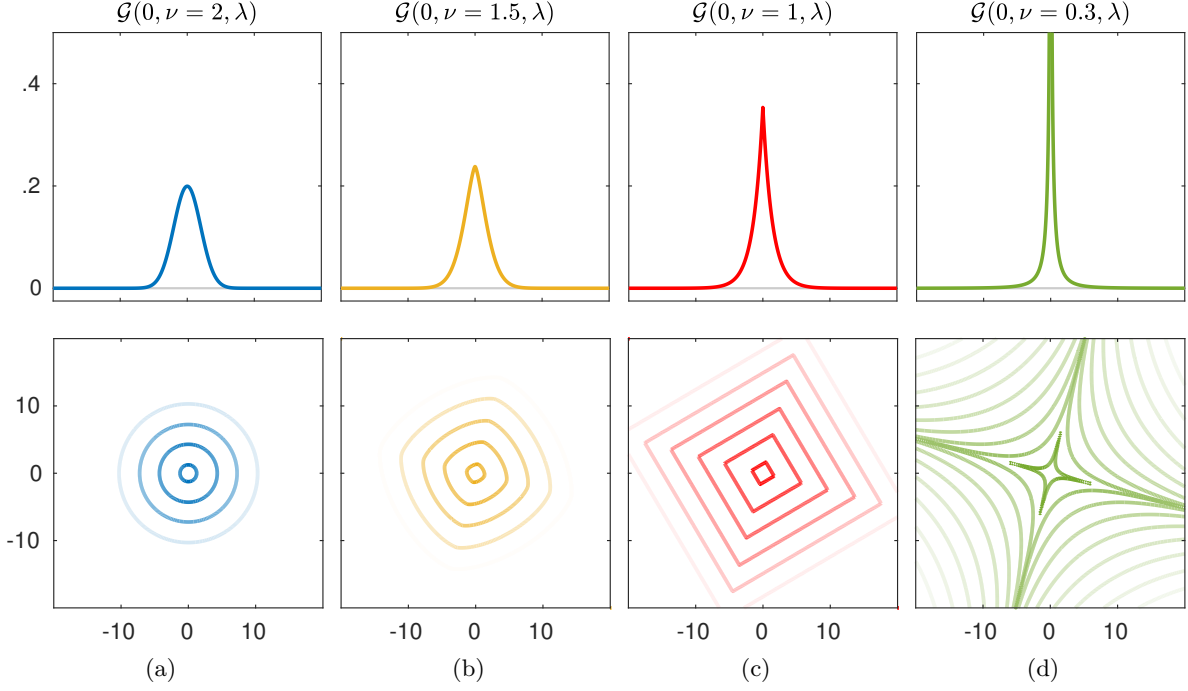
$$(2) \quad \operatorname{argmin}_{\mathbf{u} \in \mathbb{R}^n} \frac{P}{2\sigma^2} \|\mathcal{A}\mathbf{u} - \mathbf{v}\|^2 - \sum_{i=1}^N \log p(\mathcal{P}_i\mathbf{u})$$

where  $\mathcal{P}_i : \mathbb{R}^N \rightarrow \mathbb{R}^P$  is the linear operator extracting a patch with  $P$  pixels centered at the pixel with location  $i$  (typically,  $P = 8 \times 8$ ), and  $p(\cdot)$  is the *a priori* probability density function (*i.e.*, the statistical model of noiseless patches in natural images). While the first term in eq. (2) ensures that  $\mathcal{A}\mathbf{u}$  is close to the observations  $\mathbf{v}$  (this term is the negative log-likelihood under the white Gaussian noise assumption), the second term regularizes the solution  $\mathbf{u}$  by favoring an image such that all of its patches fit the *prior* model of patches in natural images.

**Optimization with half-quadratic splitting.** Problem (2) is a large optimization problem where  $\mathcal{A}$  couples all unknown pixel values of  $\mathbf{u}$  and the patch prior is often chosen non-convex. A classical technique, known as half-quadratic splitting [[22](#), [30](#)], introduces  $N$  auxiliary unknown vectors  $\mathbf{z}_i \in \mathbb{R}^P$ , and alternatively consider the penalized optimization problem, for  $\beta > 0$ , as

$$(3) \quad \operatorname{argmin}_{\substack{\mathbf{u} \in \mathbb{R}^n \\ \mathbf{z}_1, \dots, \mathbf{z}_N \in \mathbb{R}^P}} \frac{P}{2\sigma^2} \|\mathcal{A}\mathbf{u} - \mathbf{v}\|^2 + \frac{\beta}{2} \sum_{i \in \mathcal{I}} \|\mathcal{P}_i\mathbf{u} - \mathbf{z}_i\|^2 - \sum_{i \in \mathcal{I}} \log p(\mathbf{z}_i).$$

When  $\beta \rightarrow \infty$ , the problem (3) is equivalent to the original problem (2). In practice, an increasing sequence of  $\beta$  is considered, and the optimization is performed by alternating the minimization for  $\mathbf{u}$  and  $\mathbf{z}_i$ . Though little is known about the convergence of this algorithm, few iterations produce in practice remarkable results.



**Figure 1.** Illustrations of zero-mean generalized Gaussian distributions for different values of the shape parameter  $\nu$ . From left to right:  $\nu = 2, 1.5, 1, .3$ . Top: one-dimensional versions with variance  $\lambda = 2$ . Bottom: iso-contours of the two-dimensional versions with  $\mathbf{W}$  being a rotation matrix .

**Minimization with respect to  $\mathbf{u}$ .** Considering all  $\mathbf{z}_i$  to be fixed, optimizing (3) for  $\mathbf{u}$  corresponds to solving a linear inverse problem with a Tikhonov regularization, and has an explicit solution known as the linear minimum mean square estimator (or often referred to as Wiener filtering):

$$(4) \quad \hat{\mathbf{u}} = \operatorname{argmin}_{\mathbf{u} \in \mathbb{R}^n} \frac{P}{2\sigma^2} \|\mathcal{A}\mathbf{u} - \mathbf{v}\|^2 + \frac{\beta}{2} \sum_{i \in \mathcal{I}} \|\mathcal{P}_i \mathbf{u} - \hat{\mathbf{z}}_i\|^2$$

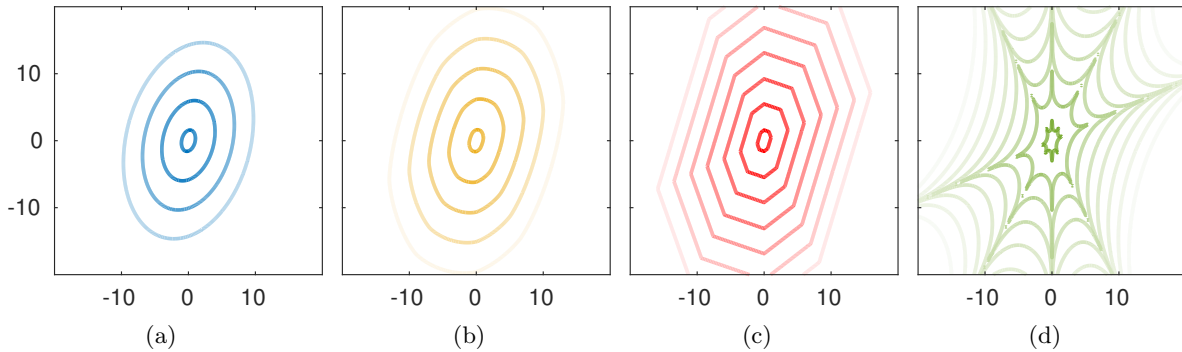
$$(5) \quad = \left( \mathcal{A}^t \mathcal{A} + \frac{\beta\sigma^2}{P} \sum_{i \in \mathcal{I}} \mathcal{P}_i^t \mathcal{P}_i \right)^{-1} \left( \mathcal{A}^t \mathbf{v} + \frac{\beta\sigma^2}{P} \sum_{i \in \mathcal{I}} \mathcal{P}_i^t \hat{\mathbf{z}}_i \right),$$

where  $\mathcal{P}_i^t \mathcal{P}_i$  is a diagonal matrix whose  $i$ -th diagonal element corresponds to the number of patches overlapping the pixel of index  $i$ .

**Minimization with respect to  $\mathbf{z}_i$ .** Considering  $\mathbf{u}$  to be fixed, optimizing (3) for  $\mathbf{z}_i$  leads to

$$(6) \quad \hat{\mathbf{z}}_i \leftarrow \operatorname{argmin}_{\mathbf{z}_i \in \mathbb{R}^P} \frac{\beta}{2} \|\tilde{\mathbf{z}}_i - \mathbf{z}_i\|^2 - \log p(\mathbf{z}_i) \quad \text{where} \quad \tilde{\mathbf{z}}_i = \mathcal{P}_i \hat{\mathbf{u}} .$$

which corresponds to the maximum *a posteriori* denoising problem under the patch prior  $p$  of a patch  $\tilde{\mathbf{z}}_i$  contaminated by Gaussian noise with variance  $1/\beta$ . The solution of this optimization problem strongly depends on the properties of the chosen patch prior.



**Figure 2.** Illustrations of the iso-lines of (a-d) four two-dimensional priors based on sparse analysis regularization with  $0 < \nu < 2$  and a dictionary  $\mathbf{\Omega} \in \mathbb{R}^{4 \times 2}$  with four directions. From left to right: the shape parameter varies respectively as  $\nu = 2, 1.5, 1, .3$ .

## 2.2. Overview of image/patch priors.

**Whitening priors.** A popular attempt in designing image priors, that can be applied to patch priors, considers a whitening transform  $\mathbf{W} \in \mathbb{R}^{P \times P}$ , *i.e.*, an orthogonal transform that decorrelates entries of  $\mathbf{z}$ . Typically, the transform  $\mathbf{W}$  can be chosen by principal component analysis on a dataset of clean patches [11], or based on some prior knowledge, *e.g.*, by assuming decorrelations in the Fourier or wavelet domain [33]. Considering decorrelated and independent transform coefficients allows for the prior distribution to be separable in the transformed domain (*i.e.*, with independent marginals). Early works on wavelet representations of images (*e.g.*, Mallat, 1989 [33] and Moulin and Liu, 1999 [39]) suggested modeling such coefficients by a zero-mean generalized Gaussian distribution (GGD) of the form

$$(7) \quad p(\mathbf{z}) \propto \exp(-\rho^\nu \|\mathbf{W}\mathbf{z}\|_\nu^\nu)$$

where  $\rho > 0$  and  $\nu \geq 0$ . Here, choosing  $1 < \nu \leq 2$  models each coefficient by a bell shape distribution with small tails (the Gaussian distribution for  $\nu = 2$ ), see first row on Figure 1(a-b). Image patches are thus assumed to be concentrated on a convex cluster (an ellipsoid for  $\nu = 2$ ), see second row on Figure 1(a-b). In this setting, the coefficients tend to get closer and closer to zero (shrinkage) as  $\rho$  increases (becomes a linear shrinkage for  $\nu = 2$ ). Consequently, this leads to diminishing edge features which yields over-smoothed images. To circumvent this issue  $\nu \leq 1$  is considered. It corresponds to modeling each coefficient by a distribution with a peak at zero (non differentiability) and large tails, see first row on Figure 1(c-d). Image patches are thus assumed to be spread on an union of orthogonal subspaces aligned with the directions of the rows of  $\mathbf{W}$ , see second row on Figure 1(c-d). In this case, the shrinkage behaves as a thresholding operator such that most coefficients are set to zeros – *sparsity* – and a few non-zero coefficients are (more or less) preserved – *robustness*. As a consequence, under this scheme, edges can be expected to be preserved and resulting images are sharper. With such a prior, the maximum *a posteriori* problem corresponding to (2) is referred to as  $\ell_\nu$  regularization on the orthogonal dictionary whose atoms are the rows of  $\mathbf{W}$ . The subsequent optimization problem is well-known to be convex for  $\nu \geq 1$ , and thus  $\nu = 1$  is a popular choice for practical applications.

*Synthesis sparse priors.* More recent works have shown that the distribution of clean image patches can be better modeled by assuming they are sparse combinations of the columns of a redundant dictionary  $\mathbf{D} \in \mathbb{R}^{P \times Q}$ ,  $Q \geq P$  (i.e., composed of linearly dependent atoms). This is typically achieved with undecimated (a.k.a, shift-invariant) wavelets [54, 8] or, more generally, frames [17, 34]. A dictionary can also be learned on a large dataset of clean patches with the k-SVD algorithm [1]. This is to suggest that the coefficients of an image are spread on a union of non-orthogonal subspaces aligned with the columns of  $\mathbf{D}$ , see Figure 2. The so-called synthesis regularization framework [20], applied to image patches in [19, 57], corresponds to the following choice of prior distribution

$$(8) \quad p(\mathbf{z}) \propto \exp(-\rho^\nu \|\boldsymbol{\alpha}\|_\nu^\nu) \quad \text{subject to} \quad \mathbf{z} = \mathbf{D}\boldsymbol{\alpha},$$

where  $\rho > 0$  and  $\nu \geq 0$ . Similar to the whitening framework, choosing  $\nu \leq 1$  enforces sparsity without penalizing large non-zero coefficients. A difficulty with such an approach is to deal with the non-orthogonality of the dictionary in the non-convex setting  $\nu < 1$ . Typically, the authors of [19, 57] consider the case  $\nu = 0$  for which greedy techniques such as orthogonal matching pursuit can be employed [10].

*Analysis sparse priors.* Alternatively, the distribution of clean images can also be captured by modeling the correlations of its patches with the rows of a redundant dictionary  $\boldsymbol{\Omega} \in \mathbb{R}^{Q \times P}$  that does not require to span the set of clean images ( $\boldsymbol{\Omega}$  can be rank deficient). This is the case with the Total-Variation model [53] that considers the gradient of an image  $\boldsymbol{\Omega} = \nabla$ , and any type of filter bank analysis. Similar to the synthesis framework, the dictionary can also be learned on a large dataset of clean patches with the analysis k-SVD algorithm [52]. This suggests that the correlations of an image are spread on a union of non-orthogonal subspaces aligned with the rows of  $\boldsymbol{\Omega}$ . The analysis regularization framework [20], corresponds to the following choice of prior

$$(9) \quad p(\mathbf{z}) \propto \exp(-\rho^\nu \|\boldsymbol{\Omega}\mathbf{z}\|_\nu^\nu)$$

where  $\rho > 0$  and  $\nu > 0$ . Again, choosing  $\nu \leq 1$  enforces correlations to be mostly zeros – *co-sparsity* – and a few large non-zero correlations. A difficulty with such an approach is to cope with the non-injectivity of  $\boldsymbol{\Omega}$ .

*GMM priors.* Rather than modeling clean patches as spread on a union of non-orthogonal sub-spaces, an alternative is to consider a union of ellipsoids (called clusters). To this end, the authors of [61, 62] suggested using a zero-mean Gaussian mixture model (GMM) prior<sup>1</sup>, that, for any patch  $\mathbf{z} \in \mathbb{R}^P$ , is given by

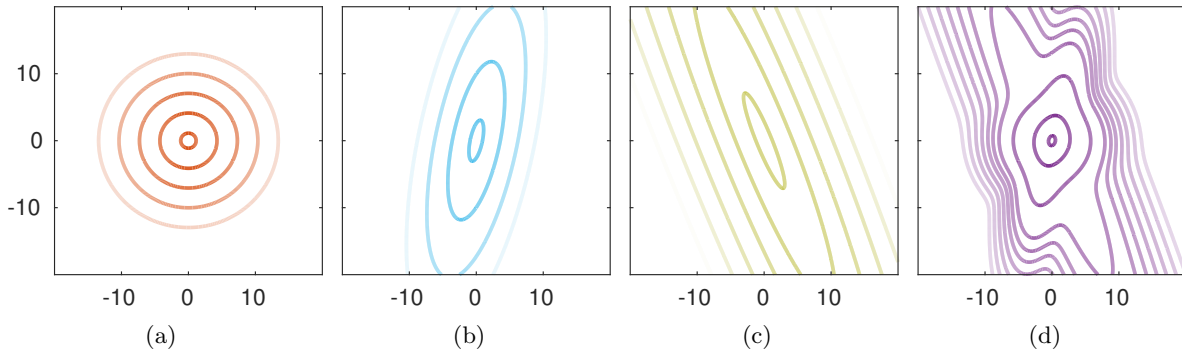
$$(10) \quad p(\mathbf{z}) = \sum_{k=1}^K w_k \mathcal{N}_P(\mathbf{z}; 0_P, \boldsymbol{\Sigma}_k)$$

where  $K$  is the number of components,  $w_k > 0$  are weights such that  $\sum_k w_k = 1$ , and  $\mathcal{N}_P(0_P, \boldsymbol{\Sigma}_k)$  denotes the multi-variate Gaussian distribution with zero-mean and covariance

---

<sup>1</sup>To enforce the zero-mean assumption, patches are first centered on zero, then denoised using eq. (6), and, finally, their initial means are added back. In fact, one can show that it corresponds to modeling  $p(\mathbf{z} - \bar{\mathbf{z}})$  with a GMM where  $\bar{z}_j = \frac{1}{P} \sum_i z_i$  for all  $1 \leq j \leq P$ .





**Figure 3.** Illustrations of the iso-lines of (a-c) three two-dimensional zero-mean Gaussian distributions and (d) their mixture with weights  $1/2$ ,  $1/3$  and  $1/6$  respectively.

$\Sigma_k \in \mathbb{R}^{P \times P}$ . An illustration of a GMM and its components is given in Figure 3. The hyper-parameters  $w_k$  and  $\Sigma_k$  can be learned using the Expectation Maximization algorithm [13] on a large dataset of clean patches [62]. As will be discussed in detail in Section 3, such a prior in the optimization problem (6) will require (i) looking for the cluster  $k^*$  that best explains the given patch  $\mathbf{z}$ , (ii) performing whitening by projecting  $\mathbf{z}$  over the main directions of that cluster (given by the eigenvectors of  $\Sigma_{k^*}$ ), and (iii) applying a linear shrinkage on the coefficients with respect to the spread of the cluster (encoded by the eigenvalues). As a function of  $\mathbf{z}$ , the resulting approach is a piece-wise linear estimator (PLE) [61]. Last but not least, when  $K = 1$ , this approach is equivalent to the whitening approach with  $\rho\mathbf{W} = \Sigma_1^{-1/2}$  and  $\nu = 2$ .

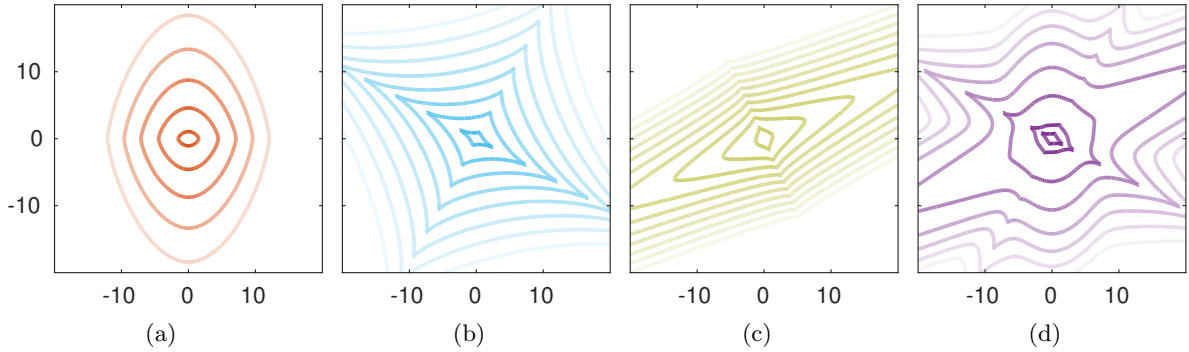
An advantage of the GMM prior over sparse priors is that all parameters of the prior model (*i.e.*,  $\Sigma_k$  and  $w_k$ ) are learned offline. As a consequence, provided that the training dataset is representative of the image to be restored, no regularization parameters are required to be tuned during runtime<sup>2</sup>. By contrast, synthesis (resp. analysis) sparse priors often consider dictionaries  $\mathbf{D}$  (resp.  $\mathbf{\Omega}$ ) with atoms of unit or prescribed norm. Consequently, the scaling of the prior controlled by  $\rho$  cannot be fixed once for all images, but will depend on the specific underlying image and the inverse problem that is being solved (*i.e.*,  $\mathcal{A}$  and  $\sigma$ ). For this reason, the parameter  $\rho$  must be tuned by the practitioner by trial and error, cross-validation or other dedicated techniques such as [24, 15, 25, 49, 12].

Note that the overview given here is mainly to provide a background for our method and is far from being exhaustive. Beyond the four aforementioned priors, many other image and patch priors have been proposed in the literature, including Gaussian scale mixtures [47], Fields-of-Experts [51], Total Generalized Variation [4] and Student mixture models [55] to just cite a few that we have omitted.

In this paper, we propose to fill the gap between sparse regularization techniques and GMM priors. To this end we suggest using a mixture of generalized Gaussian distributions

<sup>2</sup>Though all regularization parameters are learned, the solver for eq. (2) will have some free internal parameters, controlling the convergence, that require some tuning.





**Figure 4.** Illustrations of the iso-lines of (a-c) three two-dimensional zero-mean generalized Gaussian distributions with  $0 < \nu < 2$  and (d) their mixture with weights  $1/2$ ,  $1/3$  and  $1/6$  respectively.

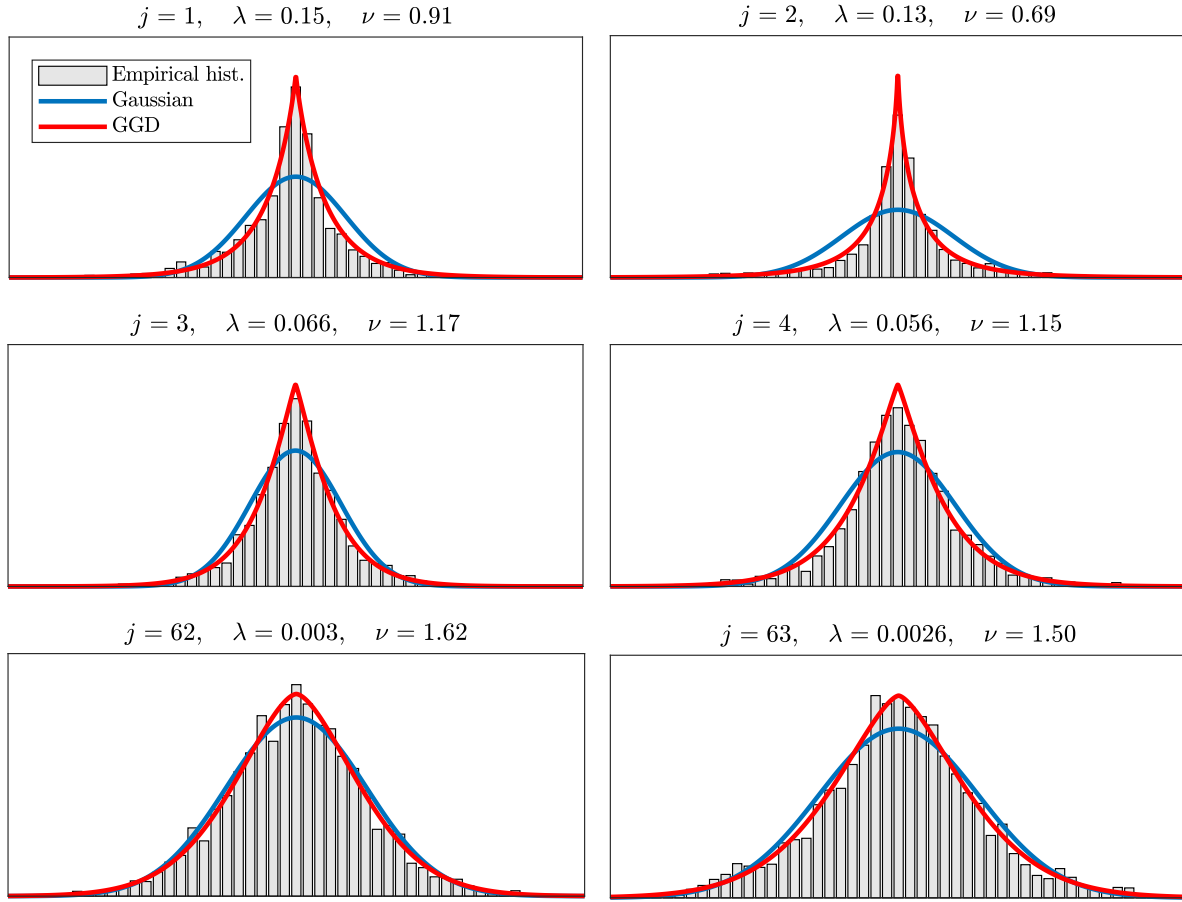
that will enable image patches to be spread over clusters that are bell shaped in some directions while being peaky in others, see [Figure 4](#). While the use of GMM priors leads to piece-wise linear estimator, our GGMM prior will lead to a piecewise non-linear shrinkage estimator.

**3. Generalized Gaussian Mixture Models.** In this paper, we aim to learn  $K$  orthogonal transforms such that each of them can map a subset (cluster) of clean patches into independent zero-mean coefficients. Instead of assuming the coefficients to be identically distributed, we consider that both the scale and the shape of their distributions may vary from one coordinate to another (within the same transform). Our motivation to assume such a highly flexible model is based on the observation illustrated in [Figure 5](#). Given one of such transform and its corresponding cluster of patches, we have displayed the histogram of the patch coefficients for six different coordinates. It can be clearly observed that the shape of the distribution varies depending on the coordinate. Some of them are peaky with heavy tails, and, therefore, would not be faithfully captured by a Gaussian distribution, as done in EPLL [62]. By contrast, some others have a bell shape, and so would not be captured properly by a peaky and heavy tailed distribution, as done by analysis sparse models [20]. This shows that one cannot simultaneously decorrelate and sparsify a cluster of clean patches for all coordinates. Since some of the coordinates reveal sparsity while some others reveal Gaussianity, we propose to use a more flexible model that can capture such variations. We propose using a multi-variate zero-mean generalized Gaussian mixture model (GGMM)

$$(11) \quad p(\mathbf{z}) = \sum_{k=1}^K w_k \mathcal{G}(\mathbf{z}; 0_P, \mathbf{\Sigma}^{(k)}, \boldsymbol{\nu}^{(k)})$$

where  $K$  is the number of components and  $w_k > 0$  are weights such that  $\sum_k w_k = 1$ . The notation  $\mathcal{G}(0_P, \mathbf{\Sigma}, \boldsymbol{\nu})$  denotes the  $P$ -dimensional generalized Gaussian distribution (GGD) with zero-mean, covariance  $\mathbf{\Sigma} \in \mathbb{R}^{P \times P}$  (symmetric positive definite) and shape parameter  $\boldsymbol{\nu} \in \mathbb{R}^P$ , whose expression is

$$(12) \quad \mathcal{G}(\mathbf{z}; 0_P, \mathbf{\Sigma}, \boldsymbol{\nu}) = \frac{\mathcal{K}}{2^{|\boldsymbol{\Sigma}_{\boldsymbol{\nu}}|^{1/2}}} \exp \left[ -\|\boldsymbol{\Sigma}_{\boldsymbol{\nu}}^{-1/2} \mathbf{z}\|_{\boldsymbol{\nu}} \right] \quad \text{with} \quad \|\mathbf{x}\|_{\boldsymbol{\nu}} = \sum_{j=1}^P |x_j|^{\nu_j},$$



**Figure 5.** Histograms of the projection of 200,000 clean patches on 6 eigenvectors  $j = 1, 2, 3, 4, 62$  and  $63$  of the covariance matrix of one component  $k$  of the mixture (with weight  $w_k = 1.3\%$ ). The contribution of each clean patch in the histograms is given by its membership values onto this component  $k$  (as obtained during the E-Step of EM). For each histogram, a generalized Gaussian distribution was adjusted by estimating the parameters  $\lambda$  and  $\nu$  by moment estimation (as obtained during the M-Step of our modified EM). For comparisons, we have also provided illustrations of the best fit obtained with a Gaussian distribution.

$$(13) \quad \text{where } \mathcal{K} = \prod_{j=1}^P \frac{\nu_j}{\Gamma(1/\nu_j)} \quad \text{and} \quad \Sigma_{\nu}^{1/2} = \Sigma^{1/2} \begin{pmatrix} \sqrt{\frac{\Gamma(1/\nu_1)}{\Gamma(3/\nu_1)}} & & & \\ & \ddots & & \\ & & \sqrt{\frac{\Gamma(1/\nu_P)}{\Gamma(3/\nu_P)}} & \\ & & & \ddots \end{pmatrix}.$$

Denoting the eigen decomposition of matrix  $\Sigma$  by  $\Sigma = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^t$  such that  $\mathbf{U} \in \mathbb{R}^{P \times P}$  is unitary and  $\mathbf{\Lambda} = \text{diag}(\lambda_1^k, \lambda_2^k, \dots, \lambda_P^k)^2$  is diagonal with positive diagonal elements  $(\lambda_j^k)^2$ ,  $\Sigma^{1/2}$  in the above expression is defined as  $\Sigma^{1/2} = \mathbf{U}\mathbf{\Lambda}^{1/2}$  and  $\Sigma^{-1/2} = \mathbf{\Lambda}^{-1/2}\mathbf{U}^t$  is its inverse.

When  $\nu$  is a constant vector with all entries equal to  $\nu_j = 2$ ,  $\mathcal{G}(0_P, \Sigma, \nu)$  is the multi-variate Gaussian distribution  $\mathcal{N}(0_P, \Sigma)$  (as used in EPLL [62]). When all  $\nu_j = 1$ , it is the multi-variate Laplacian distribution and the subsequent GGMM is a Laplacian Mixture Model (LMM). When all  $\nu_j < 1$ , it is the multi-variate hyper-Laplacian distribution and the

subsequent GGMM is a hyper-Laplacian Mixture Model (HLMM). Choosing  $K = 1$  with a constant vector  $\boldsymbol{\nu}$  corresponds to  $\ell_\nu$  regularization after whitening with  $\rho\mathbf{W} = \boldsymbol{\Sigma}_\nu^{-1/2}$ . But as motivated earlier, unlike classical multivariate GGD models [3, 45, 42], we allow for the entries of  $\boldsymbol{\nu}$  to vary from one coordinate  $j$  to another. To the best of our knowledge, the proposed work is the first one to consider this fully flexible model.

**Proposition 1.** *The multi-variate zero-mean GGD can be written as*

$$(14) \quad \mathcal{G}(\mathbf{z}; 0_P, \boldsymbol{\Sigma}, \boldsymbol{\nu}) = \prod_{j=1}^P \mathcal{G}(x_j; 0, \lambda_j, \nu_j) \quad \text{with} \quad \mathbf{x} = \mathbf{U}^t \mathbf{z}$$

$$\text{where} \quad \mathcal{G}(x; 0, \lambda, \nu) = \frac{\kappa}{2\lambda_\nu} \exp \left[ - \left( \frac{|x|}{\lambda_\nu} \right)^\nu \right]$$

$$\text{where} \quad \kappa = \frac{\nu}{\Gamma(1/\nu)} \quad \text{and} \quad \lambda_\nu = \lambda \sqrt{\frac{\Gamma(1/\nu)}{\Gamma(3/\nu)}},$$

where  $x \mapsto \mathcal{G}(x; 0, \lambda, \nu)$  is a real, even, unimodal, bounded and continuous probability density function. It is also differentiable everywhere except for  $x = 0$  when  $\nu \leq 1$ .

The proof follows directly by injecting the eigen decomposition of  $\boldsymbol{\Sigma}$  in (12) and basic properties of  $x \mapsto |x|^\nu$ . Proposition 1 shows that, for each of the  $K$  clusters, eq. (12), indeed, models a prior that is separable in a coordinate system obtained by applying the whitening transform  $\mathbf{U}^t$ . Not only is the prior separable for each coordinate  $j$ , but the shape ( $\nu_j$ ) and scale ( $\lambda_j$ ) of the distribution may vary. As observed in Figure 5, the proposed GGMM models the underlying distributions of a cluster of clean-patches much better than GGM.

**Patch denoising under GGMM prior.** We now explain why solving (6) is non-trivial when using a GGMM patch prior. In this case, for a noisy patch  $\tilde{\mathbf{z}}$  with variance  $\sigma^2$ , equation (6) becomes

$$(15) \quad \hat{\mathbf{z}} \leftarrow \operatorname{argmin}_{\mathbf{z} \in \mathbb{R}^P} \frac{1}{2\sigma^2} \|\tilde{\mathbf{z}} - \mathbf{z}\|^2 - \log \left[ \sum_{k=1}^K w_k \mathcal{G}(\mathbf{z}; 0_P, \boldsymbol{\Sigma}^{(k)}, \boldsymbol{\nu}^{(k)}) \right].$$

Due to the multi-modality of the GGMM prior, this optimization problem is highly non-convex. To circumvent this issue, we follow the strategy used by EPLL [62] in the specific case of Gaussian mixture model prior. The idea is to restrict the sum involved in the logarithm in eq. (15) to only one component  $k^*$ . If we consider the best  $k^*$  to be given (the strategy to select the best  $k^*$  will be discussed in the next section), then eq. (15) is approximated by the following simpler problem

$$(16) \quad \hat{\mathbf{z}} \leftarrow \operatorname{argmin}_{\mathbf{z} \in \mathbb{R}^P} \left\{ \frac{\|\tilde{\mathbf{z}} - \mathbf{z}\|^2}{2\sigma^2} - \log \mathcal{G}(\mathbf{z}; 0_P, \boldsymbol{\Sigma}^{(k^*)}, \boldsymbol{\nu}^{(k^*)}) = \frac{\|\tilde{\mathbf{z}} - \mathbf{z}\|^2}{2\sigma^2} + \|\boldsymbol{\Sigma}_\nu^{(k^*)}\|^{-1/2} \|\mathbf{z}\|_\nu \right\}.$$

The main advantage of this simplified version is that, by virtue of Proposition 1, the underlying optimization becomes tractable and can be separated into  $P$  one-dimensional optimization problems, as:

$$(17) \quad \hat{\mathbf{z}} = \mathbf{U} \hat{\mathbf{x}} \quad \text{where} \quad \hat{x}_j = s_{\sigma, \lambda_j}^{\nu_j}(\tilde{x}_j) \quad \text{with} \quad \tilde{\mathbf{x}} = \mathbf{U}^t \tilde{\mathbf{z}}$$

$$(18) \quad \text{and } s_{\sigma,\lambda}^{\nu}(x) \in \operatorname{argmin}_{t \in \mathbb{R}} \frac{1}{2\sigma^2}(t-x)^2 + \frac{|t|^{\nu}}{\lambda^{\nu}} \quad \text{where } \lambda_{\nu} = \lambda \sqrt{\frac{\Gamma(1/\nu)}{\Gamma(3/\nu)}}.$$

While the problem is not necessarily convex, its solution  $s_{\sigma,\lambda}^{\nu}$  is always uniquely defined almost everywhere (see, [Section 5](#)). We call this almost everywhere real function  $s_{\sigma,\lambda}^{\nu} : \mathbb{R} \rightarrow \mathbb{R}$  shrinkage function. When  $\nu = 2$ , it is a linear function that is often referred to as Wiener shrinkage. When  $\nu \neq 2$ , as we will discuss in [Section 5](#), it is a non-linear shrinkage function that can be computed in closed form for some cases or with some approximations.

Now, we address the question of finding a strategy for choosing a relevant component  $k^*$  to replace the mixture distribution inside the logarithm. The optimal component  $k^*$  can be obtained by maximizing the posterior as

$$(19) \quad k^* \in \operatorname{argmax}_{1 \leq k \leq K} p(k | \tilde{\mathbf{z}}) = \operatorname{argmax}_{1 \leq k \leq K} w_k p(\tilde{\mathbf{z}} | k) = \operatorname{argmin}_{1 \leq k \leq K} -\log w_k - \log p(\tilde{\mathbf{z}} | k)$$

where the weights of the GGMM corresponds to the prior probability  $w_k = p(k)$ . We next use the fact that the patch  $\tilde{\mathbf{z}}$  (conditioned on  $k$ ) can be expressed as  $\tilde{\mathbf{z}} = \mathbf{z} + \mathbf{n}$  where  $\mathbf{z}$  and  $\mathbf{n}$  are two independent random variables from distributions  $\mathcal{G}(0_P, \boldsymbol{\Sigma}^{(k)}, \boldsymbol{\nu}^{(k)})$  and  $\mathcal{N}(0_P, \sigma^2 \mathbf{Id}_P)$  respectively. It follows that the distribution of  $\tilde{\mathbf{z}}$  is the convolution of these latter two, and then

$$(20) \quad -\log p(\tilde{\mathbf{z}} | k) = -\log \int_{\mathbb{R}^P} \mathcal{G}(\tilde{\mathbf{z}} - \mathbf{z}; 0_P, \boldsymbol{\Sigma}^{(k)}, \boldsymbol{\nu}^{(k)}) \cdot \mathcal{N}(\mathbf{z}; 0_P, \sigma^2 \mathbf{Id}_P) d\mathbf{z}.$$

We next use [Proposition 1](#) to separate this integration problem into  $P$  one-dimensional integration problems as

$$(21) \quad -\log p(\tilde{\mathbf{z}} | k) = \sum_{j=1}^P f_{\sigma,\lambda_j}^{\nu_j^{(k)}}(\tilde{x}_j) \quad \text{with } \tilde{\mathbf{x}} = \mathbf{U}^t \tilde{\mathbf{z}},$$

$$(22) \quad \text{where } f_{\sigma,\lambda}^{\nu}(x) = -\log \int_{\mathbb{R}} \mathcal{G}(x-t; 0, \lambda, \nu) \cdot \mathcal{N}(t; 0, \sigma^2) dt.$$

We call the real function  $f_{\sigma,\lambda}^{\nu} : \mathbb{R} \rightarrow \mathbb{R}$  the discrepancy function which measures the goodness of fit of a GGD to the noisy value  $x$ . When  $\nu = 2$ , this function is quadratic with  $x$ . For  $\nu \neq 2$ , as we will discuss in [Section 4](#), it is a non-quadratic function, that can be efficiently approximated based on an in-depth analysis of its asymptotic behavior.

The next two sections are dedicated to the analysis and approximations of the discrepancy function  $f_{\sigma,\lambda}^{\nu}$  and the shrinkage function  $s_{\sigma,\lambda}^{\nu}$ , respectively.

**4. Discrepancy function: analysis and approximations.** From its definition given in eq. (22), the discrepancy function reads for  $\nu > 0$ ,  $\sigma > 0$  and  $\lambda > 0$ , as

$$(23) \quad f_{\sigma,\lambda}^{\nu}(x) = -\log \frac{1}{\sqrt{2\pi\sigma}} \frac{\nu}{2\lambda\nu\Gamma(1/\nu)} - \log \int_{-\infty}^{\infty} \exp\left(-\frac{(x-t)^2}{2\sigma^2}\right) \exp\left[-\left(\frac{|t|}{\lambda}\right)^{\nu}\right] dt.$$

It corresponds to the negative logarithm of the distribution of the sum of a zero-mean generalized Gaussian and a zero-mean Gaussian random variables. When  $\nu = 2$ , the generalized

Gaussian random variable becomes Gaussian, and the resulting distribution is also Gaussian with zero-mean and variance  $\sigma^2 + \lambda^2$ , and then

$$(24) \quad f_{\sigma,\lambda}^2(x) = \frac{1}{2} \left[ \log 2\pi + \log(\sigma^2 + \lambda^2) + \frac{x^2}{\sigma^2 + \lambda^2} \right].$$

**Remark 1.** For  $\nu = 2$ , a direct consequence of (24) is that  $-\log p(\tilde{\mathbf{z}} | k)$  is as an affine function of the Mahanalobis distance between  $\tilde{\mathbf{z}}$  and  $0_P$  for the covariance matrix  $\Sigma^{(k)} + \sigma^2 \mathbf{Id}_P$ :

$$(25) \quad -\log p(\tilde{\mathbf{z}} | k) = \frac{1}{2} \left[ P \log 2\pi + \log |\Sigma^{(k)} + \sigma^2 \mathbf{Id}_P| + \tilde{\mathbf{z}}^t (\Sigma^{(k)} + \sigma^2 \mathbf{Id}_P)^{-1} \tilde{\mathbf{z}} \right].$$

When  $\nu = 1$ , the generalized Gaussian random variable becomes Laplacian, and the distribution resulting from the convolution also has a closed form which leads to the following discrepancy function

$$(26) \quad f_{\sigma,\lambda}^1(x) = \log(2\sqrt{2}\lambda) - \frac{\sigma^2}{\lambda^2} - \log \left[ e^{\frac{\sqrt{2}x}{\lambda}} \operatorname{erfc} \left( \frac{x}{\sqrt{2}\sigma} + \frac{\sigma}{\lambda} \right) + e^{-\frac{\sqrt{2}x}{\lambda}} \operatorname{erfc} \left( -\frac{x}{\sqrt{2}\sigma} + \frac{\sigma}{\lambda} \right) \right],$$

refer to [Appendix A](#) for derivation.

To the best of our knowledge, there are no simple expressions for other values of  $\nu$ . One solution proposed by [56] is to express this in terms of the bi-variate Fox-H function [37]. This, rather cumbersome expression, is computationally demanding. In practice, this special function requires numerical integration techniques over complex lines [46], and is thus difficult to numerically evaluate it efficiently. Since, in our application, we need to evaluate this function a large number of times, we cannot utilize this solution.

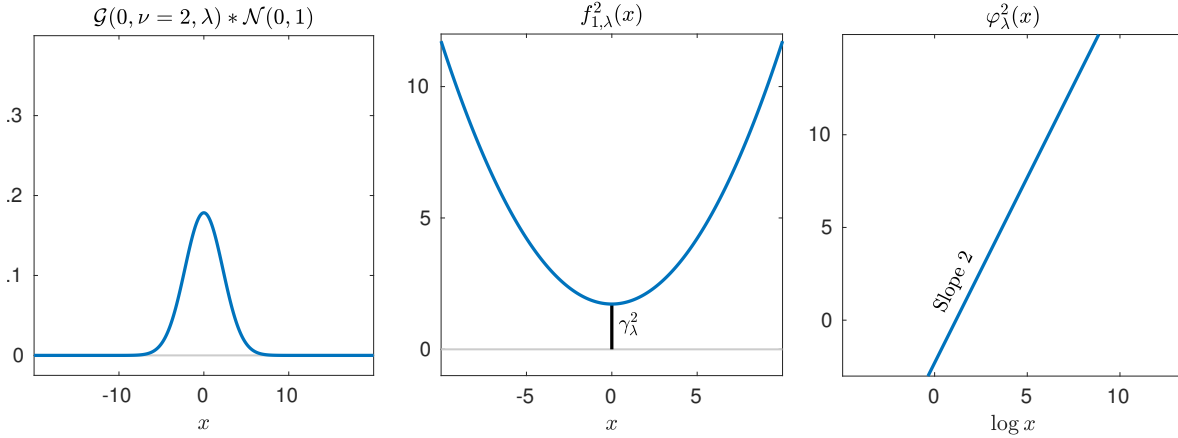
In [56], the authors have also proposed to approximate this non-trivial distribution by another GGD. For fixed values of  $\sigma$ ,  $\lambda$  and  $\nu$ , they proposed three different numerical techniques to estimate its parameters  $\lambda'$  and  $\nu'$  that best approximate either the kurtosis, the tail or the cumulative distribution function. Based on their approach, the discrepancy function  $f_{\sigma,\lambda}^\nu(x)$  would thus be a power function of the form  $x^{\nu'}$ .

In this paper, we show that  $f_{\sigma,\lambda}^\nu$  does, indeed, asymptotically behave as a power function for small and large values of  $x$ , but the exponent can be quite different for these two asymptotics. We believe that these different behaviors are important to be preserved in our application context. For this reason,  $f_{\sigma,\lambda}^\nu$  cannot be modeled as a power function through a GGD distribution. Instead, we provide an alternative solution that is able to capture the correct behavior for both of these asymptotics, and that also permits fast computation.

**4.1. Theoretical analysis.** In this section, we perform a thorough theoretical analysis of the discrepancy function, in order to approximate it accurately. Let us first introduce some basic properties regarding the discrepancy function.

**Proposition 2.** Let  $\nu > 0$ ,  $\sigma > 0$ ,  $\lambda > 0$  and  $f_{\sigma,\lambda}^\nu$  as defined in eq. (22). The following relations hold true

$$(reduction) \quad f_{\sigma,\lambda}^\nu(x) = \log \sigma + f_{1,\lambda/\sigma}^\nu(x/\sigma),$$



**Figure 6.** From left to right: the convolution of a Gaussian distribution with standard deviation  $\lambda = 2$  with a Gaussian distribution with standard deviation  $\sigma = 1$ , the corresponding discrepancy function and log-discrepancy function.

$$\begin{aligned}
 \text{(even)} & \quad f_{\sigma,\lambda}^\nu(x) = f_{\sigma,\lambda}^\nu(-x) , \\
 \text{(unimodality)} & \quad |x| \geq |y| \Leftrightarrow f_{\sigma,\lambda}^\nu(|x|) \geq f_{\sigma,\lambda}^\nu(|y|) , \\
 \text{(lower bound at 0)} & \quad \min_{x \in \mathbb{R}} f_{\sigma,\lambda}^\nu(x) = f_{\sigma,\lambda}^\nu(0) > -\infty .
 \end{aligned}$$

The proofs can be found in [Appendix B](#). Based on [Proposition 2](#), we can now express the discrepancy function  $f_{\sigma,\lambda}^\nu(x) : \mathbb{R} \rightarrow \mathbb{R}$  in terms of a constant  $\gamma_\lambda^\nu$  and another function  $\varphi_\lambda^\nu : \mathbb{R}_+^* \rightarrow \mathbb{R}$ , both of which can be parameterized by only two parameters  $\lambda > 0$  and  $\nu > 0$ , as

$$(27) \quad f_{\sigma,\lambda}^\nu(x) = \log \sigma + \gamma_{\lambda/\sigma}^\nu + \begin{cases} e^{\varphi_{\lambda/\sigma}^\nu(|x/\sigma|)} & \text{if } x \neq 0 , \\ 0 & \text{otherwise ,} \end{cases}$$

$$(28) \quad \text{where } \varphi_\lambda^\nu(x) = \log [f_{1,\lambda}^\nu(x) - \gamma_\lambda^\nu] \quad \text{and} \quad \gamma_\lambda^\nu = f_{1,\lambda}^\nu(0) .$$

We call  $\varphi_\lambda^\nu$  the log-discrepancy function.

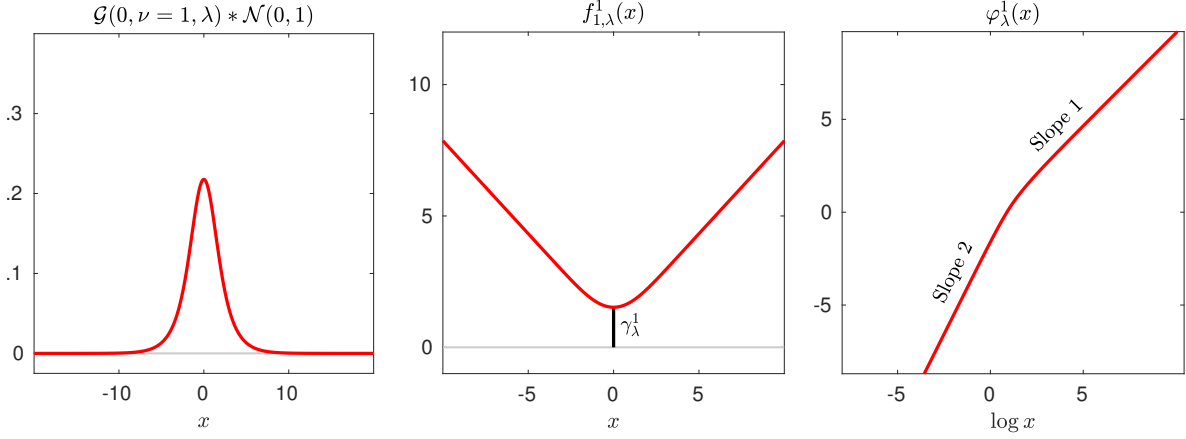
At this point, let us consider an instructive toy example for the case when  $\nu = 2$ . In this case, from eq. (24), we can deduce that the log-discrepancy function is a log-linear function (*i.e.*, a linear function of  $\log x$ )

$$(29) \quad \varphi_\lambda^2(x) = \alpha \log x + \beta ,$$

$$(30) \quad \text{and } \gamma_\lambda^2 = \frac{1}{2} [\log 2\pi + \log(1 + \lambda^2)] ,$$

$$(31) \quad \text{where } \alpha = 2 \quad \text{and} \quad \beta = -\log 2 - \log(1 + \lambda^2) .$$

Here, the slope  $\alpha = 2$  reveals the quadratic behavior of the discrepancy function. [Figure 6](#) gives an illustration of the resulting convolution (a Gaussian distribution), the discrepancy function (a quadratic function) and the log-discrepancy (a linear function with slope 2). Note that quadratic metrics are well-known to be non-robust to outliers, which is in complete agreement with the fact that Gaussian priors have thin tails.



**Figure 7.** From left to right: the convolution of a Laplacian distribution with standard deviation  $\lambda = 2$  with a Gaussian distribution with standard deviation  $\sigma = 1$ , the corresponding discrepancy function and log-discrepancy function.

Another example is the case of  $\nu = 1$ . From eq. (26), the log-discrepancy is given by

$$(32) \quad \varphi_\lambda^1(x) = \log \left[ \log \left[ 2 \operatorname{erfc} \left( \frac{1}{\lambda} \right) \right] - \log \left[ e^{\frac{\sqrt{2}x}{\lambda}} \operatorname{erfc} \left( \frac{x}{\sqrt{2}} + \frac{1}{\lambda} \right) + e^{-\frac{\sqrt{2}x}{\lambda}} \operatorname{erfc} \left( -\frac{x}{\sqrt{2}} + \frac{1}{\lambda} \right) \right] \right],$$

$$(33) \quad \text{and} \quad \gamma_\lambda^1 = \frac{1}{2} \log 2 + \log \lambda - \frac{1}{\lambda^2} - \log \left[ \operatorname{erfc} \left( \frac{1}{\lambda} \right) \right].$$

Unlike for  $\nu = 2$ , this function is not log-linear and thus  $f_{\sigma, \lambda}^1$  is not a power function. Nevertheless, as shown by the next two theorems, it is also asymptotically log-linear for small and large values of  $x$ .

**Theorem 1.** *The function  $\varphi_\lambda^1$  is asymptotically log-linear in the vicinity of 0*

$$(34) \quad \varphi_\lambda^1(x) \underset{0}{\sim} \alpha_1 \log x + \beta_1,$$

$$(35) \quad \text{where} \quad \alpha_1 = 2 \quad \text{and} \quad \beta_1 = -\log \lambda + \log \left[ \frac{1}{\sqrt{\pi}} \frac{\exp \left( -\frac{1}{\lambda^2} \right)}{\operatorname{erfc} \left( \frac{1}{\lambda} \right)} - \frac{1}{\lambda} \right].$$

The proof can be found in [Appendix C](#).

**Theorem 2.** *The function  $\varphi_\lambda^1$  is asymptotically log-linear in the vicinity of  $+\infty$*

$$(36) \quad \varphi_\lambda^1(x) \underset{\infty}{\sim} \alpha_2 \log x + \beta_2,$$

$$(37) \quad \text{where} \quad \alpha_2 = 1 \quad \text{and} \quad \beta_2 = \frac{1}{2} \log 2 - \log \lambda.$$

The proof can be found in [Appendix D](#).

**Theorem 1** and **Theorem 2** show that  $\varphi_\lambda^1$  has two different asymptotics that can be approximated by a log-linear function. Interestingly, the exponent  $\alpha_1 = 2$  in the vicinity of 0 shows that the Gaussian distribution involved in the convolution prevails over the Laplacian distribution and thus, the behavior of  $f_{\sigma, \lambda}^1$  is quadratic. Similarly, the exponent  $\alpha_2 = 1$  in the



vicinity of  $+\infty$  shows that the Laplacian distribution involved in the convolution prevails over the Gaussian distribution and the behavior of  $f_{\sigma,\lambda}^1$  is then linear. These results are supported by [Figure 7](#) which illustrates the resulting convolution, the discrepancy function (eq. (26)) and the log-discrepancy function (eq. (32)). Furthermore, the discrepancy function  $f_{\sigma,\lambda}^1$  shares a similar behavior with the well-known Huber loss function [28], known to be more robust to outliers. This is again in complete agreement with the fact that Laplacian priors have heavier tails.

In the case  $\frac{2}{3} < \nu < 2$ , even though  $\varphi_\lambda^\nu$  has no simple closed form expression, the similar conclusions can be made as a result of the next two theorems.

**Theorem 3.** *Let  $\nu > 0$ . The function  $\varphi_\lambda^\nu$  is asymptotically log-linear in the vicinity of 0*

$$\varphi_\lambda^\nu(x) \underset{0}{\sim} \alpha_1 \log x + \beta_1 ,$$

$$\text{where } \alpha_1 = 2 \quad \text{and} \quad \beta_1 = -\log 2 + \log \left( 1 - \frac{\int_{-\infty}^{\infty} t^2 e^{-\frac{t^2}{2}} \exp \left[ -\left( \frac{|t|}{\lambda^\nu} \right)^\nu \right] dt}{\int_{-\infty}^{\infty} e^{-\frac{t^2}{2}} \exp \left[ -\left( \frac{|t|}{\lambda^\nu} \right)^\nu \right] dt} \right) .$$

The proofs can be found in [Appendix E](#).

**Theorem 4.** *Let  $\frac{2}{3} < \nu < 2$ , then  $\varphi_\lambda^\nu$  is asymptotically log-linear in the vicinity of  $+\infty$*

$$\varphi_\lambda^\nu(x) \underset{\infty}{\sim} \alpha_2 \log x + \beta_2 ,$$

$$\text{where } \alpha_2 = \nu \quad \text{and} \quad \beta_2 = -\nu \log \lambda - \frac{\nu}{2} \log \frac{\Gamma(1/\nu)}{\Gamma(3/\nu)} .$$

The proofs rely on a result of Berman (1992) [2] and is detailed in [Appendix F](#).

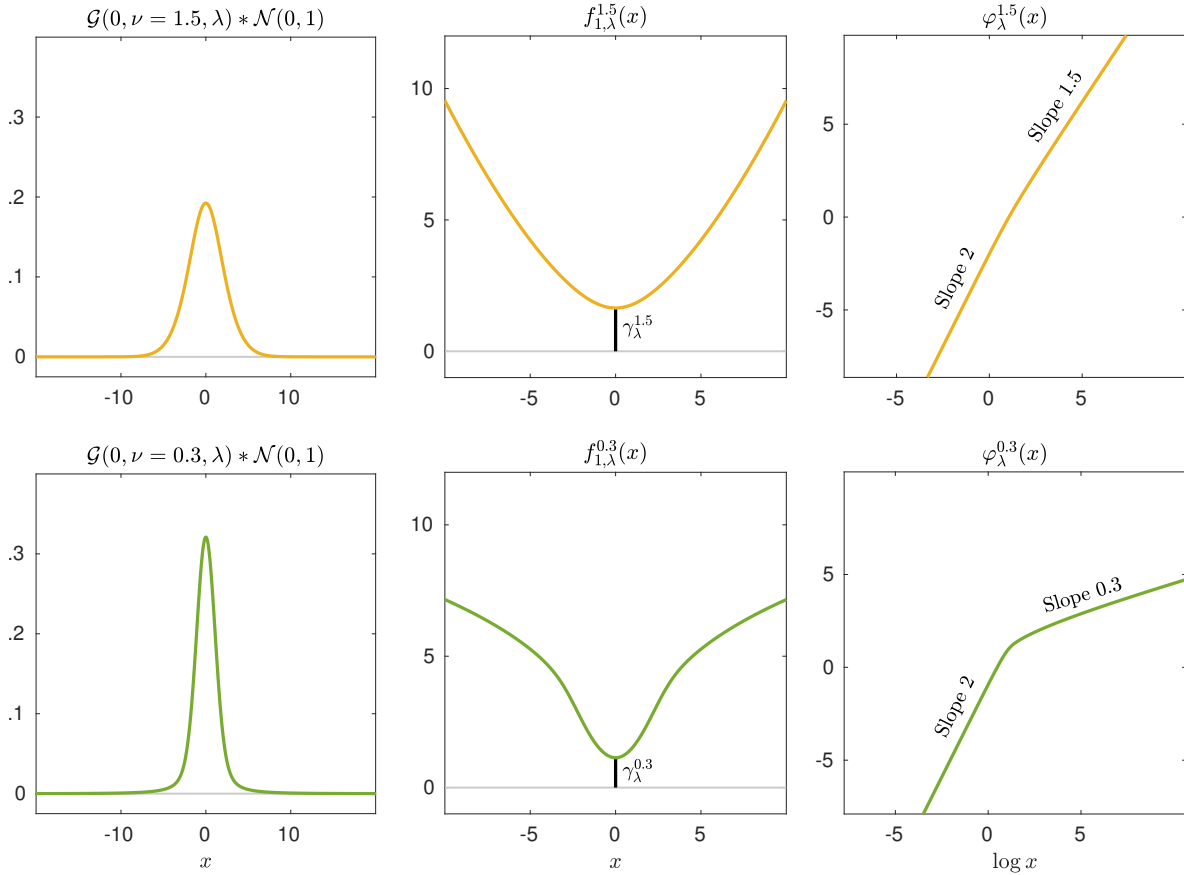
**Remark 2.** *For  $\nu > 2$ , an asymptotic log-linear behavior with  $\alpha_2 = 2$  and  $\beta_2 = -\log 2$  can be obtained using exactly the same sketch of proof as the one of [Theorem 4](#).*

**Remark 3.** *For  $\nu = 2$ , we have  $\varphi_\lambda^2$  is linear,  $\beta_1 = -\log 2 - \log(1 + \lambda^2)$  and  $\beta_2 = -\log 2 - \log \lambda^2$ , which shows that [Theorem 4](#) cannot hold true for  $\nu = 2$ .*

**Remark 4.** *For  $\nu = 1$ , [Theorem 1](#) and [Theorem 2](#) coincide with [Theorem 3](#) and [Theorem 4](#).*

**Remark 5.** *For  $0 < \nu \leq \frac{2}{3}$ , though we did not succeed in proving it, our numerical simulations also revealed a log-linear asymptotic behavior for  $x \rightarrow \infty$  in perfect agreement with the expression of  $\alpha_2$  and  $\beta_2$  given in [Theorem 4](#).*

Again, the exponent  $\alpha_1 = 2$  in the vicinity of 0 shows that the Gaussian distribution involved in the convolution prevails over the generalized Gaussian distribution and the behavior of  $f_{\sigma,\lambda}^\nu$  is then quadratic. Similarly, the exponent  $\alpha_2 = \nu$  in the vicinity of  $+\infty$  shows that the generalized Gaussian distribution involved in the convolution prevails over the Gaussian distribution and the behavior of  $f_{\sigma,\lambda}^\nu$  is then a power function of the form  $x^\nu$ . These results are supported by [Figure 8](#) that illustrates the resulting convolution, the discrepancy function and the log-discrepancy function for  $\nu = 1.5$  and  $\nu = .3$ . Moreover, the discrepancy function  $f_{\sigma,\lambda}^\nu$  with  $\nu \leq 1$  shares a similar behavior with well-known robust M-estimator loss functions

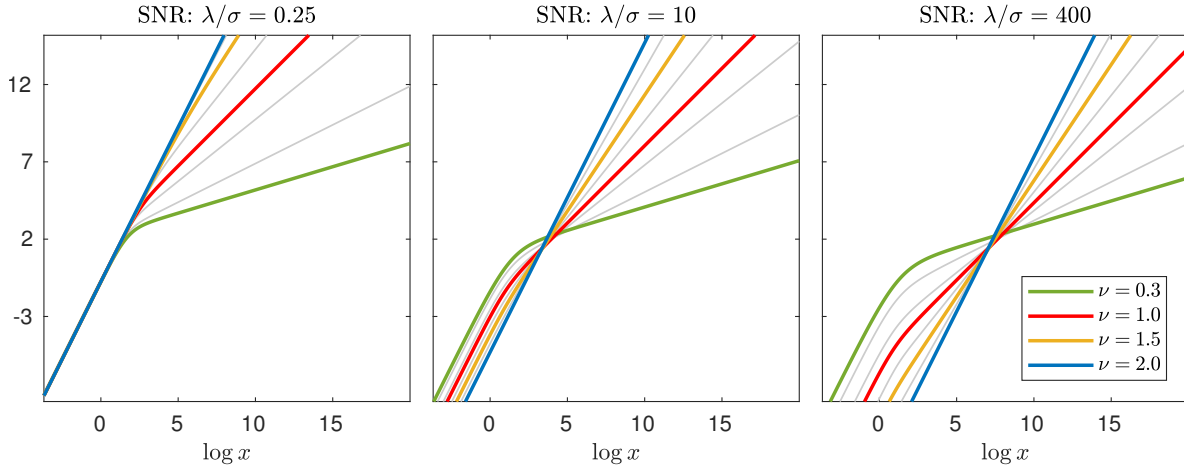


**Figure 8.** From left to right: the convolution of a generalized Gaussian distribution with standard deviation  $\lambda = 2$  with a Gaussian distribution with standard deviation  $\sigma = 1$ , the corresponding discrepancy function and log-discrepancy function. From top to bottom: the GGD has a shape parameter  $\nu = 1.5$  and  $.3$ , respectively.

[27]. In particular, the asymptotic case for  $\nu \rightarrow 0$  resembles the Tukey's bisquare loss, known to be insensitive to outliers. This is again in complete agreement with GGD priors having larger tails as  $\nu$  goes to 0.

Figure 9 shows the evolution of the log-discrepancy function for various values of  $\nu$  in the context of three different signal-to-noise ratio  $\lambda/\sigma$  (SNR). One can observe that as the SNR decreases (resp., increases), the left (resp., right) asymptotic behavior starts dominating over the right (resp., left) asymptotes. In other words, for  $\nu < 2$ , the intersection of the two asymptotes goes to  $+\infty$  (resp.,  $-\infty$ ). Last but not least, for  $0 < \nu \leq 2$ , the log-discrepancy function  $\varphi_\lambda^\nu$  is always concave and since  $\alpha_2 \leq \alpha_1$  it is thus upper-bounded by its left and right asymptotes.

From Theorem 3, Theorem 4 and Remark 5, we can now build two asymptotic log-linear approximations for  $\varphi_\lambda^\nu$ , with  $0 < \nu < 2$ , and subsequently an asymptotic power approximation for  $f_{\sigma, \lambda}^\nu$  by using the relation (27). Next, we explain the approximation process of the in-between behavior, as well as its efficient evaluation.



**Figure 9.** Illustrations of the log-discrepancy function for various  $0 < \nu < 2$  and SNR  $\lambda/\sigma$ .

**4.2. Numerical approximation.** We now describe the proposed approximation of the discrepancy function  $f_{1,\lambda}^\nu$  through approximating  $\hat{\varphi}_\lambda^\nu$  of the log-discrepancy function as

$$(38) \quad \hat{f}_{1,\lambda}^\nu(x) = \gamma_\lambda^\nu + \exp \hat{\varphi}_\lambda^\nu(x) \quad \text{where} \quad \gamma_\lambda^\nu = f_{1,\lambda}^\nu(0) .$$

Based on our previous theoretical analysis, a solution preserving the asymptotic, increasing and concave behaviors of  $\varphi_{1,\lambda}^\nu$  can be defined by making use of the following approximations

$$(39) \quad \hat{\varphi}_\lambda^\nu(x) = \alpha_1 \log |x| + \beta_1 - \mathbf{rec}(\alpha_1 \log |x| + \beta_1 - \alpha_2 \log |x| - \beta_2) ,$$

where  $\mathbf{rec}$  is a so-called rectifier function that is positive, increasing, convex and satisfies

$$(40) \quad \lim_{x \rightarrow -\infty} \mathbf{rec}(x) = 0 \quad \text{and} \quad \mathbf{rec}(x) \underset{x \rightarrow \infty}{\sim} x .$$

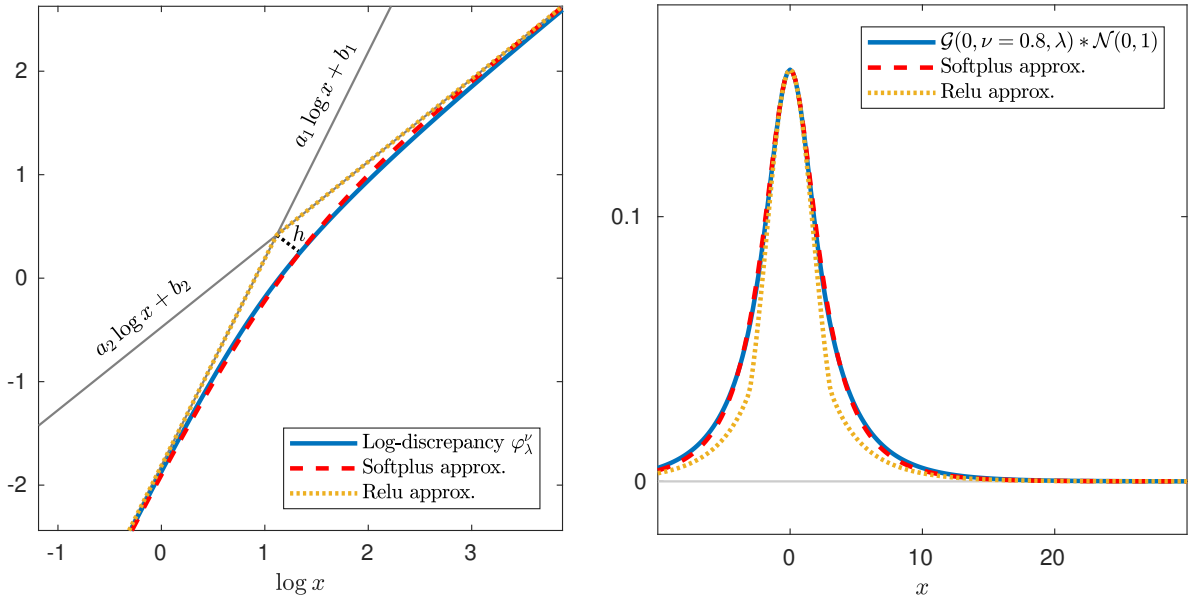
In this paper, we consider the two following rectifying functions

$$(41) \quad \mathbf{relu}(x) = \max(0, x) \quad \text{and} \quad \mathbf{softplus}(x) = h \log \left[ 1 + \exp \left( \frac{x}{h} \right) \right], \quad h > 0 ,$$

as coined respectively in [40] and [18]. Using the function  $\mathbf{relu}$  (Rectified linear unit) leads to an approximation  $\hat{\varphi}_\lambda^\nu$  that is exactly equal to the asymptotes of  $\varphi_\lambda^\nu$  with a singularity at their crossing point. In this paper, we will instead use the function  $\mathbf{softplus}$  as it allows the approximation of  $\varphi_\lambda^\nu$  to converge smoothly to the asymptotes without singularity. Its behavior is controlled by the parameter  $h > 0$ . The smaller the value of  $h$  is, the faster the convergence speed to the asymptotes.

The parameter  $h$  should be chosen such that the approximation error between  $\hat{\varphi}_\lambda^\nu(x)$  and  $\varphi_\lambda^\nu(x)$  is as small as possible. This can be done numerically by first evaluating  $\varphi_\lambda^\nu(x)$  with integration techniques for a large range of values  $x$ , and then selecting the parameter  $h$  by least square. Of course, the optimal value for  $h$  depends on the parameter  $\lambda$  and  $\nu$ .

Figure 10 gives an illustration of our approximations of the log-discrepancy and the corresponding distribution obtained with  $\mathbf{relu}$  and  $\mathbf{softplus}$ . On this figure the underlying



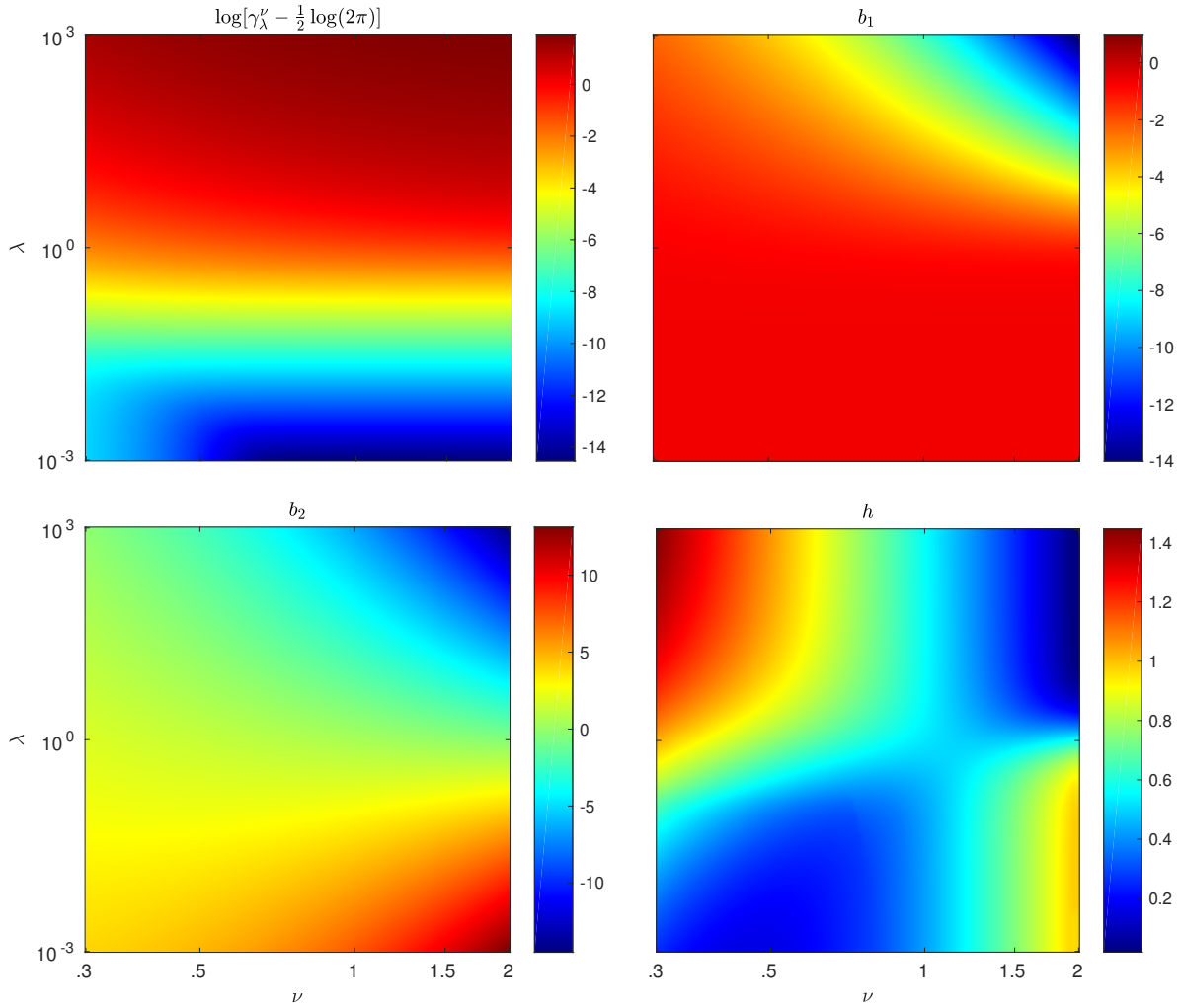
**Figure 10.** Illustrations of our approximations of  $\varphi_\lambda^\nu$  and the corresponding underlying posterior distribution  $\mathcal{N}(0, \nu, \lambda) * \mathcal{G}(0, 1)$  (where  $\nu = .8$  and  $\lambda = 4$ ). The blue curves have been obtained by evaluating the convolution using numerical integration techniques for all  $x$ . The dashed curves are obtained using the proposed `relu`- and `softplus`-based approximations that have closed-form expressions.

functions have been obtained by numerical integration for a large range of value of  $x$ . One can observe that using `softplus` provides a better approximation than `relu`.

Our approximation for  $\hat{f}_{1,\lambda}^\nu(x)$  is parameterized by six scalar values:  $\gamma_\nu^\lambda$ ,  $\alpha_1$ ,  $\beta_1$ ,  $\alpha_2$ ,  $\beta_2$  and  $h$  that depend only on the original parameters  $\lambda$  and  $\nu$ . From our previous analysis, we have that  $\alpha_1 = 2$  and  $\alpha_2 = \nu$ . The other parameters are non-linear functions of  $\lambda$  and  $\nu$ . The parameters  $\gamma_\nu^\lambda$ ,  $\beta_1$  and  $\beta_2$  require either performing numerical integration or evaluating the special function  $\Gamma$ . As discussed, the parameter  $h$  requires numerical integration for various  $x$  and then optimization. For these reasons, these values cannot be computed during runtime. Instead, we pre-compute these four parameters offline for 10,000 different combinations of  $\lambda$  and  $\nu$  values in the intervals  $[10^{-3}, 10^3]$  and  $[0.3, 2]$ , respectively (the choice for this range will be motivated in [Section 6](#)). The resulting values are then stored in four corresponding look-up-tables. During runtime, these parameters are retrieved online by bi-linear extrapolation and interpolation. The four look-up-tables are given in [Figure 11](#). We will see in [Section 6](#) that using the approximation  $\hat{f}_{1,\lambda}^\nu$  results in substantial acceleration without significant loss of performance as compared to computing  $f_{1,\lambda}^\nu$  directly by numerical integration during runtime.

**5. Shrinkage functions: analysis and approximations.** Recall that from its definition given in [eq. \(18\)](#), the shrinkage function is defined for  $\nu > 0$ ,  $\sigma > 0$  and  $\lambda > 0$ , as

$$(42) \quad s_{\sigma,\lambda}^\nu(x) \in \operatorname{argmin}_{t \in \mathbb{R}} \frac{(x-t)^2}{2\sigma^2} + \lambda_\nu^{-\nu} |t|^\nu .$$



**Figure 11.** Look-up-tables used to store the values of the parameters  $\gamma_\lambda^\nu$ ,  $b_1$ ,  $b_2$  and  $h$  for various  $.3 \leq \nu \leq 2$  and  $10^{-3} \leq \lambda \leq 10^3$ . A regular grid of 100 values has been used for  $\nu$  and a logarithmic grid of 100 values has been used for  $\lambda$ . This leads to a total of 10,000 combinations for each of the four look-up-tables.

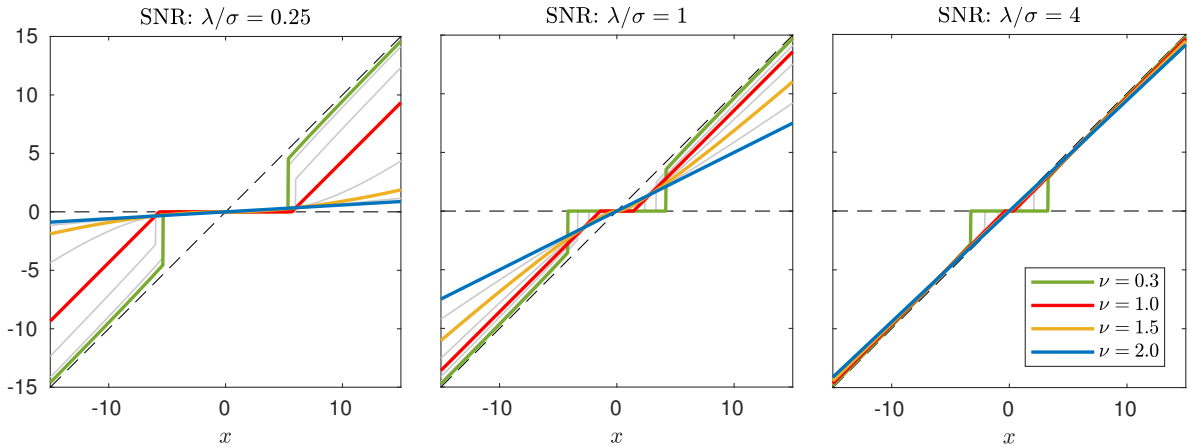
**5.1. Theoretical analysis.** Except for some particular values of  $\nu$  (see, Subsection 5.2), Problem (42) does not have explicit solutions. Nevertheless, as shown in [39], Problem (42) admits two (not necessarily distinct) solutions. One of them is implicitly characterized as

$$(43) \quad s_{\sigma,\lambda}^\nu(x) = \begin{cases} 0 & \text{if } 0 < \nu \leq 1 \text{ and } |x| \leq \tau_\lambda^\nu, \\ t^* & \text{otherwise,} \end{cases}$$

where  $t^* = x - \text{sign}(t^*)\nu\sigma^2\lambda_\nu^{-\nu}|t^*|^{\nu-1}$ ,

$$\text{and } \tau_\lambda^\nu = \begin{cases} (2-\nu)(2-2\nu)^{-\frac{1-\nu}{2-\nu}}(\sigma^2\lambda_\nu^{-\nu})^{\frac{1}{2-\nu}} & \text{if } \nu < 1, \\ \sigma^2\lambda^{-1} & \text{otherwise } (\nu = 1). \end{cases}$$

The other one is obtained by changing  $|x| \leq \tau_\lambda^\nu$  to  $|x| < \tau_\lambda^\nu$  in (43), and so they coincide for almost every  $(x, \lambda, \sigma, \nu)$ . As discussed in [39], for  $\nu > 1$ ,  $s_{\sigma,\lambda}^\nu(x)$  is differentiable, and for



**Figure 12.** Illustrations of the shrinkage function for various  $0 < \nu < 2$  and SNR  $\lambda/\sigma$ .

$\nu \leq 1$ , the shrinkage exhibits a threshold  $\tau_\lambda^\nu$  that produces sparse solutions. [Proposition 3](#) summarizes a few important properties.

**Proposition 3.** Let  $\nu > 0$ ,  $\sigma > 0$ ,  $\lambda > 0$  and  $s_{\sigma, \lambda}^\nu$  as defined in eq. (18). The following relations hold true

(reduction) 
$$s_{\sigma, \lambda}^\nu(x) = \sigma s_{1, \frac{\lambda}{\sigma}}^\nu\left(\frac{x}{\sigma}\right),$$

(odd) 
$$s_{\sigma, \lambda}^\nu(x) = -s_{\sigma, \lambda}^\nu(-x),$$

(shrinkage) 
$$s_{\sigma, \lambda}^\nu(x) \in \begin{cases} [0, x] & \text{if } x \geq 0 \\ [x, 0] & \text{otherwise} \end{cases},$$

(increasing with  $x$ ) 
$$x_1 \geq x_2 \Leftrightarrow s_{\sigma, \lambda}^\nu(x_1) \geq s_{\sigma, \lambda}^\nu(x_2),$$

(increasing with  $\lambda$ ) 
$$\lambda_1 \geq \lambda_2 \Leftrightarrow s_{\sigma, \lambda_1}^\nu(x) \geq s_{\sigma, \lambda_2}^\nu(x),$$

(kill low SNR) 
$$\lim_{\frac{\lambda}{\sigma} \rightarrow 0} s_{\sigma, \lambda}^\nu(x) = 0,$$

(keep high SNR) 
$$\lim_{\frac{\lambda}{\sigma} \rightarrow +\infty} s_{\sigma, \lambda}^\nu(x) = x.$$

The proofs can be found in [Appendix G](#). These properties show that  $s_{\sigma, \lambda}^\nu$  is indeed a shrinkage function (non-expansive). It shrinks the input coefficient  $x$  according to the model  $\nu$  and the modeled signal to noise ratio  $\frac{\lambda}{\sigma}$  (SNR). When  $x$  is small in comparison to the SNR, it is likely that its noise component dominates the underlying signal, and is, therefore, shrunk towards 0. Similarly, when  $x$  is large, it will likely be preserved. This is even more likely when  $\nu$  is small, since in this case large coefficients are favored by the *prior*. Illustrations of shrinkage functions for various SNR and  $\nu$  are given in [Figure 12](#).

**5.2. Numerical approximations.** The shrinkage function  $s_{\sigma, \lambda}^\nu$ , implicitly defined in (43) does not have a closed form expression in general. Nevertheless, for fixed values of  $x$ ,  $\sigma$ ,  $\lambda$ ,  $\nu$ ,

**Table 1**  
*Shrinkage function under generalized Gaussian priors*

$\nu$	Shrinkage $s_{\sigma,\lambda}^{\nu}(x)$	Remark
$< 1$	$\begin{cases} x - \gamma x^{\nu-1} + O(x^{2(\nu-1)}) & \text{if }  x  \geq \tau_{\lambda}^{\nu} \\ 0 & \text{otherwise} \end{cases}$	$\approx$ Hard-thresholding [39]
1	$\text{sign}(x) \max\left( x  - \frac{\sqrt{2}\sigma^2}{\lambda}, 0\right)$	Soft-thresholding [16]
4/3	$x + \gamma \left( \sqrt[3]{\frac{\zeta - x}{2}} - \sqrt[3]{\frac{\zeta + x}{2}} \right)$	[7]
3/2	$\text{sign}(x) \frac{\left(\sqrt{\gamma^2 + 4 x } - \gamma\right)^2}{4}$	[7]
2	$\frac{\lambda^2}{\lambda^2 + \sigma^2} \cdot x$	Wiener (LMMSE)

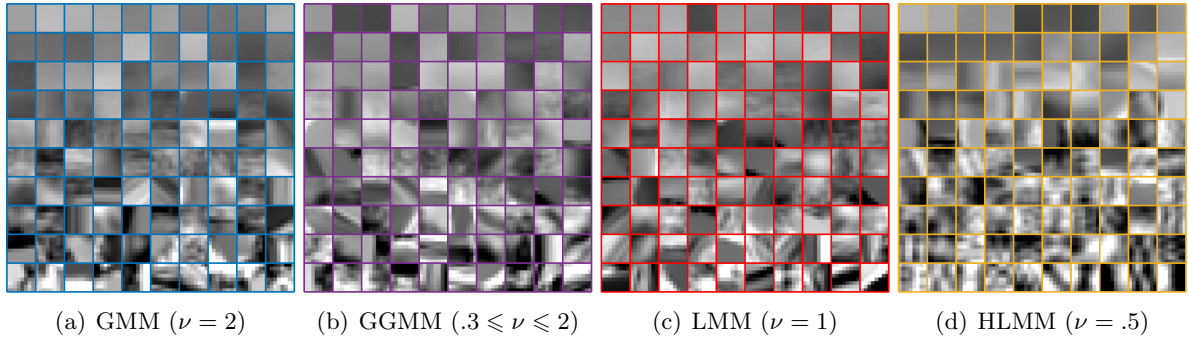
$$\text{with } \gamma = \nu\sigma^2\lambda_{\nu}^{-\nu} \quad \text{and} \quad \zeta = \sqrt{x^2 + 4\left(\frac{\gamma}{3}\right)^3}.$$

$s_{\sigma,\lambda}^{\nu}(x)$  can be estimated using iterative solvers such as Newton descent or Halley's root-finding method. These approaches converge quite fast and, in practice, reaches a satisfying solution within ten iterations. However, since in our application of interest we need to evaluate this function a large number of times, we will follow a different path in order to reduce computation time (even though we have implemented this strategy).

As discussed earlier,  $s_{\sigma,\lambda}^{\nu}$  is known in closed form for some values of  $\nu$ , more precisely:  $\nu = \{1, 4/3, 3/2, 2\}$  (as well as  $\nu = 3$  but this is out of the scope of this study), see for instance [7]. When  $\nu = 2$ , we retrieve the linear minimum mean square estimator (known in signal processing as Wiener filtering) and related to Tikhonov regularization and ridge regression. This shrinkage is linear and the slope of the shrinkage goes from 0 to 1 as the SNR increases (see Figure 12). When  $\nu = 1$ , the shrinkage is the well-known soft-thresholding [16], corresponding to the maximum a posteriori estimator under a Laplacian prior. When  $\nu < 1$ , the authors of [39] have shown that (i) the shrinkage admits a threshold with a closed-form expression (given in eq. (43)), and (ii) the shrinkage is approximately equal to hard-thresholding with an error term that vanishes when  $|x| \rightarrow \infty$ . All these expressions are summarized in Table 1.

In order to keep our algorithm as fast as possible, we propose to use the approximation of the shrinkage given for  $\nu < 1$  in [39]. Otherwise, we pick one of the four shrinkage functions corresponding to  $\nu = \{1, 4/3, 3/2, 2\}$  by nearest neighbor on the actual value of  $\nu \in [1, 2]$ . Though this approximation may seem coarse compared to the one based on iterative solvers,





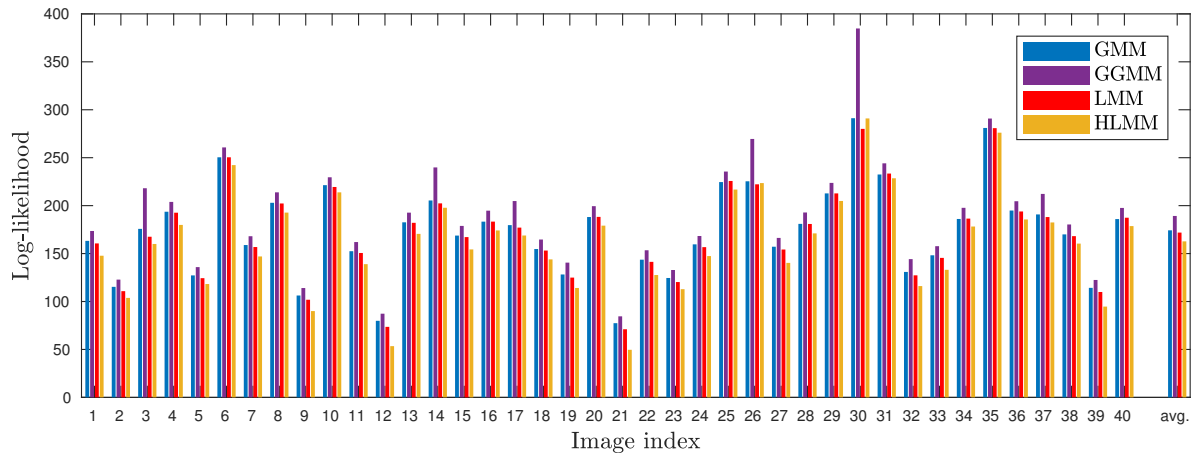
**Figure 13.** Set of 100 patches, sorted by the norm of their gradient, and generated to be independently distributed according to (from left to right) a GMM, GGMM, LMM and HLMM. For ease of visualization, only the top eigendirections corresponding to 80% of the variance have been chosen. Near-constant patches with variance smaller than  $\frac{2}{F} \|\Sigma_k\|_F^2$  have also been discarded.

we did not observe any significant loss of quality in our numerical experiments (see Section 6). Nonetheless, this alternative leads to 10 times speed-up while evaluating shrinkage.

**6. Experimental evaluation.** In this section we explain the methodology used to learn and validate the GGMM model, and present numerical experiments to compare the performance of the proposed GGMM model over existing GMM-based image denoising algorithms. To demonstrate the superiority of allowing for a flexible GGMM model, we also present results using GGMM models with fixed shape parameters,  $\nu = 1$  (Laplacian mixture model) and  $\nu = 0.5$  (Hyper-Laplacian mixture model).

**Learning.** For consistency purposes, we keep the training data and the number of mixture components in the models the same as that used in the original EPLL algorithm [62]. Specifically, we train our models on 2 million clean patches extracted from Berkeley Segmentation Dataset (BSDS) [35]. We learn  $K = 200$  zero-mean generalized Gaussian mixture components from patches of size  $8 \times 8$ . Parameter estimation is carried out using the Expectation-Maximization (EM) algorithm [13], that is known to monotonically increase the likelihood and converge to a local minimum. For applying EM to GGMM learning, we leverage standard strategies used for parameter estimation for GGD and/or GGMM that are reported in previous works [33, 3, 38, 45, 31]. We opted for a warm-start training by initializing our GGMM model with the GMM model from [62] and with initial values of shape parameters as 2. The EM algorithm is run until convergence or for a preset number of iterations (*e.g.*, 100 iterations). We noticed that a shape parameter  $\nu_j^{(k)} < .3$ , leads to numerical issues and  $\nu_j^{(k)} > 2$  leads to a local minima with several degenerate components. For this reason, we impose the constraint that the learned shape parameters satisfy  $\nu_j^{(k)} \in [.3, 2]$ . This observation is consistent with earlier works that have attempted to learn GGMM shape parameter from data [50]. Lastly, for learning Laplacian mixture model (LMM) and hyper-Laplacian mixture model (HLMM), we use the same procedure described above but force all shape parameters to be equal to 1 or 0.5, respectively.

**Model validation.** As discussed in Section 3, Figure 5 illustrates the validity of our model choices with histograms of different dimensions of a single patch cluster. It clearly shows the



**Figure 14.** Average log-likelihood of all non-overlapping patches (with subtracted mean) of each of the 40 images of our validation subset of the testing BSDS dataset for the GMM, GGMM, LMM and HLMM. The total average over the 40 images is shown in the last column.

importance of allowing the shape and scale parameter to vary across dimensions for capturing underlying patch distributions.

Since GGMM (and obviously, GMM) falls into the class of *generative* models, one can also assess the *expressivity* of a model by analyzing the variability of generated patches and its ability to generate relevant image features (edges, texture elements etc.). This can be tested by selecting a component  $k$  of the GGMM (or GMM) with probability  $w_k$  and sampling patches from the GGD (or GD) as described in [41]. Figure 13 presents a collage of 100 patches independently generated by this procedure using GMM, GGMM, LMM and HLMM. As observed, patches generated from GGMM show greater balance between strong/faint edges, constant patches and subtle textures than the models that use constant shape parameters such as GMM, LMM and HLMM.

The superiority of our GGMM model over GMM, LMM or HLMM models can also be illustrated by comparing the log-likelihood (LL) achieved by these models over a set of clean patches from natural images. Note that, to maintain objectivity, the models have to be tested on data that is different than the dataset used during training. To this end, we compute the LL of the four above-mentioned models on all non-overlapping patches of 40 randomly selected images extracted from BSDS testing set [35], which is a different set than the training images used in the EM algorithm (parameter estimation/model learning). One can observe that not only GGMM is a better fit than GMM, LMM and HLMM on average for the 40 images, but it is also a better fit on each single image.

**Denoising evaluation.** Following the verification of the model, we provide a thorough evaluation of our GGMM prior in denoising task by comparing its performance against EPLL that uses a GMM prior [62] and with our LMM and HLMM models explained above. For the ease of comparison, we utilize the pipeline that was prescribed for the original EPLL [62] algorithm. Specifically, we subject the noisy images to 5 iterations of the denoising procedure with the estimate from one iteration fed in as the input for the next one. The iterative procedure is initialized using first estimate  $\hat{\mathbf{x}} = \mathbf{y}$  and with the parameter  $\beta$  set to  $\frac{1}{\sigma^2} \{1, 4, 8, 16, 32\}$  for

Table 2

Image denoising performance comparison of EPLL algorithm with different priors. PSNR and SSIM values are obtained on the BSDS test set (average over 60 images), and on six standard images corrupted with 3 different levels of noise. BM3D algorithm results are also included for reference purposes.

$\sigma$	Algo.	BSDS	barbara	cameraman	hill	house	lena	mandrill	Avg.
PSNR									
5	BM3D	37.33	38.30	38.33	36.03	39.79	38.70	35.24	37.37
	GMM	37.26	37.58	38.13	<b>35.92</b>	38.83	38.50	35.21	37.27
	GGMM	<b>37.33</b>	37.72	<b>38.18</b>	<b>35.92</b>	<b>38.91</b>	<b>38.52</b>	<b>35.22</b>	<b>37.34</b>
	LMM	37.31	<b>37.83</b>	38.15	35.87	38.90	38.49	35.17	37.32
	HLMM	36.85	37.42	37.70	35.40	38.33	38.08	34.75	36.86
20	BM3D	29.38	31.81	30.39	28.59	33.79	33.02	26.62	29.50
	GMM	29.35	29.74	30.10	28.45	32.77	32.37	26.62	29.41
	GGMM	<b>29.43</b>	30.08	<b>30.22</b>	<b>28.50</b>	33.01	32.59	<b>26.68</b>	<b>29.50</b>
	LMM	29.29	<b>30.21</b>	30.01	28.38	<b>33.23</b>	<b>32.72</b>	26.46	29.37
	HLMM	28.48	29.34	29.00	27.71	32.54	32.11	25.47	28.56
60	BM3D	24.81	26.26	25.39	24.49	28.73	28.16	21.71	24.90
	GMM	24.54	23.78	25.10	24.24	27.33	27.19	<b>21.56</b>	24.57
	GGMM	<b>24.63</b>	<b>23.96</b>	<b>25.27</b>	<b>24.28</b>	27.65	27.49	21.48	<b>24.67</b>
	LMM	24.55	23.88	25.10	24.27	<b>27.73</b>	<b>27.55</b>	21.33	24.59
	HLMM	23.95	23.17	23.80	23.86	26.97	26.98	20.65	23.98
SSIM									
5	BM3D	.9619	.9643	.9617	.9507	.9567	.9435	.9585	.9614
	GMM	.9627	.9617	<b>.9608</b>	<b>.9511</b>	<b>.9474</b>	<b>.9435</b>	<b>.9598</b>	<b>.9619</b>
	GGMM	<b>.9628</b>	<b>.9618</b>	.9604	.9503	.9461	.9424	.9592	<b>.9619</b>
	LMM	.9620	.9615	.9599	.9488	.9443	.9408	.9576	.9611
	HLMM	.9562	.9571	.9570	.9408	.9359	.9352	.9501	.9553
20	BM3D	.8236	.9050	.8687	.7791	.8729	.8769	.7962	.8260
	GMM	<b>.8313</b>	.8691	.8681	<b>.7813</b>	.8582	.8641	<b>.8048</b>	<b>.8322</b>
	GGMM	.8296	<b>.8728</b>	<b>.8691</b>	.7777	.8614	.8679	.8010	.8307
	LMM	.8176	.8729	.8599	.7639	<b>.8638</b>	<b>.8687</b>	.7849	.8192
	HLMM	.7884	.8502	.8426	.7325	.8577	.8594	.7383	.7907
60	BM3D	.6373	.7562	.7572	.5802	.7917	.7753	.5012	.6424
	GMM	<b>.6205</b>	.6455	.7172	<b>.5675</b>	.7386	.7281	<b>.5014</b>	<b>.6232</b>
	GGMM	.6176	<b>.6512</b>	<b>.7318</b>	.5586	.7564	.7416	.4769	.6208
	LMM	.6120	.6475	.7299	.5577	<b>.7596</b>	<b>.7443</b>	.4614	.6155
	HLMM	.5775	.6117	.7025	.5268	.7487	.7338	.3843	.5812

each iteration, respectively. To reduce the computation time of *all* EPLL-based algorithms, we utilize the random patch overlap procedure introduced by [44]. That is, instead of extracting all patches at each iteration, a randomly selected but different subset of overlapping patches consisting of only 3% of all possible patches is processed in each iteration. For the

sake of reproducibility, we made our MATLAB/MEX-C implementation available online at <https://bitbucket.org/cdeledalle/ggmm-epll>.

The evaluation is carried out on classical images such as *Barbara*, *Cameraman*, *Hill*, *House*, *Lena* and *Mandrill*, and on 60 images taken from BSDS testing set [35] (the original BSDS test data contains 100 images, the other 40 were used for model validation experiments). All image have been corrupted independently by additive white Gaussian noise with standard deviation  $\sigma = 5, 20$  and  $60$  (with pixel values between  $[0, 255]$ ). The EPLL algorithm using mixture of Gaussian, generalized Gaussian, Laplacian and hyper-Laplacian priors are indicated as GMM, GGMM, LMM and HLMM in Table 2. BM3D algorithm [9] also included for reference purposes. To stay with the focus of this paper, *i.e.*, on the effect of image priors on EPLL-based algorithms, BM3D will be excluded from our performance comparison discussions. The denoising performance of the algorithms are measured in terms of Peak Signal to Noise Ratio (PSNR) and Structural SIMilarity (SSIM) [58]. As can be observed in Table 2, in general, GGMM prior provides better PSNR performance than the three other priors. In terms of SSIM, GGMM prior is comparable to GMM. The differences in denoising performance can also be verified visually in Figure 15, Figure 16 and Figure 17. The denoised images obtained using GGMM prior show much fewer edge artifacts as compared to GMM-EPLL results. On the other hand, GGMM prior is also able to better preserve textures than LMM and HLMM.

*Prior fitness for image denoising.* In this work, we have considered non-blind image denoising. That is, the noise standard deviation is assumed to be known. In this setting, if the restoration model is accurate, one should ideally achieve optimal restoration performance when using the true degradation. To verify this, we conducted a denoising task with image corrupted with noise with standard deviation  $\sigma = 20$ . We used GMM, LMM and GGMM priors in the restoration framework with assumed  $\sigma$  values ranging from 15 to 30. Figure 18 shows the evolution of average restoration performance over 40 images from BSDS testing set (kept aside for validation, as mentioned above) with varying noise variances. GGMM prior attains its best performance when the noise variance used in the restoration model matches with the ground truth  $\sigma = 20$ . In contrast to GGMM, GMM (resp., LMM) reaches its best performance at a larger (resp., lower) value of  $\sigma$  than the correct noise used during degradation. This is because GMM tends to under-smooth clean patches (resp., over-smooth) so that a larger (resp., lower) value of  $\sigma$  is required to compensate the mismatch between the assumed prior and the actual distribution in the restoration model. This indicates that GGMM is a better option to model distribution of image patches than GMM or LMM.

*Influence of our approximations.* All previous experiments using GGMM patch priors were conducted based on the two proposed approximations introduced in Section 4 and Section 5. In Figure 19, we provide a quantitative illustration of the speed-ups provided by these approximations and their effect on denoising performance. Figure 19.(a) shows the result obtained by calculating original discrepancy function via numerical integration and the shrinkage function via Halley’s root-finding method. This makes the denoising process extremely slow and takes 6 hours and 20 minutes for denoising an image of size  $128 \times 128$  pixels. The approximated discrepancy function provides 4 orders of magnitude speed-up with no perceivable drop in performance (Figure 19.(b)). In addition, incorporating the shrinkage approximation provides further acceleration that allows the denoising to complete in less than 2 seconds with a very minor drop in PSNR/SSIM. Although the shrinkage approximation provides an acceler-

ation of ten-fold to the shrinkage calculation step, please note that the  $10\times$  speed-up is not reflected in the denoising process due to the major bottleneck caused by discrepancy function calculation. The approximately  $15,000\times$  speed-up realized without any perceivable drop in denoising performance underscores the efficacy of our proposed approximations.

**7. Conclusions and Discussion.** In this work, we suggest using a mixture of generalized Gaussians for modeling the patch distribution of clean images. We then provide a detailed study of the challenges that one encounters when using a highly flexible GGMM prior for image restoration in place of a more common GMM prior. We identify the two main bottlenecks in the restoration procedure using EPLL and GGMM namely the patch classification step and the shrinkage step. One of the main contributions of this paper, is the thorough theoretical analysis of the classification problem allowing us to introduce an asymptotically accurate approximation being computationally efficient. In order to tackle the shrinkage step, we collate and extend the existing solutions for the shrinkage function under GGMM prior.

Our numerical experiments indicate that our flexible patch GGMM is a better fit for each single image than GMM and other mixtures with constant shape parameters such as LMM or HLMM. In image denoising tasks, we have shown that using GGMM priors often outperforms GMM when used in the EPLL framework. Nevertheless, we believe the performance of GGMM prior in these scenarios falls short of its expected potential. Given GGMM is systematically a better prior than GMM, one would expect the GGMM-EPLL to outperform GMM-EPLL consistently. We postulate that even though the GGMM prior is improving the quality of the global solution, the half quadratic splitting strategy used in EPLL has no guarantee to return a better solution due to the non-convexity of the underlying problem. For this reason, we will focus our future work in designing specific optimization strategies for GGMM-EPLL leveraging the better expressivity of the model. Also, preliminary results on inverse problems indicates that degradation- and prior-specific algorithmic design should be employed to make GGMM-EPLL more competitive against GMM-EPLL (otherwise they both are on par). The development of such a formulation designed to adapt according to the properties of chosen priors and presented degradations is still under investigation.

**Acknowledgments.** The authors would like to thank Charles Dossal and Nicolas Papadakis for fruitful discussions.

#### Appendix A. Proof of equation (26).

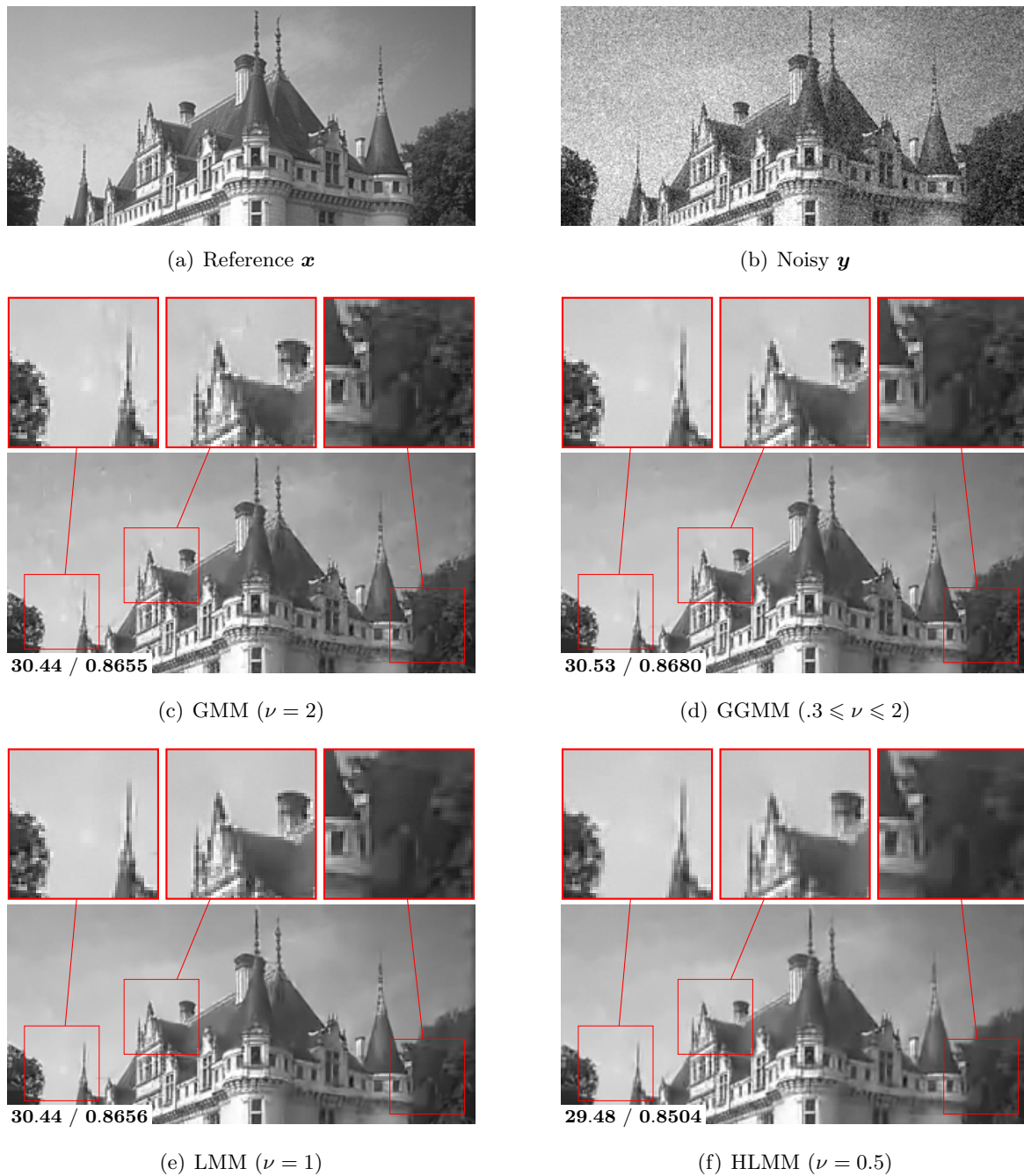
*Proof.* For  $\nu = 1$ , using  $\Gamma(1) = 1$  and  $\Gamma(3) = 2$ , we obtain

$$(44) \quad f_{\sigma,\lambda}^1(x) = \log(2\sqrt{\pi}\sigma\lambda) - \log \int_{-\infty}^{\infty} e^{-\frac{t^2}{2\sigma^2} - \frac{\sqrt{2}|x-t|}{\lambda}} dt ,$$

$$(45) \quad = \log(2\sqrt{\pi}\sigma\lambda) - \log \left[ e^{-\frac{\sqrt{2}x}{\lambda}} \int_{-\infty}^x e^{-\frac{t^2}{2\sigma^2} + \frac{\sqrt{2}t}{\lambda}} dt + e^{\frac{\sqrt{2}x}{\lambda}} \int_x^{\infty} e^{-\frac{t^2}{2\sigma^2} - \frac{\sqrt{2}t}{\lambda}} dt \right] .$$

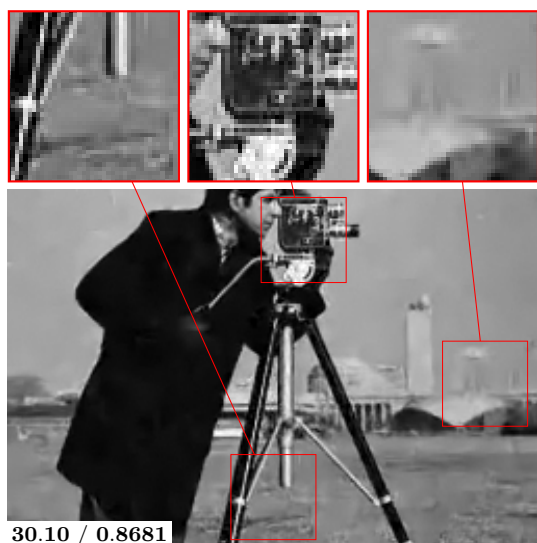
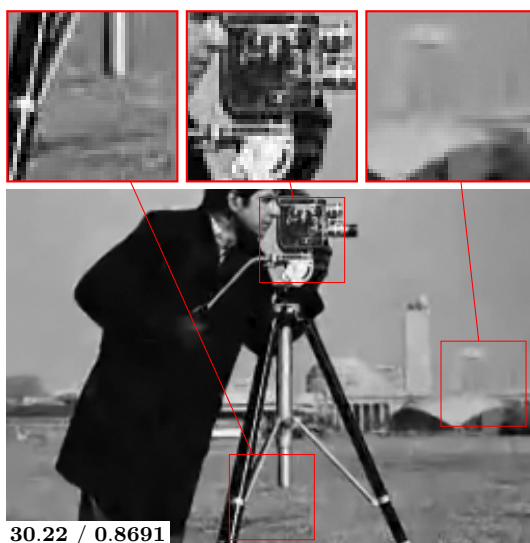
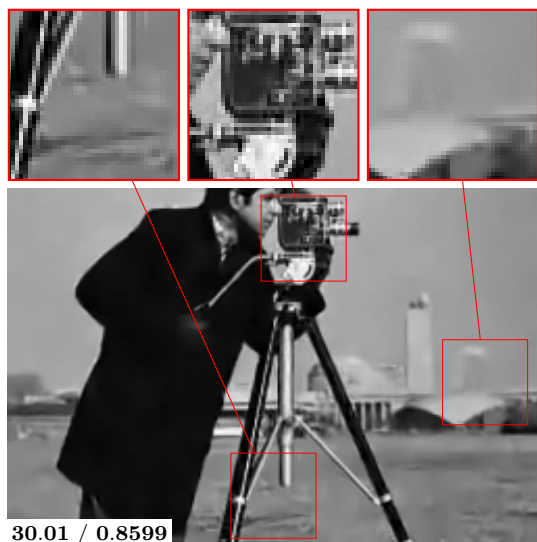
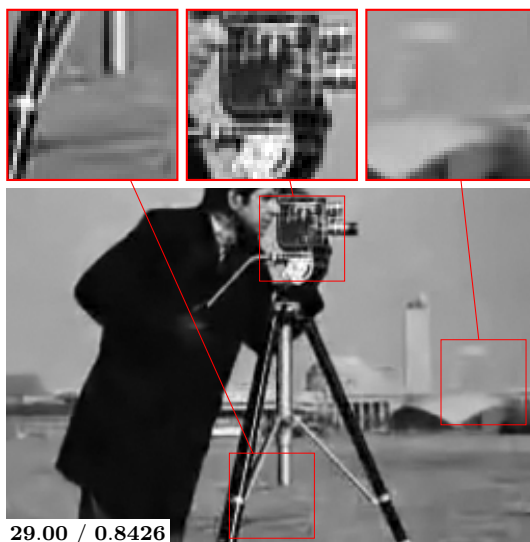
Since  $\operatorname{erf}'(t) = \frac{2e^{-t^2}}{\sqrt{\pi}}$ , it follows that for  $a > 0$  and  $b > 0$

$$(46) \quad \frac{\partial}{\partial t} \left[ -\frac{\sqrt{\pi a}}{2} e^{\frac{a}{4b^2}} \operatorname{erf} \left( -\frac{t}{\sqrt{a}} - \frac{\sqrt{a}}{2b} \right) \right] = e^{-\frac{t^2}{a} - \frac{t}{b}} .$$



**Figure 15.** (a) Close in on the image Castle from the BSDS testing dataset, (b) a noisy version degraded by additive white Gaussian noise with standard deviation  $\sigma = 20$  and (c)-(f) results of EPLL under four patch priors: GMM, GGMM, LMM and HLMM, respectively. PSNR and SSIM are given in the bottom-left corner.



(a) Reference  $\mathbf{x}$ (b) Noisy  $\mathbf{y}$ (c) GMM ( $\nu = 2$ )(d) GGMM ( $.3 \leq \nu \leq 2$ )(e) LMM ( $\nu = 1$ )(f) HLMM ( $\nu = 0.5$ )

**Figure 16.** (a) Close in on the standard image *Cameraman*. (b) a noisy version degraded by additive white Gaussian noise with standard deviation  $\sigma = 20$  and (c)-(f) results of EPLL under four patch priors: GMM, GGMM, LMM and HLMM, respectively. PSNR and SSIM are given in the bottom-left corner.





(a) Reference  $x$



(b) Noisy  $y$



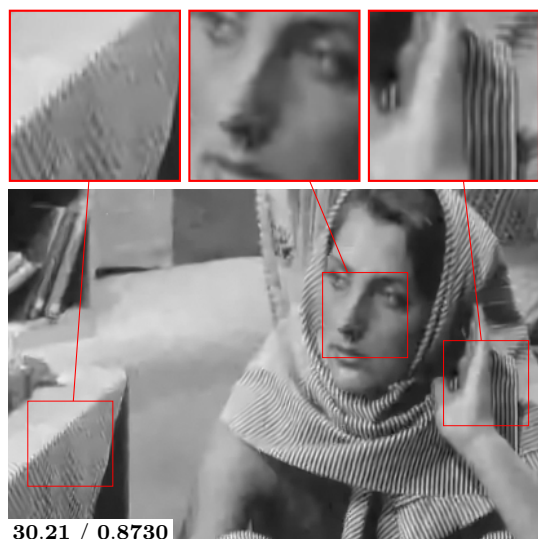
29.74 / 0.8691

(c) GMM ( $\nu = 2$ )



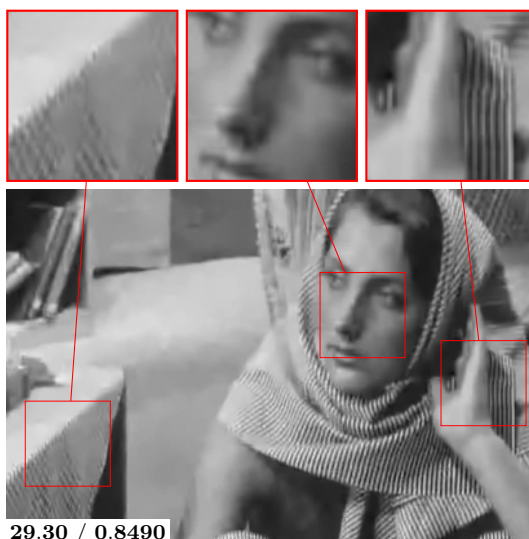
30.07 / 0.8728

(d) GGMM ( $.3 \leq \nu \leq 2$ )



30.21 / 0.8730

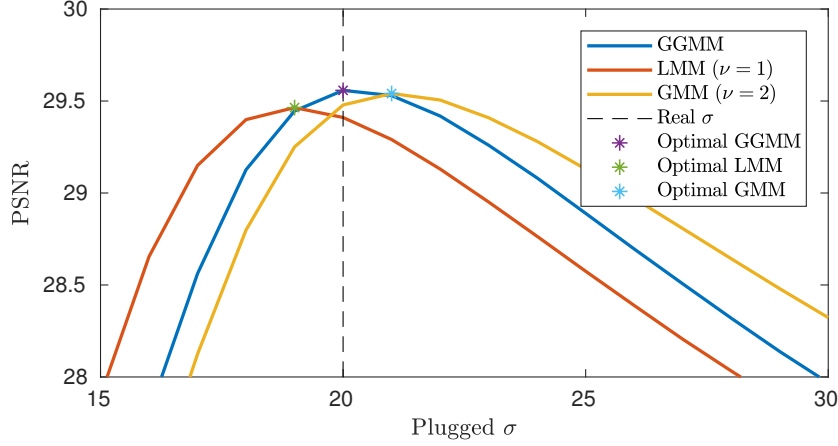
(e) LMM ( $\nu = 1$ )



29.30 / 0.8490

(f) HLMM ( $\nu = 0.5$ )

**Figure 17.** (a) Close in on the standard image Barbara. (b) a noisy version degraded by additive white Gaussian noise with standard deviation  $\sigma = 20$  and (c)-(f) results of EPLL under four patch priors: GMM, GGMM, LMM and HLMM, respectively. PSNR and SSIM are given in the bottom-left corner.



**Figure 18.** Evolution of performance of EPLL with a GGMM, LMM ( $\nu = 1$ ) and a GMM ( $\nu = 2$ ) under misspecification of the noise standard deviation  $\sigma$ . Performances are measured in terms of PSNR on the BSDS dataset corrupted by a Gaussian noise with standard deviation  $\sigma = 20$ . For each of the three priors, EPLL has been run assuming  $\sigma$  was ranging from 15 to 30.

Therefore we have with  $a = 2\sigma^2$  and  $b = \lambda/\sqrt{2}$

$$(47) \quad \int_{-\infty}^x e^{-\frac{t^2}{2\sigma^2} + \frac{\sqrt{2}t}{\lambda}} dt = \frac{-\sqrt{\pi}\sigma e^{\frac{\sigma^2}{\lambda^2}}}{\sqrt{2}} \left[ \operatorname{erf} \left( -\frac{t}{\sqrt{2}\sigma} + \frac{\sigma}{\lambda} \right) \right]_{-\infty}^x = \frac{\sqrt{\pi}\sigma e^{\frac{\sigma^2}{\lambda^2}}}{\sqrt{2}} \operatorname{erfc} \left( -\frac{x}{\sqrt{2}\sigma} + \frac{\sigma}{\lambda} \right),$$

since  $\lim_{t \rightarrow \infty} \operatorname{erf}(t) = 1$  and  $\operatorname{erfc}(t) = 1 - \operatorname{erf}(t)$ . Similarly, we get

$$(48) \quad \int_x^{\infty} e^{-\frac{t^2}{2\sigma^2} - \frac{\sqrt{2}t}{\lambda}} dt = \frac{-\sqrt{\pi}\sigma e^{\frac{\sigma^2}{\lambda^2}}}{\sqrt{2}} \operatorname{erfc} \left( \frac{x}{\sqrt{2}\sigma} + \frac{\sigma}{\lambda} \right).$$

Plugging these two last expressions in (45) and rearranging the terms conclude the proof. ■

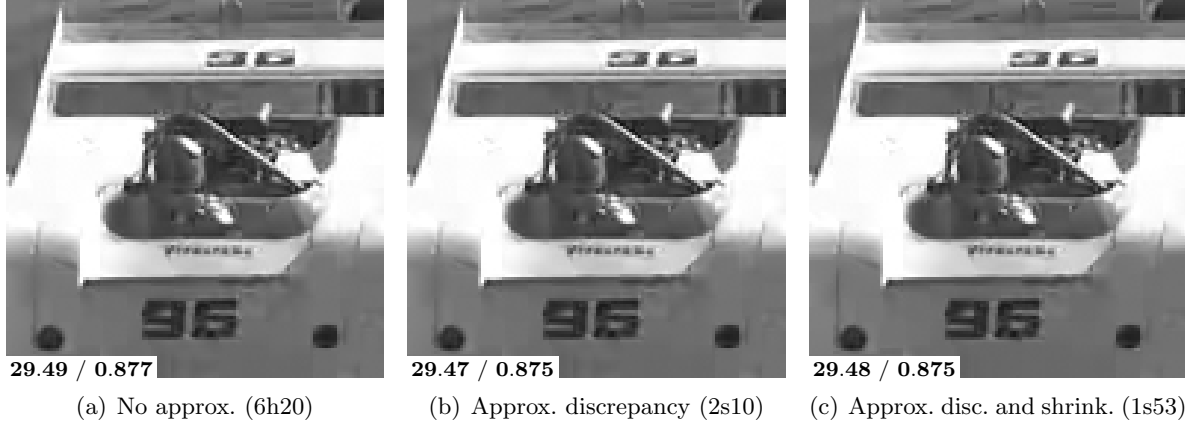
## Appendix B. Proof of Proposition 2.

*Proof.* Starting from the definition of  $f_{\sigma,\lambda}^\nu$  and using the change of variable  $t \rightarrow \sigma t$ , eq. (reduction) follows as

$$(49) \quad f_{\sigma,\lambda}^\nu(x) = -\log \int_{\mathbb{R}} \sigma \frac{\kappa}{2\lambda^\nu} \exp \left[ -\left( \frac{\sigma|t|}{\lambda^\nu} \right)^\nu \right] \frac{1}{\sqrt{2\pi}\sigma} \exp \left[ -\frac{(x - \sigma t)^2}{2\sigma^2} \right] dt,$$

$$= \log \sigma - \log \int_{\mathbb{R}} \frac{\kappa}{2\lambda^\nu/\sigma} \exp \left[ -\left( \frac{|t|}{\lambda^\nu/\sigma} \right)^\nu \right] \frac{1}{\sqrt{2\pi}} \exp \left[ -\frac{(x/\sigma - t)^2}{2} \right] dt = \log \sigma + f_{1,\lambda/\sigma}^\nu \left( \frac{x}{\sigma} \right).$$

Properties (even) and (unimodality) hold since the convolution of two real even unimodal distributions is even unimodal [59, 48]. Property (lower bound at 0) follows from (even), (unimodality) and the fact that the convolution of continuous and bounded real functions are continuous and bounded. ■



**Figure 19.** Results obtained by GGMM-EPLL on a  $128 \times 128$  cropped image of the BSDS testing dataset damaged by additive white Gaussian noise with  $\sigma = 20$ . These are obtained respectively by (a) evaluating the classification and shrinkage problem with numerical solvers (numerical integration and Halley's root-finding method), (b) approximating the classification problem only, (c) approximating both problems. PSNR and SSIM are given in the bottom-left corner. Running time are indicated on the captions: our accelerations lead to a speed-up of  $\times 15,000$  and  $\times 1.4$  respectively.

### Appendix C. Proof of Theorem 1.

**Lemma C.1.** Let  $a > 0$  and  $b > 0$ . For  $x$  in the vicinity of 0, we have

$$(50) \quad \frac{1}{2abx} \log \left[ \frac{\operatorname{erfc}(ax+b) + e^{-4abx} \operatorname{erfc}(-ax+b)}{2 \operatorname{erfc}(b)} \right] = -1 + \left( ab - \frac{ae^{-b^2}}{\operatorname{erfc}(b)\sqrt{\pi}} \right) x + o(x).$$

*Proof.* Since  $\operatorname{erfc}'(x) = -\frac{2e^{-x^2}}{\sqrt{\pi}}$  and  $\operatorname{erfc}''(x) = \frac{2xe^{-x^2}}{\sqrt{\pi}}$ , using second order Taylor's expansion for  $x$  in the vicinity of 0, it follows that

$$(51) \quad \operatorname{erfc}(ax+b) \underset{0}{\sim} \operatorname{erfc}(b) - \frac{2ae^{-b^2}}{\sqrt{\pi}}x + \frac{2a^2be^{-b^2}}{\sqrt{\pi}}x^2,$$

$$(52) \quad \text{and } \operatorname{erfc}(-ax+b) \underset{0}{\sim} \operatorname{erfc}(b) + \frac{2ae^{-b^2}}{\sqrt{\pi}}x + \frac{2a^2be^{-b^2}}{\sqrt{\pi}}x^2.$$

We next make the following deductions

$$(53) \quad e^{-4abx} \operatorname{erfc}(-ax+b) \underset{0}{\sim} e^{-4abx} \operatorname{erfc}(b) + e^{-4abx} \frac{2ae^{-b^2}}{\sqrt{\pi}}x + e^{-4abx} \frac{2a^2be^{-b^2}}{\sqrt{\pi}}x^2,$$

$$\operatorname{erfc}(ax+b) + e^{-4abx} \operatorname{erfc}(-ax+b) \underset{0}{\sim} (1+e^{-4abx}) \left( \operatorname{erfc}(b) + \frac{2a^2be^{-b^2}}{\sqrt{\pi}}x^2 \right) - (1-e^{-4abx}) \frac{2ae^{-b^2}}{\sqrt{\pi}}x,$$

$$\underbrace{\frac{\operatorname{erfc}(ax+b) + e^{-4abx} \operatorname{erfc}(-ax+b)}{2 \operatorname{erfc} b}}_{A(x)} \underset{0}{\sim} \frac{1+e^{-4abx}}{2} \left( 1 + \frac{2a^2be^{-b^2}}{\operatorname{erfc}(b)\sqrt{\pi}}x^2 \right) - (1-e^{-4abx}) \frac{ae^{-b^2}}{\operatorname{erfc}(b)\sqrt{\pi}}x.$$

The left-hand side  $A(x)$  of this last equation is then, in the vicinity of  $x = 0$ , equals to

$$(54) \quad A(x) = \frac{1 + e^{-4abx}}{2} \left( 1 + \frac{2a^2be^{-b^2}}{\operatorname{erfc}(b)\sqrt{\pi}}x^2 \right) - (1 - e^{-4abx}) \frac{ae^{-b^2}}{\operatorname{erfc}(b)\sqrt{\pi}}x + o(x^2).$$

We next use second-order Taylor's expansion for  $e^{-4abx}$ , leading to

$$(55) \quad A(x) = (1 - 2abx + 4a^2b^2x^2 + o(x^2)) \left( 1 + \frac{2a^2be^{-b^2}}{\operatorname{erfc}(b)\sqrt{\pi}}x^2 \right) \\ - (4abx - 8a^2b^2x^2 + o(x^2)) \frac{ae^{-b^2}}{\operatorname{erfc}(b)\sqrt{\pi}}x + o(x^2),$$

$$(56) \quad = 1 - 2abx + \left( 4a^2b^2 - \frac{2a^2be^{-b^2}}{\operatorname{erfc}(b)\sqrt{\pi}} \right) x^2 + o(x^2).$$

By using the second-order Taylor's expansion of  $\log(1+x)$ , it follows that

$$(57) \quad \log[A(x)] = -2abx + \left( 4a^2b^2 - \frac{2a^2be^{-b^2}}{\operatorname{erfc}(b)\sqrt{\pi}} \right) x^2 - 2a^2b^2x^2 + o(x^2).$$

Dividing both sides by  $2abx$  then concludes the proof,

$$(58) \quad \frac{1}{2abx} \log \left[ \frac{\operatorname{erfc}(ax+b) + e^{-4abx} \operatorname{erfc}(-ax+b)}{2 \operatorname{erfc} b} \right] = -1 + \left( ab - \frac{ae^{-b^2}}{\operatorname{erfc}(b)\sqrt{\pi}} \right) x + o(x). \blacksquare$$

*Proof of Theorem 1.* We first rewrite  $\varphi_\lambda^1(x)$  as

$$(59) \quad \varphi_\lambda^1(x) = \log \left[ \log \left[ 2 \operatorname{erfc} \left( \frac{1}{\lambda} \right) \right] - \log \left[ e^{\frac{\sqrt{2}x}{\lambda}} \operatorname{erfc} \left( \frac{x}{\sqrt{2}} + \frac{1}{\lambda} \right) + e^{-\frac{\sqrt{2}x}{\lambda}} \operatorname{erfc} \left( -\frac{x}{\sqrt{2}} + \frac{1}{\lambda} \right) \right] \right], \\ = \log \left[ \frac{\sqrt{2}x}{\lambda} \right] + \log \left[ -1 - \frac{\lambda}{\sqrt{2}x} \log \left[ \frac{\operatorname{erfc} \left( \frac{x}{\sqrt{2}} + \frac{1}{\lambda} \right) + e^{-\frac{2\sqrt{2}x}{\lambda}} \operatorname{erfc} \left( -\frac{x}{\sqrt{2}} + \frac{1}{\lambda} \right)}{2 \operatorname{erfc} \left( \frac{1}{\lambda} \right)} \right] \right].$$

Next, using [Lemma C.1](#) with  $a = 1/\sqrt{2}$  and  $b = 1/\lambda$ , it follows that

$$(60) \quad \varphi_\lambda^1(x) = \log \left[ \frac{\sqrt{2}x}{\lambda} \right] + \log \left[ - \left( \frac{1}{\sqrt{2}\lambda} - \frac{e^{-\frac{1}{\lambda^2}}}{\sqrt{2} \operatorname{erfc}(\frac{1}{\lambda})\sqrt{\pi}} \right) x + o(x) \right],$$

$$(61) \quad = \log \left[ \frac{x^2}{\lambda} \right] + \log \left[ \frac{e^{-\frac{1}{\lambda^2}}}{\operatorname{erfc}(\frac{1}{\lambda})\sqrt{\pi}} - \frac{1}{\lambda} + o(1) \right],$$

$$(62) \quad = \log \left[ \frac{x^2}{\lambda} \right] + \log \left[ \frac{1}{\sqrt{\pi}} \frac{e^{-\frac{1}{\lambda^2}}}{\operatorname{erfc}(\frac{1}{\lambda})} - \frac{1}{\lambda} \right] + o(1),$$

where the last equation follows from the first-order Taylor expansion of  $\log(a+x)$ .  $\blacksquare$

### Appendix D. Proof of Theorem 2.

Lemma D.1. Let  $a > 0$  and  $b > 0$ . For  $x$  in the vicinity of  $+\infty$ , we have

$$(63) \quad \frac{1}{2abx} \log \left[ \frac{\operatorname{erfc}(ax+b) + e^{-4abx} \operatorname{erfc}(-ax+b)}{2} \right] = -2 + o(1) .$$

*Proof.* We have  $\lim_{x \rightarrow +\infty} \operatorname{erfc}(x) = 0$  and  $\lim_{x \rightarrow +\infty} \operatorname{erfc}(-x) = 2$ , it follows that

$$(64) \quad \operatorname{erfc}(-ax+b) \underset{+\infty}{\sim} 2 ,$$

$$(65) \quad e^{-4abx} \operatorname{erfc}(-ax+b) \underset{+\infty}{\sim} 2e^{-4abx} ,$$

$$(66) \quad \frac{\operatorname{erfc}(ax+b) + e^{-4abx} \operatorname{erfc}(-ax+b)}{2} \underset{+\infty}{\sim} e^{-4abx} ,$$

$$(67) \quad \log \left[ \frac{\operatorname{erfc}(ax+b) + e^{-4abx} \operatorname{erfc}(-ax+b)}{2} \right] \underset{+\infty}{\sim} -4abx ,$$

$$(68) \quad \frac{1}{2abx} \log \left[ \operatorname{erfc}(ax+b) + e^{-4abx} \operatorname{erfc}(-ax+b) \right] \underset{+\infty}{\sim} -2 ,$$

where we have used the knowledge that  $f \sim g$  implies that  $\log f \sim \log g$ . ■

*Proof of Theorem 2.* By writing  $\varphi_\lambda^1$  as in eq. (59) and using Lemma D.1 with  $a = 1/\sqrt{2}$  and  $b = 1/\lambda$ , it follows that

$$\varphi_\lambda^1(x) = \log \left[ \frac{\sqrt{2}x}{\lambda} \right] + \log \left[ 1 + \frac{\lambda}{\sqrt{2}x} \log \operatorname{erfc} \left( \frac{1}{\lambda} \right) + o(1) \right] = \log \left[ \frac{\sqrt{2}x}{\lambda} \right] + o(1) . \quad \blacksquare$$

### Appendix E. Proof of Theorem 3.

*Proof.* We first decompose the term  $\exp\left(-\frac{(x-t)^2}{2}\right)$  involved in the definition of  $f_{1,\lambda}^\nu$  and rewrite  $e^{xt}$  using its power series

$$(69) \quad f_{1,\lambda}^\nu(x) = -\log \frac{1}{\sqrt{2\pi}} \frac{\nu}{2\lambda_\nu \Gamma(1/\nu)} - \log \int_{-\infty}^{\infty} \exp\left(-\frac{(x-t)^2}{2}\right) \exp\left[-\left(\frac{|t|}{\lambda_\nu}\right)^\nu\right] dt ,$$

$$(70) \quad = -\log \frac{1}{\sqrt{2\pi}} \frac{\nu}{2\lambda_\nu \Gamma(1/\nu)} - \log \int_{-\infty}^{\infty} e^{-\frac{x^2}{2}} e^{xt} e^{-\frac{t^2}{2}} \exp\left[-\left(\frac{|t|}{\lambda_\nu}\right)^\nu\right] dt ,$$

$$(71) \quad = -\log \frac{1}{\sqrt{2\pi}} \frac{\nu}{2\lambda_\nu \Gamma(1/\nu)} + \frac{x^2}{2} - \log \int_{-\infty}^{\infty} \sum_{k=0}^{\infty} \frac{(xt)^k}{k!} e^{-\frac{t^2}{2}} \exp\left[-\left(\frac{|t|}{\lambda_\nu}\right)^\nu\right] dt .$$

For  $x$  in the vicinity of 0, we can consider  $|x| \leq 1$ , and then

$$(72) \quad \int_{-\infty}^{\infty} \sum_{k=0}^{\infty} \left| \frac{(xt)^k}{k!} e^{-\frac{t^2}{2}} \exp\left[-\left(\frac{|t|}{\lambda_\nu}\right)^\nu\right] \right| dt \leq \int_{-\infty}^{\infty} \sum_{k=0}^{\infty} \frac{|t|^k}{k!} e^{-\frac{t^2}{2}} \exp\left[-\left(\frac{|t|}{\lambda_\nu}\right)^\nu\right] dt ,$$

$$(73) \quad \leq \int_{-\infty}^{\infty} e^{-\frac{t^2}{2} + |t| - \left(\frac{|t|}{\lambda_\nu}\right)^\nu} dt < \infty .$$

Then, Fubini's theorem applies and we get

$$(74) \quad f_{1,\lambda}^\nu(x) = -\log \frac{1}{\sqrt{2\pi}} \frac{\nu}{2\lambda\nu\Gamma(1/\nu)} + \frac{x^2}{2} - \log \sum_{k=0}^{\infty} \int_{-\infty}^{\infty} \frac{(xt)^k}{k!} e^{-\frac{t^2}{2}} \exp \left[ -\left(\frac{|t|}{\lambda\nu}\right)^\nu \right] dt ,$$

$$(75) \quad = -\log \frac{1}{\sqrt{2\pi}} \frac{\nu}{2\lambda\nu\Gamma(1/\nu)} + \frac{x^2}{2} - \log \sum_{k=0}^{\infty} \frac{x^k}{k!} \int_{-\infty}^{\infty} t^k e^{-\frac{t^2}{2}} \exp \left[ -\left(\frac{|t|}{\lambda\nu}\right)^\nu \right] dt .$$

By definition, we have

$$(76) \quad \gamma_\lambda^\nu \triangleq f_{1,\lambda}^\nu(0) = -\log \frac{1}{\sqrt{2\pi}} \frac{\nu}{2\lambda\nu\Gamma(1/\nu)} - \log \int_{-\infty}^{\infty} \exp \left( -\frac{t^2}{2} \right) \exp \left[ -\left(\frac{|t|}{\lambda\nu}\right)^\nu \right] dt .$$

Moreover, when  $k$  is odd, we have

$$(77) \quad \int_{-\infty}^{\infty} t^k e^{-\frac{t^2}{2}} \exp \left[ -\left(\frac{|t|}{\lambda\nu}\right)^\nu \right] dt = 0 .$$

Using third-order Taylor expansion of  $\log(1+x)$  for  $x$  in the vicinity of 0, it follows that

$$(78) \quad f_{1,\lambda}^\nu(x) = \gamma_\lambda^\nu + \frac{x^2}{2} - \log \left( 1 + \frac{x^2}{2} \frac{\int_{-\infty}^{\infty} t^2 e^{-\frac{t^2}{2}} \exp \left[ -\left(\frac{|t|}{\lambda\nu}\right)^\nu \right] dt}{\int_{-\infty}^{\infty} e^{-\frac{t^2}{2}} \exp \left[ -\left(\frac{|t|}{\lambda\nu}\right)^\nu \right] dt} + o(x^3) \right) ,$$

$$(79) \quad = \gamma_\lambda^\nu + \frac{x^2}{2} \left( 1 - \frac{\int_{-\infty}^{\infty} t^2 e^{-\frac{t^2}{2}} \exp \left[ -\left(\frac{|t|}{\lambda\nu}\right)^\nu \right] dt}{\int_{-\infty}^{\infty} e^{-\frac{t^2}{2}} \exp \left[ -\left(\frac{|t|}{\lambda\nu}\right)^\nu \right] dt} \right) + o(x^3) .$$

Finally, using first-order Taylor's expansion for  $\log(1+x)$ , we conclude the proof as

$$(80) \quad \varphi_\lambda^\nu(x) = \log \left[ \frac{x^2}{2} \left( 1 - \frac{\int_{-\infty}^{\infty} t^2 e^{-\frac{t^2}{2}} \exp \left[ -\left(\frac{|t|}{\lambda\nu}\right)^\nu \right] dt}{\int_{-\infty}^{\infty} e^{-\frac{t^2}{2}} \exp \left[ -\left(\frac{|t|}{\lambda\nu}\right)^\nu \right] dt} + o(x) \right) \right] ,$$

$$(81) \quad = 2 \log x - \log 2 + \log \left( 1 - \frac{\int_{-\infty}^{\infty} t^2 e^{-\frac{t^2}{2}} \exp \left[ -\left(\frac{|t|}{\lambda\nu}\right)^\nu \right] dt}{\int_{-\infty}^{\infty} e^{-\frac{t^2}{2}} \exp \left[ -\left(\frac{|t|}{\lambda\nu}\right)^\nu \right] dt} \right) + o(x) . \quad \blacksquare$$

#### Appendix F. Proof of Theorem 4.

We first recall in a Lemma, a result extracted from Corollary 3.3 in [2].

**Lemma F.1 (Berman).** *Let  $p$  and  $q$  be differentiable real probability density functions. Define for  $x$  large enough  $u(x) = p^{-1}(q(x))$  and define  $v$  and  $w$  as*

$$(82) \quad v(x) = -\frac{\partial}{\partial x} \log p(x) \quad \text{and} \quad w(x) = -\frac{\partial}{\partial x} \log q(x) .$$

Assume  $v$  and  $w$  are positive continuous function and regularly oscillating, i.e.:

$$(83) \quad \lim_{\substack{x, x' \rightarrow \infty \\ x/x' \rightarrow 1}} \frac{v(x)}{v(x')} = 1 \quad \text{and} \quad \lim_{\substack{x, x' \rightarrow \infty \\ x/x' \rightarrow 1}} \frac{w(x)}{w(x')} = 1 .$$

Suppose that we have

$$(84) \quad \lim_{x \rightarrow \infty} \frac{w(x)}{v(x)} = 0 \quad \text{and} \quad \lim_{x \rightarrow \infty} u(x)w(x) = +\infty ,$$

then, for  $x \rightarrow \infty$ , we have

$$(85) \quad \log \int_{-\infty}^{+\infty} p(x-t)q(t) dt \sim \log q(x) .$$

*Proof of Theorem 4.* Using the definition of the discrepancy function,

$$(86) \quad f_{1,\lambda}^\nu(x) = -\log \int_{\mathbb{R}} \mathcal{G}(t; 0, \lambda, \nu) \cdot \mathcal{N}(x-t; 0, 1) dt = -\log \int_{-\infty}^{+\infty} p(x-t)q(t) dt .$$

with  $q(x) = \mathcal{G}(x; 0, \lambda, \nu)$  and  $p(x) = \mathcal{N}(x; 0, 1)$ . We have

$$(87) \quad v(x) = -\frac{\partial}{\partial x} \log p(x) = 2x \quad \text{and} \quad w(x) = -\frac{\partial}{\partial x} \log q(x) = \frac{\nu x^{\nu-1}}{\lambda_\nu^\nu} .$$

Remark that  $v$  and  $w$  are positive continuous, and

$$(88) \quad \lim_{\substack{x, x' \rightarrow \infty \\ x/x' \rightarrow 1}} \frac{v(x)}{v(x')} = \lim_{\substack{x, x' \rightarrow \infty \\ x/x' \rightarrow 1}} \frac{x}{x'} = 1 \quad \text{and} \quad \lim_{\substack{x, x' \rightarrow \infty \\ x/x' \rightarrow 1}} \frac{w(x)}{w(x')} = \lim_{\substack{x, x' \rightarrow \infty \\ x/x' \rightarrow 1}} \left(\frac{x}{x'}\right)^{\nu-1} = 1 .$$

For  $x > 0$  large enough and  $y > 0$  small enough

$$(89) \quad y = p(x) \Leftrightarrow y = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right) \Leftrightarrow x = \sqrt{-\log(2\pi) - 2 \log y} .$$

Then, from [Proposition 1](#), we have

$$(90) \quad u(x)w(x) = \frac{\nu x^{\nu-1}}{\lambda_\nu^\nu} \sqrt{-\log(2\pi) - 2 \log q(x)}$$

$$(91) \quad = \frac{\nu x^{\nu-1}}{\lambda_\nu^\nu} \sqrt{-\log(2\pi) - 2 \log \left[\frac{\kappa}{2\lambda_\nu}\right] + 2 \left(\frac{x}{\lambda_\nu}\right)^\nu} \sim \frac{\nu \sqrt{2} x^{\frac{3}{2}\nu-1}}{\lambda_\nu^{\frac{3}{2}\nu}} ,$$

and thus, as  $\nu > \frac{2}{3}$ ,  $\lim_{x \rightarrow \infty} u(x)w(x) = \infty$ . Moreover, for  $\nu < 2$ , we have

$$(92) \quad \lim_{x \rightarrow \infty} \frac{w(x)}{v(x)} = \lim_{x \rightarrow \infty} \frac{\nu}{2\lambda_\nu} x^{\nu-2} = 0 .$$

It follows that [Lemma F.1](#) applies, and then for large  $x$

$$(93) \quad f_{1,\lambda}^\nu(x) \sim -\log q(x) \sim \left(\frac{x}{\lambda_\nu}\right)^\nu .$$

Using that  $\lambda_\nu = \lambda \sqrt{\frac{\Gamma(1/\nu)}{\Gamma(3/\nu)}}$ , we conclude the proof since

$$(94) \quad \varphi_\lambda^\nu(x) = \log [f_{1,\lambda}^\nu(x) - \gamma_\lambda^\nu] \sim \nu \log x - \nu \log \lambda_\nu . \quad \blacksquare$$

### Appendix G. Proof of Proposition 3.

*Proof.* Starting from the definition of  $s_{\sigma,\lambda}^\nu$  and using the change of variable  $t \rightarrow \sigma t$ , eq. [\(reduction\)](#) follows as

$$(95) \quad s_{\sigma,\lambda}^\nu(x) = \operatorname{argmin}_{t \in \mathbb{R}} \frac{(x-t)^2}{2\sigma^2} + \lambda_\nu^{-\nu} |t|^\nu = \sigma \operatorname{argmin}_{t \in \mathbb{R}} \frac{(x-\sigma t)^2}{2\sigma^2} + \lambda_\nu^{-\nu} |\sigma t|^\nu ,$$

$$(96) \quad = \sigma \operatorname{argmin}_{t \in \mathbb{R}} \frac{(x/\sigma - t)^2}{2} + (\lambda_\nu/\sigma)^{-\nu} |t|^\nu = \sigma s_{1,\lambda/\sigma}^\nu(x/\sigma) .$$

For eq. [\(odd\)](#), we use the change of variable  $t \rightarrow -t$

$$(97) \quad s_{\sigma,\lambda}^\nu(-x) = \operatorname{argmin}_{t \in \mathbb{R}} \frac{(-x-t)^2}{2\sigma^2} + \lambda_\nu^{-\nu} |t|^\nu = -\operatorname{argmin}_{t \in \mathbb{R}} \frac{(-x+t)^2}{2\sigma^2} + \lambda_\nu^{-\nu} |t|^\nu ,$$

$$(98) \quad = -\operatorname{argmin}_{t \in \mathbb{R}} \frac{(x-t)^2}{2\sigma^2} + \lambda_\nu^{-\nu} |t|^\nu = -s_{\sigma,\lambda}^\nu(x) .$$

We now prove eq. [\(shrinkage\)](#). Let  $t = s_{\sigma,\lambda}^\nu(x)$ , and since  $t$  minimizes the objective, then

$$(99) \quad \lambda_\nu^{-\nu} |t|^\nu \leq \frac{(x-t)^2}{2\sigma^2} + \lambda_\nu^{-\nu} |t|^\nu \leq \frac{(x-x)^2}{2\sigma^2} + \lambda_\nu^{-\nu} |x|^\nu = \lambda_\nu^{-\nu} |x|^\nu ,$$

which implies that  $|t| \leq |x|$ . Let  $x > 0$  and assume  $t = s_{\sigma,\lambda}^\nu(x) < 0$ . Since  $t$  minimizes the objective, then

$$(100) \quad \frac{(x-t)^2}{2\sigma^2} + \lambda_\nu^{-\nu} |t|^\nu \leq \frac{(x+t)^2}{2\sigma^2} + \lambda_\nu^{-\nu} |t|^\nu$$

which implies that  $-xt \leq xt$  and leads to a contradiction. Then for  $x > 0$ ,  $s_{\sigma,\lambda}^\nu(x) \in [0, x]$ , which concludes the proof since  $s_{\sigma,\lambda}^\nu$  is odd.

We now prove [\(increasing with x\)](#). Let  $x_1 > x_2$  and define  $t_1 = s_{\sigma,\lambda}^\nu(x_1)$  and  $t_2 = s_{\sigma,\lambda}^\nu(x_2)$ . Since  $t_1$  and  $t_2$  minimize their respective objectives, the following two statements hold

$$(101) \quad \frac{(x_1 - t_1)^2}{2\sigma^2} + \lambda_\nu^{-\nu} |t_1|^\nu \leq \frac{(x_1 - t_2)^2}{2\sigma^2} + \lambda_\nu^{-\nu} |t_2|^\nu ,$$

$$(102) \quad \text{and} \quad \frac{(x_2 - t_2)^2}{2\sigma^2} + \lambda_\nu^{-\nu} |t_2|^\nu \leq \frac{(x_2 - t_1)^2}{2\sigma^2} + \lambda_\nu^{-\nu} |t_1|^\nu .$$



Summing both inequalities lead to

$$(103) \quad (x_1 - t_1)^2 + (x_2 - t_2)^2 \leq (x_1 - t_2)^2 + (x_2 - t_1)^2 ,$$

$$(104) \quad \Rightarrow -2x_1t_1 - 2x_2t_2 \leq -2x_1t_2 - 2x_2t_1 ,$$

$$(105) \quad \Rightarrow t_1(x_1 - x_2) \geq t_2(x_1 - x_2) \quad \Rightarrow \quad t_1 \geq t_2 \quad (\text{since } x_1 > x_2) .$$

We now prove (**increasing with  $\lambda$** ). Let  $\lambda_1 > \lambda_2$  and define  $t_1 = s_{\sigma, \lambda_1}^\nu(x)$  and  $t_2 = s_{\sigma, \lambda_2}^\nu(x)$ . Since  $t_1$  and  $t_2$  minimize their respective objectives, the following expressions hold

$$(106) \quad \frac{(x - t_1)^2}{2\sigma^2} + \lambda_{\nu, 1}^{-\nu} |t_1|^\nu \leq \frac{(x - t_2)^2}{2\sigma^2} + \lambda_{\nu, 1}^{-\nu} |t_2|^\nu ,$$

$$(107) \quad \text{and} \quad \frac{(x - t_2)^2}{2\sigma^2} + \lambda_{\nu, 2}^{-\nu} |t_2|^\nu \leq \frac{(x - t_1)^2}{2\sigma^2} + \lambda_{\nu, 2}^{-\nu} |t_1|^\nu .$$

Again, summing both inequalities lead to

$$(108) \quad \lambda_{\nu, 1}^{-\nu} |t_1|^\nu + \lambda_{\nu, 2}^{-\nu} |t_2|^\nu \leq \lambda_{\nu, 1}^{-\nu} |t_2|^\nu + \lambda_{\nu, 2}^{-\nu} |t_1|^\nu ,$$

$$(109) \quad \Rightarrow (\lambda_{\nu, 1}^{-\nu} - \lambda_{\nu, 2}^{-\nu}) |t_1|^\nu \leq (\lambda_{\nu, 1}^{-\nu} - \lambda_{\nu, 2}^{-\nu}) |t_2|^\nu ,$$

$$(110) \quad \Rightarrow |t_1|^\nu \geq |t_2|^\nu \quad (\text{since } \lambda_1 > \lambda_2 \text{ and } \nu > 0) .$$

We now prove (**keep high SNR**). Consider  $x > 0$ . Since  $\lambda \mapsto s_{\sigma, \lambda}^\nu(x)$  is a monotonic function and  $s_{\sigma, \lambda}^\nu(x) \in [0, x]$  for all  $\lambda$ , it converges for  $\lambda \rightarrow \infty$  to a value  $\omega \in [0, x]$ . Assume  $0 < \omega < x$  and let  $0 < \varepsilon < \max(\omega, x - \omega)$ . By definition of the limit, for  $\lambda$  big enough

$$(111) \quad 0 < \omega - \varepsilon < t \triangleq s_{\sigma, \lambda}^\nu(x) < \omega + \varepsilon .$$

It follows that  $x - t > x - (\omega + \varepsilon) > 0$ , and then

$$(112) \quad \frac{(x - (\omega + \varepsilon))^2}{2\sigma^2} + \lambda_\nu^{-\nu} |\omega - \varepsilon|^\nu < \frac{(x - t)^2}{2\sigma^2} + \lambda_\nu^{-\nu} |t|^\nu .$$

Moreover, since  $\omega + \varepsilon \neq x$ , we have for  $\lambda$  big enough

$$(113) \quad \lambda_\nu^{-\nu} |x|^\nu < \frac{(x - (\omega + \varepsilon))^2}{2\sigma^2} + \lambda_\nu^{-\nu} |\omega - \varepsilon|^\nu .$$

Combining the two last inequalities shows that

$$(114) \quad \frac{(x - x)^2}{2\sigma^2} + \lambda_\nu^{-\nu} |x|^\nu < \frac{(x - t)^2}{2\sigma^2} + \lambda_\nu^{-\nu} |t|^\nu ,$$

which is in contradiction with the fact that  $t$  minimizes the objective. As a consequence,  $\omega = x$ , which concludes the proof since  $s_{\sigma, \lambda}^\nu(x)$  is odd and satisfies (**reduction**).

We now prove (**kill low SNR**). Consider  $x > 0$ . Since  $\lambda \mapsto s_{\sigma, 1/\lambda}^\nu(x)$  is a monotonic function and  $s_{\sigma, \lambda}^\nu(x) \in [0, x]$  for all  $\lambda$ , it converges for  $\lambda \rightarrow 0^+$  to a value  $\omega \in [0, x]$ . Assume  $0 < \omega < x$  and let  $0 < \varepsilon < \max(\omega, x - \omega)$ . Again, we have for  $\lambda$  small enough

$$(115) \quad \frac{(x - (\omega + \varepsilon))^2}{2\sigma^2} + \lambda_\nu^{-\nu} |\omega - \varepsilon|^\nu < \frac{(x - t)^2}{2\sigma^2} + \lambda_\nu^{-\nu} |t|^\nu .$$

Moreover, since  $\omega \neq \varepsilon$ , we have for  $\lambda$  small enough

$$(116) \quad \frac{x^2}{2\sigma^2} < \frac{(x - (\omega + \varepsilon))^2}{2\sigma^2} + \lambda_\nu^{-\nu} |\omega - \varepsilon|^\nu .$$

Combining the two last inequalities shows that

$$(117) \quad \frac{(x - 0)^2}{2\sigma^2} + \lambda_\nu^{-\nu} |0|^\nu < \frac{(x - t)^2}{2\sigma^2} + \lambda_\nu^{-\nu} |t|^\nu ,$$

which is in contradiction with the fact that  $t$  minimizes the objective. As a consequence,  $\omega = 0$ , which concludes the proof since  $s_{\sigma,\lambda}^\nu(x)$  is odd and satisfies (reduction). ■

## REFERENCES

- [1] M. AHARON, M. ELAD, AND A. BRUCKSTEIN, *k-SVD: An algorithm for designing overcomplete dictionaries for sparse representation*, IEEE Transactions on signal processing, 54 (2006), pp. 4311–4322.
- [2] S. M. BERMAN, *The tail of the convolution of densities and its application to a model of HIV-latency time*, The Annals of Applied Probability, 2 (1992), pp. 481–502.
- [3] L. BOUBCHIR AND J. FADILI, *Multivariate statistical modeling of images with the curvelet transform*, in ISSPA, 2005, pp. 747–750.
- [4] K. BREDIES, K. KUNISCH, AND T. POCK, *Total generalized variation*, SIAM Journal on Imaging Sciences, 3 (2010), pp. 492–526.
- [5] A. BUADES, B. COLL, AND J.-M. MOREL, *A non-local algorithm for image denoising*, in Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on, vol. 2, IEEE, 2005, pp. 60–65.
- [6] S. G. CHANG, B. YU, AND M. VETTERLI, *Adaptive wavelet thresholding for image denoising and compression*, IEEE transactions on image processing, 9 (2000), pp. 1532–1546.
- [7] C. CHAUX, P. L. COMBETTES, J.-C. PESQUET, AND V. R. WAJS, *A variational formulation for frame-based inverse problems*, Inverse Problems, 23 (2007), p. 1495.
- [8] R. R. COIFMAN AND D. L. DONOHO, *Translation-invariant de-noising*, Wavelets and statistics, 103 (1995), pp. 125–150.
- [9] K. DABOV, A. FOI, V. KATKOVNIK, AND K. EGIAZARIAN, *Image denoising by sparse 3-d transform-domain collaborative filtering*, IEEE Transactions on image processing, 16 (2007), pp. 2080–2095.
- [10] G. DAVIS, S. MALLAT, AND Z. ZHANG, *Adaptive time-frequency decompositions with matching pursuit*, Optical Engineering, 33 (1994).
- [11] C.-A. DELEDALLE, J. SALMON, A. S. DALALYAN, ET AL., *Image denoising with patch based pca: local versus global.*, in BMVC, vol. 81, 2011, pp. 425–455.
- [12] C.-A. DELEDALLE, S. VAITER, J. FADILI, AND G. PEYRÉ, *Stein unbiased gradient estimator of the risk (sugar) for multiple parameter selection*, SIAM Journal on Imaging Sciences, 7 (2014), pp. 2448–2487.
- [13] 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 (methodological), (1977), pp. 1–38.
- [14] M. N. DO AND M. VETTERLI, *Wavelet-based texture retrieval using generalized gaussian density and kullback-leibler distance*, IEEE transactions on image processing, 11 (2002), pp. 146–158.
- [15] D. L. DONOHO AND I. M. JOHNSTONE, *Adapting to unknown smoothness via wavelet shrinkage*, Journal of the american statistical association, 90 (1995), pp. 1200–1224.
- [16] D. L. DONOHO AND J. M. JOHNSTONE, *Ideal spatial adaptation by wavelet shrinkage*, biometrika, 81 (1994), pp. 425–455.
- [17] R. J. DUFFIN AND A. C. SCHAEFFER, *A class of nonharmonic fourier series*, Transactions of the American Mathematical Society, 72 (1952), pp. 341–366.
- [18] C. DUGAS, Y. BENGIO, F. BÉLISLE, C. NADEAU, AND R. GARCIA, *Incorporating second-order functional knowledge for better option pricing*, in Advances in neural information processing systems, 2001, pp. 472–478.

- [19] M. ELAD AND M. AHARON, *Image denoising via sparse and redundant representations over learned dictionaries*, IEEE Transactions on Image processing, 15 (2006), pp. 3736–3745.
- [20] M. ELAD, P. MILANFAR, AND R. RUBINSTEIN, *Analysis versus synthesis in signal priors*, Inverse problems, 23 (2007), p. 947.
- [21] R. FERGUS, B. SINGH, A. HERTZMANN, S. T. ROWEIS, AND W. T. FREEMAN, *Removing camera shake from a single photograph*, in ACM transactions on graphics (TOG), vol. 25, ACM, 2006, pp. 787–794.
- [22] D. GEMAN AND C. YANG, *Nonlinear image recovery with half-quadratic regularization*, IEEE Transactions on Image Processing, 4 (1995), pp. 932–946.
- [23] S. GEMAN AND D. GEMAN, *Stochastic relaxation, gibbs distributions, and the bayesian restoration of images*, IEEE Transactions on pattern analysis and machine intelligence, (1984), pp. 721–741.
- [24] G. H. GOLUB, M. HEATH, AND G. WAHBA, *Generalized cross-validation as a method for choosing a good ridge parameter*, Technometrics, 21 (1979), pp. 215–223.
- [25] P. C. HANSEN, *The L-curve and its use in the numerical treatment of inverse problems*, IMM, Department of Mathematical Modelling, Technical University of Denmark, 1999.
- [26] A. HOUDARD, C. BOUVEYRON, AND J. DELON, *High-Dimensional Mixture Models For Unsupervised Image Denoising (HDMI)*. Preprint hal-01544249, Aug. 2017, <https://hal.archives-ouvertes.fr/hal-01544249>.
- [27] P. J. HUBER, *Robust statistics*, in International Encyclopedia of Statistical Science, Springer, 2011, pp. 1248–1251.
- [28] P. J. HUBER ET AL., *Robust estimation of a location parameter*, The Annals of Mathematical Statistics, 35 (1964), pp. 73–101.
- [29] H. JI, S. HUANG, Z. SHEN, AND Y. XU, *Robust video restoration by joint sparse and low rank matrix approximation*, SIAM Journal on Imaging Sciences, 4 (2011), pp. 1122–1142.
- [30] D. KRISHNAN AND R. FERGUS, *Fast image deconvolution using hyper-laplacian priors*, in Advances in Neural Information Processing Systems, 2009, pp. 1033–1041.
- [31] R. KRUPIŃSKI, *Approximated fast estimator for the shape parameter of generalized gaussian distribution for a small sample size*, Bulletin of the Polish Academy of Sciences Technical Sciences, 63 (2015), pp. 405–411.
- [32] M. LEBRUN, A. BUADES, AND J.-M. MOREL, *A nonlocal bayesian image denoising algorithm*, SIAM Journal on Imaging Sciences, 6 (2013), pp. 1665–1688.
- [33] S. MALLAT, *A theory for multiresolution signal decomposition: the wavelet representation*, IEEE transactions on pattern analysis and machine intelligence, 11 (1989), pp. 674–693.
- [34] S. MALLAT, *A wavelet tour of signal processing: the sparse way*, Academic press, 2008.
- [35] D. MARTIN, C. FOWLKES, D. TAL, AND J. MALIK, *A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics*, in Computer Vision, 2001. ICCV 2001. Proceedings. Eighth IEEE International Conference on, vol. 2, IEEE, 2001, pp. 416–423.
- [36] T. MICHAELI AND M. IRANI, *Blind deblurring using internal patch recurrence*, in European Conference on Computer Vision, Springer, 2014, pp. 783–798.
- [37] P. MITTAL AND K. GUPTA, *An integral involving generalized function of two variables*, in Proceedings of the Indian academy of sciences-section A, vol. 75, Springer, 1972, pp. 117–123.
- [38] O. M. M. MOHAMED AND M. JAÏDANE-SAÏDANE, *Generalized gaussian mixture model*, in Signal Processing Conference, 2009 17th European, IEEE, 2009, pp. 2273–2277.
- [39] P. MOULIN AND J. LIU, *Analysis of multiresolution image denoising schemes using generalized gaussian and complexity priors*, IEEE Transactions on Information Theory, 45 (1999), pp. 909–919.
- [40] V. NAIR AND G. E. HINTON, *Rectified linear units improve restricted boltzmann machines*, in Proceedings of the 27th international conference on machine learning (ICML-10), 2010, pp. 807–814.
- [41] M. NARDON AND P. PIANCA, *Simulation techniques for generalized gaussian densities*, Journal of Statistical Computation and Simulation, 79 (2009), pp. 1317–1329.
- [42] M. NIKNEJAD, J. M. BIOCAS-DIAS, AND M. A. FIGUEIREDO, *Class-specific image denoising using importance sampling*, in IEEE International Conference on Image Processing, IEEE, 2017.
- [43] V. PAPPYAN AND M. ELAD, *Multi-scale patch-based image restoration*, IEEE Transactions on image processing, 25 (2016), pp. 249–261.
- [44] S. PARAMESWARAN, C.-A. DELEDALLE, L. DENIS, AND T. Q. NGUYEN, *Accelerating gmm-based patch*

- priors for image restoration: Three ingredients for a  $100\times$  speed-up, arXiv preprint arXiv:1710.08124, (2017).
- [45] F. PASCAL, L. BOMBRUN, J.-Y. TOURNERET, AND Y. BERTHOUMIEU, *Parameter estimation for multivariate generalized gaussian distributions*, IEEE Transactions on Signal Processing, 61 (2013), pp. 5960–5971.
  - [46] K. P. PEPPAS, *A new formula for the average bit error probability of dual-hop amplify-and-forward relaying systems over generalized shadowed fading channels*, IEEE Wireless Communications Letters, 1 (2012), pp. 85–88.
  - [47] J. PORTILLA, V. STRELA, M. J. WAINWRIGHT, AND E. P. SIMONCELLI, *Image denoising using scale mixtures of gaussians in the wavelet domain*, IEEE Transactions on Image processing, 12 (2003), pp. 1338–1351.
  - [48] S. PURKAYASTHA, *Simple proofs of two results on convolutions of unimodal distributions*, Statistics & probability letters, 39 (1998), pp. 97–100.
  - [49] S. RAMANI, T. BLU, AND M. UNSER, *Monte-carlo sure: A black-box optimization of regularization parameters for general denoising algorithms*, IEEE Transactions on Image Processing, 17 (2008), pp. 1540–1554.
  - [50] A. A. ROENKO, V. V. LUKIN, I. DJUROVIC, AND M. SIMEUNOVIC, *Estimation of parameters for generalized gaussian distribution*, in Communications, Control and Signal Processing (ISCCSP), 2014 6th International Symposium on, IEEE, 2014, pp. 376–379.
  - [51] S. ROTH AND M. J. BLACK, *Fields of experts: A framework for learning image priors*, in Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on, vol. 2, IEEE, 2005, pp. 860–867.
  - [52] R. RUBINSTEIN, T. PELEG, AND M. ELAD, *Analysis K-SVD: A dictionary-learning algorithm for the analysis sparse model*, IEEE Transactions on Signal Processing, 61 (2013), pp. 661–677.
  - [53] L. I. RUDIN, S. OSHER, AND E. FATEMI, *Nonlinear total variation based noise removal algorithms*, Physica D: Nonlinear Phenomena, 60 (1992), pp. 259–268.
  - [54] M. J. SHENSA, *The discrete wavelet transform: wedding the a trous and mallat algorithms*, IEEE Transactions on signal processing, 40 (1992), pp. 2464–2482.
  - [55] G. SINGH AND J. SINGH, *Adapted non-gaussian mixture model for image denoising*, Journal of Basic and Applied Engineering Research, 4 (2017).
  - [56] H. SOURY AND M.-S. ALOUINI, *New results on the sum of two generalized gaussian random variables*, in Signal and Information Processing (GlobalSIP), 2015 IEEE Global Conference on, IEEE, 2015, pp. 1017–1021.
  - [57] J. SULAM AND M. ELAD, *Expected patch log likelihood with a sparse prior*, Energy Minimization Methods in Computer Vision and Pattern Recognition, Lecture Notes in Computer Science, (2015), pp. 99–111.
  - [58] Z. WANG, A. C. BOVIK, H. R. SHEIKH, AND E. P. SIMONCELLI, *Image quality assessment: from error visibility to structural similarity*, IEEE transactions on image processing, 13 (2004), pp. 600–612.
  - [59] A. WINTNER, *Asymptotic distributions and infinite convolutions. edwards brothers*, Ann Arbor, MI. In Bickel, PJ, (1938).
  - [60] J. YANG, Y. ZHANG, AND W. YIN, *An efficient tvl1 algorithm for deblurring multichannel images corrupted by impulsive noise*, SIAM Journal on Scientific Computing, 31 (2009), pp. 2842–2865.
  - [61] G. YU, G. SAPIRO, AND S. MALLAT, *Solving inverse problems with piecewise linear estimators: From gaussian mixture models to structured sparsity*, IEEE Transactions on Image Processing, 21 (2012), pp. 2481–2499.
  - [62] D. ZORAN AND Y. WEISS, *From learning models of natural image patches to whole image restoration*, in International Conference on Computer Vision, IEEE, November 2011, pp. 479–486, <https://doi.org/10.1109/ICCV.2011.6126278>.