

Back-Projection based Fidelity Term for Ill-Posed Linear Inverse Problems

Tom Tirer and Raja Giryes

Abstract—Ill-posed linear inverse problems appear in many image processing applications, such as deblurring, super-resolution and compressed sensing. Many restoration strategies involve minimizing a cost function, which is composed of fidelity and prior terms, balanced by a regularization parameter. While a vast amount of research has been focused on different prior models, the fidelity term is almost always chosen to be the least squares (LS) objective, that encourages fitting the linearly transformed optimization variable to the observations. In this paper, we examine a different fidelity term, which has been implicitly used by the recently proposed iterative denoising and backward projections (IDBP) framework. This term encourages agreement between the projection of the optimization variable onto the row space of the linear operator and the pseudo-inverse of the linear operator (“back-projection”) applied on the observations. We analytically examine the difference between the two fidelity terms for Tikhonov regularization and identify cases (such as a badly conditioned linear operator) where the new term has an advantage over the standard LS one. Moreover, we demonstrate empirically that the behavior of the two induced cost functions for sophisticated convex and non-convex priors, such as total-variation, BM3D, and deep generative models, correlates with the obtained theoretical analysis.

Index Terms—Inverse problems, image restoration, image deblurring, image super-resolution, compressed sensing, total variation, non-convex priors, BM3D, deep generative models.

I. INTRODUCTION

Inverse problems appear in many fields of science and engineering, where the goal is to recover a signal from its observations that are obtained by some acquisition process. In image processing, the observations are usually a degraded version of the latent image, which may be noisy, blurred, downsampled, or all together. Such observation models, and others, can be formulated by a linear model

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{e}, \quad (1)$$

where $\mathbf{x} \in \mathbb{R}^n$ represents the unknown original image, $\mathbf{y} \in \mathbb{R}^m$ represents the observations, \mathbf{A} is an $m \times n$ degradation matrix (sometimes also referred to as the measurement matrix) and $\mathbf{e} \in \mathbb{R}^m$ is a noise vector. For example, this model corresponds to the problem of denoising [1]–[4] when \mathbf{A} is the $n \times n$ identity matrix \mathbf{I}_n ; inpainting [5]–[7] when \mathbf{A} is an $m \times n$ sampling matrix (i.e. a selection of m rows of \mathbf{I}_n); deblurring [8], [9] when \mathbf{A} is a blur operator; super-resolution [10], [11] if \mathbf{A} is a composite operator of blurring (e.g. anti-aliasing filtering) and down-sampling; and compressed sensing

when \mathbf{A} is a (random) measurement matrix ($m \ll n$) and the signal is sparse under some basis representation [12]–[14] or resides in a general union of low-dimensional subspaces [15], [16].

The inverse problems represented by (1) are usually ill-posed, i.e. the measurements do not suffice for obtaining a successful reconstruction. Therefore, a vast amount of research has focused on designing good prior models for natural images. In fact, many of the methods for the problems mentioned above differ only in their prior assumptions and not in the way that they enforce fidelity to the observations.

To be more formal, a common strategy for recovering \mathbf{x} aims at minimizing a cost function of the form

$$f(\tilde{\mathbf{x}}) = \ell(\tilde{\mathbf{x}}) + \beta s(\tilde{\mathbf{x}}), \quad (2)$$

where $\ell(\tilde{\mathbf{x}})$ is a fidelity term, $s(\tilde{\mathbf{x}})$ is a prior term (can be also referred to as the regularizer), β is a positive scalar that controls the level of regularization, and $\tilde{\mathbf{x}}$ is the optimization variable. Many different prior functions are used in the literature, whether explicitly, e.g. total-variation (TV) [1], or implicitly, e.g. BM3D [4] and deep generative models [17]. Yet, most of the works use a typical least squares (LS) fidelity term

$$\ell_{LS}(\tilde{\mathbf{x}}) \triangleq \frac{1}{2} \|\mathbf{y} - \mathbf{A}\tilde{\mathbf{x}}\|_2^2, \quad (3)$$

where $\|\cdot\|_2$ stands for the Euclidean norm. The frequent usage of this term is probably also motivated by the fact that it can be derived from the negative log-likelihood function, under the assumption that the noise \mathbf{e} is a vector of i.i.d. Gaussian random variables $e_i \sim \mathcal{N}(0, \sigma_e^2)$. However, note that, in general, maximum likelihood estimation has optimality properties only when the number of measurements is *much larger* than the number of unknown variables, which is obviously *not the case* in ill-posed problems.

In this paper, we examine a different fidelity term, which has been implicitly used by the recently proposed iterative denoising and backward projections (IDBP) framework [18] (we elaborate on this method in the appendix). Under the practical assumptions that $m \leq n$ and $\text{rank}(\mathbf{A}) = m$, we examine the fidelity term

$$\ell_{BP}(\tilde{\mathbf{x}}) \triangleq \frac{1}{2} \|\mathbf{A}^\dagger \mathbf{y} - \mathbf{A}^\dagger \mathbf{A}\tilde{\mathbf{x}}\|_2^2, \quad (4)$$

where $\mathbf{A}^\dagger \triangleq \mathbf{A}^T(\mathbf{A}\mathbf{A}^T)^{-1}$ is the pseudoinverse of the full row-rank matrix \mathbf{A} . Note that $\mathbf{P}_A \triangleq \mathbf{A}^\dagger \mathbf{A}$ is an orthogonal projection onto the row space of \mathbf{A}^\dagger , and that \mathbf{A}^\dagger can be

¹In row space of \mathbf{A} , we mean the subspace spanned by the rows of \mathbf{A} .

interpreted as a "back-projection" (BP) from $\mathbf{A}\mathbb{R}^n$ back to \mathbb{R}^n . Therefore, the fidelity (4) encourages agreement between $\mathbf{P}_A\tilde{\mathbf{x}}$ —the projection of $\tilde{\mathbf{x}}$ onto the row space of \mathbf{A} , and $\mathbf{A}^\dagger\mathbf{y}$ —the back-projection of the measurements. In general, this is different than $\ell_{LS}(\tilde{\mathbf{x}})$ that encourages agreement between $\mathbf{A}\tilde{\mathbf{x}}$ and \mathbf{y} . Note that in the noiseless case, i.e. when $\mathbf{y} = \mathbf{A}\mathbf{x}$, the terms in (3) and (4) are translated to fitting $\mathbf{A}\tilde{\mathbf{x}}$ to $\mathbf{A}\mathbf{x}$ and $\mathbf{P}_A\tilde{\mathbf{x}}$ to $\mathbf{P}_A\mathbf{x}$, respectively.

Note that for some inverse problems $\ell_{LS}(\tilde{\mathbf{x}})$ and $\ell_{BP}(\tilde{\mathbf{x}})$ may coincide. For example, in image inpainting, where \mathbf{A} is a selection of m rows of \mathbf{I}_n , we have that $\mathbf{A}^\dagger = \mathbf{A}^T$ is an $n \times m$ matrix that merely pads with $n - m$ zeros the vector on which it is applied, and so $\|\mathbf{A}^\dagger(\mathbf{y} - \mathbf{A}\tilde{\mathbf{x}})\|_2^2 = \|\mathbf{y} - \mathbf{A}\tilde{\mathbf{x}}\|_2^2$. Therefore, we specifically focus on three popular inverse problems: super-resolution, deblurring and certain compressed sensing scenarios, where the two fidelity terms, $\ell_{LS}(\tilde{\mathbf{x}})$ and $\ell_{BP}(\tilde{\mathbf{x}})$, are indeed very different.

Contribution. This work makes a first attempt towards characterizing for which observation model \mathbf{A} and prior $s(\tilde{\mathbf{x}})$ it is better to use each of the following objectives:

$$f_{LS}(\tilde{\mathbf{x}}) \triangleq \frac{1}{2}\|\mathbf{y} - \mathbf{A}\tilde{\mathbf{x}}\|_2^2 + \beta s(\tilde{\mathbf{x}}), \quad (5)$$

$$f_{BP}(\tilde{\mathbf{x}}) \triangleq \frac{1}{2}\|\mathbf{A}^\dagger\mathbf{y} - \mathbf{A}^\dagger\mathbf{A}\tilde{\mathbf{x}}\|_2^2 + \beta s(\tilde{\mathbf{x}}). \quad (6)$$

Particularly, for $s(\tilde{\mathbf{x}})$ being the Tikhonov regularization (the ℓ_2 prior), where closed-form solutions exist, we derive analytical expressions for the estimations' mean square error (MSE) that allow to examine which fidelity term is preferable. For example, we show that in the noiseless case $f_{BP}(\tilde{\mathbf{x}})$ yields provably better restoration than $f_{LS}(\tilde{\mathbf{x}})$ if the condition number of $\mathbf{A}\mathbf{A}^T$ (i.e. the ratio between the largest and smallest squared singular values of \mathbf{A}) is large, e.g. in typical super-resolution problems.

For sophisticated convex and non-convex priors, such as TV [1], BM3D [4], and DCGAN [19], analytical analysis is intractable. Therefore, we perform an intensive empirical study, where we use the same optimization method (FISTA [20] or ADAM [21]) to minimize each of the two different cost functions. Interestingly, we demonstrate that the behavior for the sophisticated priors strongly correlates with properties for which we establish concrete mathematical reasoning in the case of ℓ_2 priors.

Another contribution of the paper that is deferred to the appendix is showing that IDBP framework [18], which has achieved excellent results for deblurring [18], [22] and super-resolution [23] is in fact the proximal gradient method [20], [24] (popularized under the name ISTA) applied on $f_{BP}(\tilde{\mathbf{x}})$. This derivation of IDBP is completely different, and arguably simpler, than the way it is developed in [18].

The paper is organized as follows. Section II includes mathematical analysis of the two cost functions for the case of ℓ_2 -type priors. The analytical results are verified in Section III. In Section IV the two cost functions are empirically examined for different sophisticated priors. Section V concludes the paper.

II. MATHEMATICAL ANALYSIS FOR ℓ_2 PRIORS

In this section, we analyze the performance of the new cost function (6) and compare it to (5) for a type of ℓ_2 prior functions, for which the closed-form solutions of (5) and (6) lead to a tractable performance analysis. We start with specifying the required assumptions, then we derive the estimators and expressions for their expected mean square error. Finally, the error expressions are compared and several observations are stated.

A. Assumptions

In order to allow a concrete mathematical comparison between $f_{BP}(\tilde{\mathbf{x}})$ and $f_{LS}(\tilde{\mathbf{x}})$, in the theoretical analysis we restrict our discussion to ℓ_2 prior functions of the form $s(\tilde{\mathbf{x}}) = \frac{1}{2}\|\mathbf{D}\tilde{\mathbf{x}}\|_2^2 = \frac{1}{2}\tilde{\mathbf{x}}^T\mathbf{D}^T\mathbf{D}\tilde{\mathbf{x}}$, where $\mathbf{D}^T\mathbf{D}$ is a positive-definite matrix. This prior is often referred to as Tikhonov regularization and is one of the most widely used methods to solve ill-posed inverse problems. Yet, for obtaining analytical results, we further focus on a more specific type of this prior—we require that both \mathbf{A} and \mathbf{D} have the same right singular vectors. Let us define the singular value decomposition (SVD) of the $m \times n$ matrix $\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{V}^T$, where \mathbf{U} is an $m \times m$ orthogonal matrix whose columns are the left singular vectors, $\mathbf{\Lambda}$ is an $m \times n$ rectangular diagonal matrix with nonzero singular values $\{\lambda_i\}_{i=1}^m$ on the diagonal, and \mathbf{V} is an $n \times n$ orthogonal matrix whose columns are the right singular vectors. The property that $\{\lambda_i\}_{i=1}^m$ are strictly positive follows from our assumptions in Section I, that $m \leq n$ and $\text{rank}(\mathbf{A}) = m$. For \mathbf{D} , essentially, we assume that $\mathbf{D}^T\mathbf{D} = \mathbf{V}\mathbf{\Gamma}^2\mathbf{V}^T \succ 0$, where $\mathbf{\Gamma}^2$ is an $n \times n$ diagonal matrix of nonzero eigenvalues $\{\gamma_i^2\}_{i=1}^n$.

The assumption above is required because, as far as we know, currently there is no known analytical expression for the eigen-decomposition of arbitrary matrices $\mathbf{A}^T\mathbf{A} + \mathbf{D}^T\mathbf{D}$ which is required for our analysis [25]. Yet, this assumption holds in some practical cases, e.g. if \mathbf{A} and \mathbf{D} are circulant matrices (and thus diagonalized by the discrete Fourier transform), or if $\mathbf{D} = \mathbf{I}_n$ (i.e. least-norm regularization).

B. Performance analysis

Let us start with obtaining closed-form expressions for the estimators $\hat{\mathbf{x}}_{LS}$ and $\hat{\mathbf{x}}_{BP}$, which minimize $f_{LS}(\tilde{\mathbf{x}})$ and $f_{BP}(\tilde{\mathbf{x}})$, respectively. Due to the convexity of the cost functions, this is done simply by equating their gradients to zero

$$\begin{aligned} \nabla f_{LS}(\tilde{\mathbf{x}}) &= -\mathbf{A}^T(\mathbf{y} - \mathbf{A}\tilde{\mathbf{x}}) + \beta\mathbf{D}^T\mathbf{D}\tilde{\mathbf{x}} = 0 \\ &\Rightarrow \hat{\mathbf{x}}_{LS} = (\mathbf{A}^T\mathbf{A} + \beta\mathbf{D}^T\mathbf{D})^{-1}\mathbf{A}^T\mathbf{y}, \end{aligned} \quad (7)$$

$$\begin{aligned} \nabla f_{BP}(\tilde{\mathbf{x}}) &= -\mathbf{P}_A(\mathbf{A}^\dagger\mathbf{y} - \mathbf{P}_A\tilde{\mathbf{x}}) + \beta\mathbf{D}^T\mathbf{D}\tilde{\mathbf{x}} = 0 \\ &\Rightarrow \hat{\mathbf{x}}_{BP} = (\mathbf{P}_A + \beta\mathbf{D}^T\mathbf{D})^{-1}\mathbf{A}^\dagger\mathbf{y}. \end{aligned} \quad (8)$$

In (8) we use the properties $\mathbf{P}_A \triangleq \mathbf{A}^\dagger\mathbf{A} = \mathbf{P}_A^T = \mathbf{P}_A^2$ and $\mathbf{P}_A\mathbf{A}^\dagger = \mathbf{A}^\dagger$. We turn to compute the expected mean square errors (MSEs) of the estimators, conditioned on \mathbf{x} , under the assumptions that $\mathbb{E}[e] = 0$ and $\mathbb{E}[ee^T] = \sigma_e^2\mathbf{I}_m$. To ease formulations, we define the $n - m$ zero eigenvalues of $\mathbf{A}^T\mathbf{A}$ (i.e. zeros in the diagonal of $\mathbf{\Lambda}^T\mathbf{\Lambda}$) by $\{\lambda_i^2\}_{i=m+1}^n$.

The computation of the MSE of $\hat{\mathbf{x}}_{LS}$ is given by

$$\begin{aligned}
MSE_{LS} &= \mathbb{E} \|\hat{\mathbf{x}}_{LS} - \mathbf{x}\|_2^2 \\
&= \mathbb{E} \|(\mathbf{A}^T \mathbf{A} + \beta \mathbf{D}^T \mathbf{D})^{-1} \mathbf{A}^T (\mathbf{A} \mathbf{x} + \mathbf{e}) - \mathbf{x}\|_2^2 \\
&= \|((\mathbf{A}^T \mathbf{A} + \beta \mathbf{D}^T \mathbf{D})^{-1} \mathbf{A}^T \mathbf{A} - \mathbf{I}_n) \mathbf{x}\|_2^2 \\
&\quad + 2\mathbb{E} [\mathbf{e}]^T \mathbf{A} (\mathbf{A}^T \mathbf{A} + \beta \mathbf{D}^T \mathbf{D})^{-2} \mathbf{A}^T \mathbf{A} \mathbf{x} \\
&\quad - 2\mathbb{E} [\mathbf{e}]^T \mathbf{A} (\mathbf{A}^T \mathbf{A} + \beta \mathbf{D}^T \mathbf{D})^{-1} \mathbf{x} \\
&\quad + \mathbb{E} [\mathbf{e}^T \mathbf{A} (\mathbf{A}^T \mathbf{A} + \beta \mathbf{D}^T \mathbf{D})^{-2} \mathbf{A}^T \mathbf{e}] \\
&= \|((\mathbf{A}^T \mathbf{A} + \beta \mathbf{D}^T \mathbf{D})^{-1} \mathbf{A}^T \mathbf{A} - \mathbf{I}_n) \mathbf{x}\|_2^2 \\
&\quad + \text{Tr}((\mathbf{A}^T \mathbf{A} + \beta \mathbf{D}^T \mathbf{D})^{-2} \mathbf{A}^T \mathbb{E} [\mathbf{e} \mathbf{e}^T] \mathbf{A}) \\
&= \|((\mathbf{A}^T \mathbf{A} + \beta \mathbf{D}^T \mathbf{D})^{-1} \mathbf{A}^T \mathbf{A} - \mathbf{I}_n) \mathbf{x}\|_2^2 \\
&\quad + \sigma_e^2 \text{Tr}((\mathbf{A}^T \mathbf{A} + \beta \mathbf{D}^T \mathbf{D})^{-2} \mathbf{A}^T \mathbf{A}) \\
&= \|\mathbf{V} ((\mathbf{\Lambda}^T \mathbf{\Lambda} + \beta \mathbf{\Gamma}^2)^{-1} \mathbf{\Lambda}^T \mathbf{\Lambda} - \mathbf{I}_n) \mathbf{V}^T \mathbf{x}\|_2^2 \\
&\quad + \sigma_e^2 \text{Tr}(\mathbf{V} (\mathbf{\Lambda}^T \mathbf{\Lambda} + \beta \mathbf{\Gamma}^2)^{-2} \mathbf{\Lambda}^T \mathbf{\Lambda} \mathbf{V}^T) \\
&= \sum_{i=1}^n \left(\frac{\lambda_i^2}{\lambda_i^2 + \beta \gamma_i^2} - 1 \right)^2 [\mathbf{V}^T \mathbf{x}]_i^2 + \sigma_e^2 \sum_{i=1}^n \frac{\lambda_i^2}{(\lambda_i^2 + \beta \gamma_i^2)^2}. \tag{9}
\end{aligned}$$

The second equality follows from substituting (1) in (7), the fourth equality uses $\mathbb{E}[\mathbf{e}] = \mathbf{0}$ and the cyclic property of trace, the fifth equality uses $\mathbb{E}[\mathbf{e} \mathbf{e}^T] = \sigma_e^2 \mathbf{I}_m$, the sixth equality is obtained by substituting the eigen-decompositions of $\mathbf{A}^T \mathbf{A}$ and $\mathbf{D}^T \mathbf{D}$, and the last equality follows from the fact that \mathbf{V} is an orthogonal matrix. Therefore, by defining the (squared) bias and variance terms as

$$\begin{aligned}
bias_{LS}^2 &\triangleq \sum_{i=1}^m \underbrace{\left(\frac{\beta \gamma_i^2}{\lambda_i^2 + \beta \gamma_i^2} \right)^2}_{\triangleq bias_{LS}^{2(i)}} [\mathbf{V}^T \mathbf{x}]_i^2 + \sum_{i=m+1}^n [\mathbf{V}^T \mathbf{x}]_i^2, \\
var_{LS} &\triangleq \sum_{i=1}^m \underbrace{\frac{\sigma_e^2}{\lambda_i^2 (1 + \beta \gamma_i^2 / \lambda_i^2)}}_{\triangleq var_{LS}^{(i)}}, \tag{10}
\end{aligned}$$

we may write the error as

$$MSE_{LS} = bias_{LS}^2 + var_{LS}. \tag{11}$$

Note that the bias depends on the original image \mathbf{x} and not on the noise, and the opposite holds for the variance. Yet, both terms are affected by the structure of \mathbf{A} . The regularization parameters $\beta, \{\gamma_i\}$ introduce a tradeoff: increasing them reduces the variance but increases the bias.

To ease the computation of the MSE of $\hat{\mathbf{x}}_{BP}$, let us also define an indicator function $1_{i \leq m}$ that is equal to 1 if $i \leq m$ and 0 otherwise, and an $n \times n$ diagonal matrix $\mathbf{I}_{i \leq m}$ with $\{1_{i \leq m}\}_{i=1}^n$ on its diagonal. The following identities are used

$$\begin{aligned}
\mathbf{P}_A &= \mathbf{V} \mathbf{I}_{i \leq m} \mathbf{V}^T, \\
\mathbf{A}^\dagger &= \mathbf{V} \mathbf{\Lambda}^T (\mathbf{\Lambda} \mathbf{\Lambda}^T)^{-1} \mathbf{U}^T, \\
\mathbf{A}^\dagger \mathbf{A}^{\dagger T} &= \mathbf{V} \mathbf{\Lambda}^T (\mathbf{\Lambda} \mathbf{\Lambda}^T)^{-2} \mathbf{\Lambda} \mathbf{V}^T. \tag{12}
\end{aligned}$$

Now, we get

$$\begin{aligned}
MSE_{BP} &= \mathbb{E} \|\hat{\mathbf{x}}_{BP} - \mathbf{x}\|_2^2 \\
&= \mathbb{E} \|(\mathbf{P}_A + \beta \mathbf{D}^T \mathbf{D})^{-1} \mathbf{A}^\dagger (\mathbf{A} \mathbf{x} + \mathbf{e}) - \mathbf{x}\|_2^2 \\
&= \|((\mathbf{P}_A + \beta \mathbf{D}^T \mathbf{D})^{-1} \mathbf{P}_A - \mathbf{I}_n) \mathbf{x}\|_2^2 \\
&\quad + 2\mathbb{E} [\mathbf{e}]^T \mathbf{A}^{\dagger T} (\mathbf{P}_A + \beta \mathbf{D}^T \mathbf{D})^{-2} \mathbf{P}_A \mathbf{x} \\
&\quad - 2\mathbb{E} [\mathbf{e}]^T \mathbf{A}^{\dagger T} (\mathbf{P}_A + \beta \mathbf{D}^T \mathbf{D})^{-1} \mathbf{x} \\
&\quad + \mathbb{E} [\mathbf{e}^T \mathbf{A}^{\dagger T} (\mathbf{P}_A + \beta \mathbf{D}^T \mathbf{D})^{-2} \mathbf{A}^\dagger \mathbf{e}] \\
&= \|((\mathbf{P}_A + \beta \mathbf{D}^T \mathbf{D})^{-1} \mathbf{P}_A - \mathbf{I}_n) \mathbf{x}\|_2^2 \\
&\quad + \text{Tr}((\mathbf{P}_A + \beta \mathbf{D}^T \mathbf{D})^{-2} \mathbf{A}^\dagger \mathbb{E} [\mathbf{e} \mathbf{e}^T] \mathbf{A}^{\dagger T}) \\
&= \|((\mathbf{P}_A + \beta \mathbf{D}^T \mathbf{D})^{-1} \mathbf{P}_A - \mathbf{I}_n) \mathbf{x}\|_2^2 \\
&\quad + \sigma_e^2 \text{Tr}((\mathbf{P}_A + \beta \mathbf{D}^T \mathbf{D})^{-2} \mathbf{A}^\dagger \mathbf{A}^{\dagger T}) \\
&= \|\mathbf{V} ((\mathbf{I}_{i \leq m} + \beta \mathbf{\Gamma}^2)^{-1} \mathbf{I}_{i \leq m} - \mathbf{I}_n) \mathbf{V}^T \mathbf{x}\|_2^2 \\
&\quad + \sigma_e^2 \text{Tr}(\mathbf{V} (\mathbf{I}_{i \leq m} + \beta \mathbf{\Gamma}^2)^{-2} \mathbf{\Lambda}^T (\mathbf{\Lambda} \mathbf{\Lambda}^T)^{-2} \mathbf{\Lambda} \mathbf{V}^T) \\
&= \sum_{i=1}^n \left(\frac{1_{i \leq m}}{1_{i \leq m} + \beta \gamma_i^2} - 1 \right)^2 [\mathbf{V}^T \mathbf{x}]_i^2 + \sigma_e^2 \sum_{i=1}^n \frac{\lambda_i^{-2} 1_{i \leq m}}{(1_{i \leq m} + \beta \gamma_i^2)^2}. \tag{13}
\end{aligned}$$

The second equality follows from substituting (1) in (8), the fourth equality uses $\mathbb{E}[\mathbf{e}] = \mathbf{0}$ and the cyclic property of trace, the fifth equality uses $\mathbb{E}[\mathbf{e} \mathbf{e}^T] = \sigma_e^2 \mathbf{I}_m$, the sixth equality is obtained by substituting the eigen-decompositions of \mathbf{P}_A , $\mathbf{D}^T \mathbf{D}$ and $\mathbf{A}^\dagger \mathbf{A}^{\dagger T}$, and the last equality uses the orthogonality of \mathbf{V} . Therefore, by defining

$$\begin{aligned}
bias_{BP}^2 &\triangleq \sum_{i=1}^m \underbrace{\left(\frac{\beta \gamma_i^2}{1 + \beta \gamma_i^2} \right)^2}_{\triangleq bias_{BP}^{2(i)}} [\mathbf{V}^T \mathbf{x}]_i^2 + \sum_{i=m+1}^n [\mathbf{V}^T \mathbf{x}]_i^2, \\
var_{BP} &\triangleq \sum_{i=1}^m \underbrace{\frac{\sigma_e^2}{\lambda_i^2 (1 + \beta \gamma_i^2)^2}}_{\triangleq var_{BP}^{(i)}}, \tag{14}
\end{aligned}$$

we have that

$$MSE_{BP} = bias_{BP}^2 + var_{BP}. \tag{15}$$

Comparing (10) and (14) we may notice the following. First, the term $bias_{BP}^2$ handles small $\{\lambda_i\}_{i=1}^m$ (i.e. singular values of \mathbf{A} that are smaller than 1) better than $bias_{LS}^2$. However, var_{BP} handles such small singular values worse than var_{LS} . The opposite holds for singular values that are greater than 1. This behavior can be formulated as the following observation.

Observation 1. For $\lambda_i < 1$ we have that $bias_{BP}^{2(i)} < bias_{LS}^{2(i)}$ but $var_{BP}^{(i)} > var_{LS}^{(i)}$. And, for $\lambda_i > 1$ we have that $bias_{BP}^{2(i)} > bias_{LS}^{2(i)}$ but $var_{BP}^{(i)} < var_{LS}^{(i)}$.

Notice that in the noiseless case $\sigma_e = 0$, implying that $MSE_{LS} = bias_{LS}^2$ and $MSE_{BP} = bias_{BP}^2$. This leads us to the following observation for the noiseless case.

Observation 2. In a noiseless scenario, the relation between $\sum_{i=1}^m bias_{BP}^{2(i)}$ and $\sum_{i=1}^m bias_{LS}^{2(i)}$, dictates the relation between MSE_{BP} and MSE_{LS} . In particular, if all the singular values

of \mathbf{A} are smaller than 1, then $MSE_{BP} < MSE_{LS}$, and if all the singular values of \mathbf{A} are greater than 1, then $MSE_{BP} > MSE_{LS}$.

Note that Observation 2 holds for any given setting of β that is used by the two estimators. Therefore, these relations between MSE_{BP} and MSE_{LS} hold also when β is tuned for best performance of each estimator.

In practice, a different value of β can be preferred for the different cost functions. Let us denote by β_{LS} and β_{BP} the regularization parameter in $\ell_{LS}(\tilde{\mathbf{x}})$ and $\ell_{BP}(\tilde{\mathbf{x}})$, respectively, and let the singular values of \mathbf{A} be indexed in a descending order, i.e. $\lambda_1 \geq \dots \geq \lambda_m$. Comparing MSE_{BP} and MSE_{LS} with $\beta_{BP} \neq \beta_{LS}$ leads to an additional observation for the noiseless case, which is in favor of the BP cost.

Observation 3. *In a noiseless scenario, for any β_{LS} and $\beta_{BP} = \beta_{LS}/\lambda_1^2$, we have that $MSE_{BP} \leq MSE_{LS}$. If in addition $[\mathbf{V}^T \mathbf{x}]_i \neq 0$ for some indices $2 \leq i \leq m$, then $MSE_{BP} < MSE_{LS}$ unless $\lambda_i = \lambda_1$ for all these indices.*

Proof. Since $\beta_{BP} = \beta_{LS}/\lambda_1^2$, we have that $\frac{\beta_{BP}\gamma_i^2}{1+\beta_{BP}\gamma_i^2} = \frac{\beta_{LS}\gamma_i^2}{\lambda_1^2 + \beta_{LS}\gamma_i^2}$. Therefore,

$$\begin{aligned} \sum_{i=1}^m bias_{BP}^{2(i)} &= \sum_{i=1}^m \left(\frac{\beta_{BP}\gamma_i^2}{1+\beta_{BP}\gamma_i^2} \right)^2 [\mathbf{V}^T \mathbf{x}]_i^2 \\ &= \sum_{i=1}^m \left(\frac{\beta_{LS}\gamma_i^2}{\lambda_1^2 + \beta_{LS}\gamma_i^2} \right)^2 [\mathbf{V}^T \mathbf{x}]_i^2 \\ &\leq \sum_{i=1}^m \left(\frac{\beta_{LS}\gamma_i^2}{\lambda_i^2 + \beta_{LS}\gamma_i^2} \right)^2 [\mathbf{V}^T \mathbf{x}]_i^2 = \sum_{i=1}^m bias_{LS}^{2(i)}. \end{aligned} \quad (16)$$

If $[\mathbf{V}^T \mathbf{x}]_i \neq 0$ for some indices $2 \leq i \leq m$, it is easy to see that the inequality is strict unless $\lambda_i = \lambda_1$ for these indices. Finally, recall that in the noiseless case the relation between $\sum_{i=1}^m bias_{BP}^{2(i)}$ and $\sum_{i=1}^m bias_{LS}^{2(i)}$, dictates the relation between MSE_{BP} and MSE_{LS} . \square

Even though Observation 2 and Observation 3 consider the noiseless case, note that they cover events where the gap between $bias_{LS}^2$ and $bias_{BP}^2$ may be substantial enough to dictate the relationship between the MSEs also when the noise level is moderate. For example, if $\beta_{BP} = \beta_{LS}$ and all the singular values are much smaller than 1 then the 'in particular'-part in Observation 2 implies that $\sum_{i=1}^m bias_{BP}^{2(i)}$ is much smaller than $\sum_{i=1}^m bias_{LS}^{2(i)}$. Another example, if $\beta_{BP} = \beta_{LS}/\lambda_1^2$ and the condition number of $\mathbf{A}\mathbf{A}^T$, i.e. the ratio λ_1^2/λ_m^2 , is very large, then Observation 3 implies that $\sum_{i=1}^m bias_{BP}^{2(i)}$ is much smaller than $\sum_{i=1}^m bias_{LS}^{2(i)}$.

C. Discussion and implications for priors beyond ℓ_2

As can be seen in (10) and (14), for the discussed Tikhonov regularization the bias term of each estimator is minimized if

$\beta \rightarrow 0$, and in this case $bias_{LS}^2$ tends to $bias_{BP}^2$. This means that the performance gap in the noiseless case, which is stated in Observation 2 and Observation 3, tends to zero for $\beta \rightarrow 0$. However, note that we consider here ℓ_2 priors mainly as a surrogate to complex priors which are hard to analyze. As we demonstrate in Section IV, the results that are obtained for sophisticated priors, such as TV, BM3D and DCGAN, indeed strongly correlate with the observations above (especially with Observation 3 that implies an advantage of BP for badly conditioned $\mathbf{A}\mathbf{A}^T$). For such priors, the optimal value of β for each fidelity term is significantly above 0 even in the noiseless case (contrary to ℓ_2 priors), and the gap between the best recoveries is significant as well.

Another motivation for connecting the above analysis to other priors comes from recognizing attributes that distinguish between the LS and BP fidelity terms regardless of the prior used with them. Let us focus on the noiseless case, where $\mathbf{y} = \mathbf{A}\mathbf{x}$. In this case, (5) and (6) can be written as

$$\begin{aligned} f_{LS}(\tilde{\mathbf{x}}) &= \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{A}\tilde{\mathbf{x}}\|_2^2 + \beta s(\tilde{\mathbf{x}}) \\ &= \frac{1}{2} (\mathbf{x} - \tilde{\mathbf{x}})^T \mathbf{A}^T \mathbf{A} (\mathbf{x} - \tilde{\mathbf{x}}) + \beta s(\tilde{\mathbf{x}}), \end{aligned} \quad (17)$$

$$\begin{aligned} f_{BP}(\tilde{\mathbf{x}}) &= \frac{1}{2} \|\mathbf{A}^\dagger \mathbf{A}\mathbf{x} - \mathbf{A}^\dagger \mathbf{A}\tilde{\mathbf{x}}\|_2^2 + \beta s(\tilde{\mathbf{x}}) \\ &= \frac{1}{2} (\mathbf{x} - \tilde{\mathbf{x}})^T \mathbf{P}_A (\mathbf{x} - \tilde{\mathbf{x}}) + \beta s(\tilde{\mathbf{x}}). \end{aligned} \quad (18)$$

Under our SVD notations, we have $\mathbf{A}^T \mathbf{A} = \sum_{i=1}^m \lambda_i^2 \mathbf{v}_i \mathbf{v}_i^T$ and $\mathbf{P}_A = \sum_{i=1}^m \mathbf{v}_i \mathbf{v}_i^T$, where \mathbf{v}_i is the right singular vector of \mathbf{A} associated with the singular value λ_i . Therefore, we get

$$f_{LS}(\tilde{\mathbf{x}}) = \frac{1}{2} \sum_{i=1}^m \lambda_i^2 |\mathbf{v}_i^T (\mathbf{x} - \tilde{\mathbf{x}})|^2 + \beta s(\tilde{\mathbf{x}}), \quad (19)$$

$$f_{BP}(\tilde{\mathbf{x}}) = \frac{1}{2} \sum_{i=1}^m |\mathbf{v}_i^T (\mathbf{x} - \tilde{\mathbf{x}})|^2 + \beta s(\tilde{\mathbf{x}}). \quad (20)$$

Note that $f_{BP}(\tilde{\mathbf{x}})$ equally weighs all $\{|\mathbf{v}_i^T (\mathbf{x} - \tilde{\mathbf{x}})|^2\}_{i=1}^m$, contrary to $f_{LS}(\tilde{\mathbf{x}})$ that weighs them according to $\{\lambda_i^2\}$. As in inverse problems one (typically) cares about minimizing the MSE, an *intuition* that minimizing (20) may have an advantage over minimizing (19) for *general* priors, comes from the similarity between the BP fidelity term and formulating the MSE as $\|\tilde{\mathbf{x}} - \mathbf{x}\|_2^2 = \sum_{i=1}^n |\mathbf{v}_i^T (\mathbf{x} - \tilde{\mathbf{x}})|^2$ (note that the sum here goes over all the n basis vectors in \mathbf{V}). For ℓ_2 priors, we indeed have shown in Section II-B that this "equal weighting" strategy translates to the fact that $\{bias_{BP}^{2(i)}\}$ do not depend on $\{\lambda_i^2\}$, contrary to $\{bias_{LS}^{2(i)}\}$, which later yields the MSE advantage of BP over LS in Observation 3. For ℓ_2 priors, we have obtained analytical results and tradeoffs also for the noisy case. For other priors, we empirically show in Section IV correlation to the above analytical findings.

An important factor that is not taken into account in the above analysis is optimization, since for ℓ_2 priors there is a closed-form solution. Yet, for sophisticated priors iterative optimization schemes are inevitable, and the regularization parameter has an effect which is similar to the step size

in these schemes. In such cases, extremely low value of β inherently results in a massive slowdown in the convergence for convex priors [26], [27] and/or bad local minima for non-convex priors. Taking a numerical optimization point of view, in the sequel we empirically show that $\hat{\mathbf{x}}_{BP}$ is superior to $\hat{\mathbf{x}}_{LS}$ even for ℓ_2 priors with $\beta \rightarrow 0$, if few iterations of conjugate gradients are used instead of the closed-form expressions (7) and (8). This implementation choice may be preferable in high-dimensional problems when it is not possible to invert the matrices. The advantage of BP in this case follows from the fact that the eigenvalues of \mathbf{P}_A are only 1 (in the row space of \mathbf{A}) and 0 (in the null space of \mathbf{A}), while $\mathbf{A}^T \mathbf{A}$ may have very different eigenvalues in general, and conjugate gradients (among other methods) performs better when the eigenvalues are clustered [28]. In Section IV we provide empirical evidence that BP requires less iterations than LS also for other optimization schemes and priors.

III. EXPERIMENTS WITH ℓ_2 PRIORS

In this section, we discuss the implications of the analytical results from Section II and verify them for specific observation models: super-resolution and compressed sensing. In the first, all the singular values of \mathbf{A} are smaller than 1 and the condition number of $\mathbf{A}\mathbf{A}^T$ is large, while in the latter it is possible that all singular values are greater than 1 and that the condition number is very moderate. We also discuss the typical deblurring problem, which is highly ill-conditioned. In this case, \mathbf{A}^\dagger in $\hat{\mathbf{x}}_{BP}$ has to be regularized due to the large number of near zero singular values, and (13) needs to be modified accordingly.

Throughout this section, we use the closed-form estimators in (7) and (8) to restore the images. The empirical performance of these two estimators is presented by markers, while the analytical expressions from (11) and (15) are plotted in solid curves. Different colors are used to distinguish between the two fidelity terms that are used for the estimation.

A. Super-resolution

Let us consider the super-resolution (SR) task, where \mathbf{A} is a composite operator of blurring (e.g. anti-aliasing filtering) followed by down-sampling. Note that the largest singular value of a typical low-pass filtering operation is 1, and it is associated with the DC (i.e. the magnitude of the Fourier coefficient that is associated with zero frequency). The rest of the singular values are smaller than 1. The subsequent operator is subsampling, which inevitably reduces the energy of the signal (as $m < n$). Therefore, essentially, all the singular values of \mathbf{A} are smaller than 1. Accordingly, the condition number of $\mathbf{A}\mathbf{A}^T$ is large. These properties are demonstrated in Fig. 1a for SR with scale factor 3 and Gaussian filter of size 7×7 and standard deviation 1.6 (used in many works, e.g. [11], [23], [29]), which is performed on a 64×64 image (thus $n = 4096$ and $m = 484$). We consider such a small image to allow computing the SVD of \mathbf{A} (our analytic expressions require both $\{\lambda_i^2\}_{i=1}^m$ and \mathbf{V}).

We verify our analytical results for the SRx3 scenario mentioned above, and two cases: $\sigma_e = 0$ and Gaussian

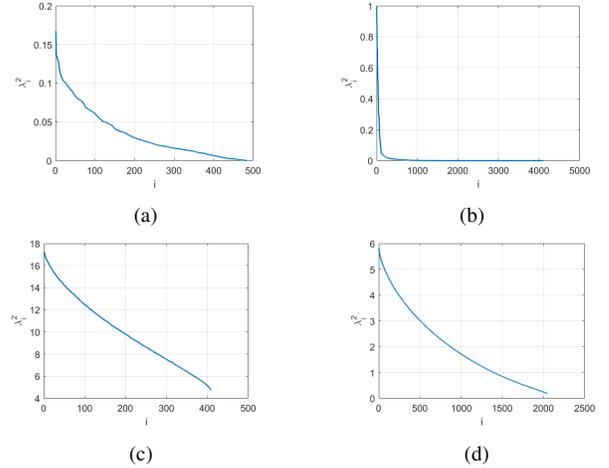


Fig. 1: The (squared) singular values of \mathbf{A} applied on a 64×64 image for: (a) SRx3 with 7×7 Gaussian filter ($\frac{\lambda_1^2}{\lambda_m^2} = 2.93e3$); (b) blurring with 9×9 uniform filter ($\frac{\lambda_1^2}{\lambda_m^2} = 1.46e7$); (c) CS with $m = 0.1n$ Gaussian measurements and Haar basis ($\frac{\lambda_1^2}{\lambda_m^2} = 3.63$); (d) CS with $m = 0.5n$ Gaussian measurements and Haar basis ($\frac{\lambda_1^2}{\lambda_m^2} = 33.36$).

noise with $\sigma_e = \sqrt{2}$. The experiments are performed on the *cameraman* image, resized to 64×64 pixels. In the noisy case, we average the results over 5 noise realizations. We have observed similar results for other images as well. We use the ℓ_2 prior $s(\tilde{\mathbf{x}}) = \frac{1}{2} \|\tilde{\mathbf{x}}\|_2^2$, which satisfies the assumptions ($\mathbf{D} = \mathbf{I}_n$ and $\gamma_i = 1$).

The PSNR² results are presented in Fig. 2 and validate the analytical expressions. For $\sigma_e = 0$, $\hat{\mathbf{x}}_{BP}$ is better than $\hat{\mathbf{x}}_{LS}$ for any value of the parameter β , as implied by Observation 2 since all the singular values of \mathbf{A} are smaller than 1 (Fig. 1a). The rather large gap in favor of BP also agrees with Observation 3 that predicts it when the ratio λ_1^2/λ_m^2 is large. The fact that BP at $\beta/\lambda_1^2 = 5.97\beta$ outperforms LS at β , further verifies Observation 3. For $\sigma_e = \sqrt{2}$, the gap between the estimators is reduced because var_{BP} is worse than var_{LS} at handling the small singular values, as mentioned in Observation 1.

To demonstrate the numerical optimization advantage of the BP cost over the LS cost for $\beta \rightarrow 0$ (where the gap between the bias terms in (10) and (14) tends to 0), we repeat the experiments above for very small values of β . However, this time instead of inverting the matrices in (7) and (8) we obtain the estimators using the conjugate gradient method. The results are presented in Fig. 3. Remarkably, a single iteration is enough for obtaining the exact BP estimator (for ℓ_2 prior).

B. Compressed sensing

Contrary to SR scenarios, in compressed sensing (CS) the condition number of $\mathbf{A}\mathbf{A}^T$ is moderate and the singular values of \mathbf{A} may be larger than 1. Consider the commonly examined scenario where \mathbf{A} is the multiplication of an $m \times n$

²The PSNR for a recovery $\hat{\mathbf{x}}$ of a uint8 image $\mathbf{x} \in \mathbb{R}^n$ is computed as $10 \log_{10} \left(\frac{255^2}{\frac{1}{n} \|\hat{\mathbf{x}} - \mathbf{x}\|_2^2} \right)$.

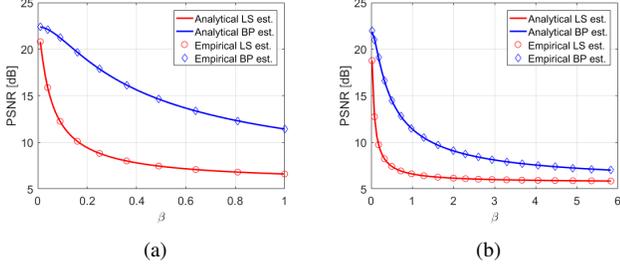


Fig. 2: Super-resolution with Gaussian filter and scale factor of 3, using ℓ_2 prior. PSNR (for *cameraman*) vs. β (regularization parameter), for (a) $\sigma_e = 0$, and (b) $\sigma_e = \sqrt{2}$.

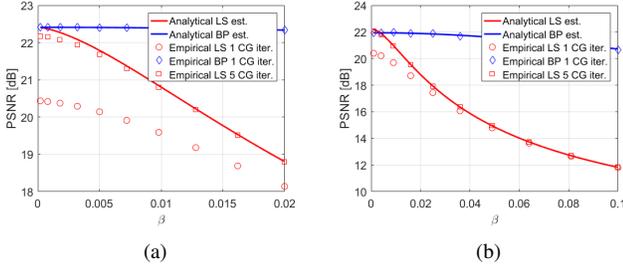


Fig. 3: Super-resolution with Gaussian filter and scale factor of 3, using ℓ_2 prior and iterations of conjugate gradients instead of matrix inversion. PSNR (for *cameraman*) vs. β (regularization parameter), for (a) $\sigma_e = 0$, and (b) $\sigma_e = \sqrt{2}$. Note that the LS cost requires more CG iterations than the BP cost to attend the solution (solid line).

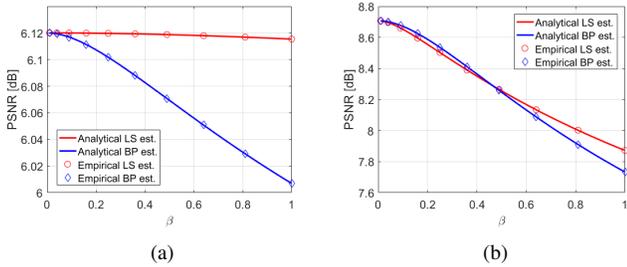


Fig. 4: Compressed sensing with Gaussian measurements and Haar basis, using ℓ_2 prior. PSNR (for *cameraman*) vs. β (regularization parameter), for (a) $m = 0.1n$, and (b) $m = 0.5n$.

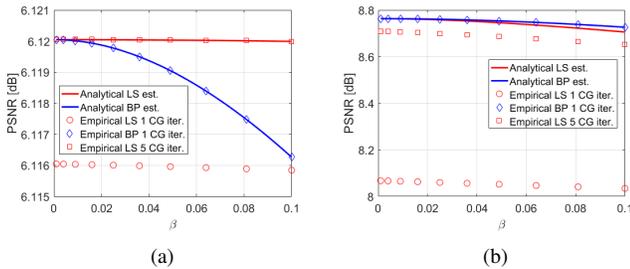


Fig. 5: Compressed sensing with Gaussian measurements and Haar basis, using ℓ_2 prior and iterations of conjugate gradients instead of matrix inversion. PSNR (for *cameraman*) vs. β (regularization parameter), for (a) $m = 0.1n$, and (b) $m = 0.5n$. Note that the LS cost requires more CG iterations than the BP cost to attend the solution (solid line).

Gaussian measurement matrix (whose i.i.d. entries are drawn from $\mathcal{N}(0, 1/m)$) with an $n \times n$ Haar wavelet basis. We have observed that for high compression, e.g. $m/n = 0.1$, all the singular values are larger than 1 and the condition number is very small, as demonstrated in Fig. 1c. However, for lower compression, e.g. $m/n = 0.5$, there are also singular values smaller than 1 and the condition number increases, as demonstrated in Fig. 1d.

We verify our analytical results for these two compression ratios (both with $\sigma_e = 0$). The experiments are performed on the same 64×64 version of *cameraman* image, and we use again the ℓ_2 prior $s(\tilde{\mathbf{x}}) = \frac{1}{2} \|\tilde{\mathbf{x}}\|_2^2$. The results are presented in Fig. 4 and validate the analytical expressions. For $m/n = 0.1$, $\hat{\mathbf{x}}_{LS}$ is better than $\hat{\mathbf{x}}_{BP}$ for any value of β , as implied by Observation 2 since all the singular values of \mathbf{A} are greater than 1 (Fig. 1c). To verify Observation 3 in this case, see that BP at $\beta/\lambda_1^2 = 0.058\beta$ has (slightly) higher PSNR than LS at β , e.g. for $\beta = 1$. We have verified this also for very large values of β (not presented here)—both curves decrease and reach a similar plateau at high β , yet BP at β/λ_1^2 indeed has higher PSNR than LS at β , but the difference is extremely small. Interestingly, for $m/n = 0.5$, where some singular values of \mathbf{A} are smaller than 1 (Fig. 1d), $\hat{\mathbf{x}}_{BP}$ gets better results than $\hat{\mathbf{x}}_{LS}$. Also in this scenario, it can be verified that BP at $\beta/\lambda_1^2 = 0.171\beta$ has higher PSNR than LS at β , as implied by Observation 3. The fact that the gap between BP and LS for $m/n = 0.1$ and $m/n = 0.5$ has been changed in favor of BP in the latter agrees with the derivation of Observation 3 in (16) that links the advantage of BP to an increased λ_1^2/λ_m^2 ratio.

We demonstrate again the numerical optimization advantage of the BP cost over the LS cost by repeating the experiments above for very small values of β , while using the conjugate gradient method instead of matrix inversion. The results are presented in Fig. 5. It can be seen again that for the ℓ_2 prior a single iteration is enough for obtaining the exact BP estimator.

We find it necessary to emphasize that compressed sensing scenarios require a sparsity-inducing prior, e.g. $s(\tilde{\mathbf{x}}) = \|\tilde{\mathbf{x}}\|_1$ or TV prior, rather than an ℓ_2 prior, for which both estimators exhibit poor results (i.e. very low PSNR). However, our purpose here is merely to validate our analysis, which applies only to ℓ_2 priors, for a case in which all the singular values are greater than 1 and/or the condition number is small.

Finally, note that for Gaussian \mathbf{A} there is no efficient way to implement the operators \mathbf{A} and \mathbf{A}^T for large dimensions. Therefore, in practice, taking \mathbf{A} to be the subsampled Fourier transform is more common, e.g. in sparse MRI [30]. However, note that for this acquisition model \mathbf{A}^\dagger is simply the Hermitian transpose of \mathbf{A} (this property follows from the fact that the subsampled Fourier transform is a tight frame [31]), which together with the unitarity of the Fourier transform leads to $\|\mathbf{A}^\dagger(\mathbf{y} - \mathbf{A}\tilde{\mathbf{x}})\|_2^2 = \|\mathbf{y} - \mathbf{A}\tilde{\mathbf{x}}\|_2^2$. This means that the two cost functions coincide, which is also implied by the fact that in this case all the singular values of \mathbf{A} are 1 and thus (9) is identical to (13). Therefore, we do not make a comparison for this case.

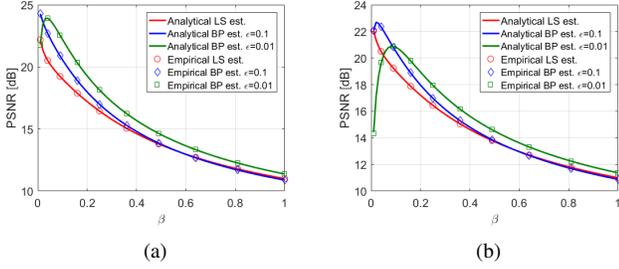


Fig. 6: Deblurring with uniform 9×9 blur kernel, using ℓ_2 prior. PSNR (for *cameraman*) vs. β (regularization parameter), for (a) $\sigma_e = \sqrt{0.3}$, and (b) $\sigma_e = \sqrt{2}$.

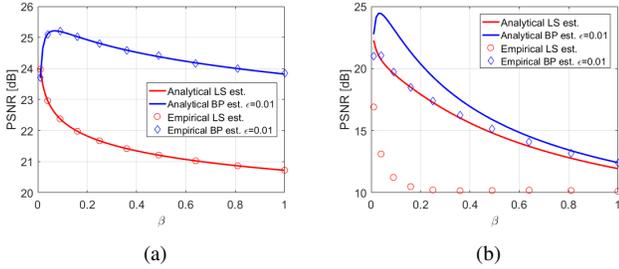


Fig. 7: Deblurring with uniform 9×9 blur kernel and $\sigma_e = \sqrt{0.3}$, using priors with different $D^T D$. PSNR (for *cameraman*) vs. β (regularization parameter), for (a) circulant $D^T D$, and (b) non-circulant $D^T D$. Note that $D^T D \neq V\Gamma^2 V^T$ in (b).

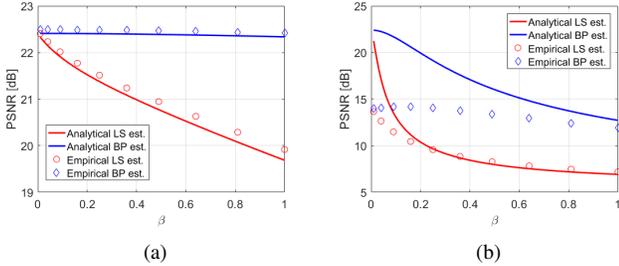


Fig. 8: Super-resolution with Gaussian filter and scale factor of 3 and $\sigma_e = 0$, using priors with different $D^T D$. PSNR (for *cameraman*) vs. β (regularization parameter), for (a) circulant $D^T D$, and (b) non-circulant $D^T D$. Note that $D^T D \neq V\Gamma^2 V^T$ in (a) and (b).

C. Deblurring

In the deblurring problem, \mathbf{A} is a square ($m = n$) ill-conditioned matrix that performs blurring (i.e. filtering by a blur kernel). Typically, the blur kernel coefficients are normalized such that their sum is 1. Thus, the largest singular value of \mathbf{A} is 1 (associated with the DC), and many other singular values are near 0. Accordingly, the condition number of $\mathbf{A}\mathbf{A}^T$ is extremely large. These properties are demonstrated in Fig. 1b for uniform kernel of size 9×9 (used in many works, e.g. [8], [9], [18]).

Note that if one uses $\hat{\mathbf{x}}_{BP}$, exactly as defined in (8), both Observation 2 and Observation 3 imply an advantage of $\hat{\mathbf{x}}_{BP}$ over $\hat{\mathbf{x}}_{LS}$ in the noiseless case due to small singular values and a large condition number, respectively. However, since in the deblurring problem \mathbf{A} is not rank-deficient but rather (very) ill-conditioned, deblurring scenarios always assume that

the measurements are noisy (typically with low noise levels). Therefore, it is required to regularize the inversion of $\mathbf{A}\mathbf{A}^T$ in \mathbf{A}^\dagger in order to mitigate the effect of near zero $\{\lambda_i\}$ on the variance of $\hat{\mathbf{x}}_{BP}$. A common regularized inversion is diagonal loading: inverting $\mathbf{A}\mathbf{A}^T + \epsilon \mathbf{I}_n$ instead of $\mathbf{A}\mathbf{A}^T$, where ϵ is a parameter. This is equivalent to replacing λ_i^2 with $\lambda_i^2 + \epsilon$ in the eigen-decomposition of $\mathbf{A}\mathbf{A}^T$.

For $\hat{\mathbf{x}}_{BP}$ with such a regularized inversion, it is not hard to repeat the computations in (13) and obtain a very similar result, where $1_{i \leq m}$ is replaced with $\lambda_i^2 / (\lambda_i^2 + \epsilon)$ and $\lambda_i^{-2} 1_{i \leq m}$ is replaced with $\lambda_i^2 / (\lambda_i^2 + \epsilon)^2$. Formally, we get

$$\begin{aligned} MSE_{BP} &= \sum_{i=1}^n \left(\frac{\lambda_i^2 / (\lambda_i^2 + \epsilon)}{\lambda_i^2 / (\lambda_i^2 + \epsilon) + \beta \gamma_i^2} - 1 \right)^2 [\mathbf{V}^T \mathbf{x}]_i^2 \\ &+ \sigma_e^2 \sum_{i=1}^n \frac{\lambda_i^2 / (\lambda_i^2 + \epsilon)^2}{(\lambda_i^2 / (\lambda_i^2 + \epsilon) + \beta \gamma_i^2)^2}, \\ &= \sum_{i=1}^n \left(\frac{\beta \gamma_i^2}{\lambda_i^2 / (\lambda_i^2 + \epsilon) + \beta \gamma_i^2} \right)^2 [\mathbf{V}^T \mathbf{x}]_i^2 \\ &+ \sigma_e^2 \sum_{i=1}^n \frac{1}{\lambda_i^2 (1 + \beta \gamma_i^2 (\lambda_i^2 + \epsilon) / \lambda_i^2)^2}. \end{aligned} \quad (21)$$

Therefore, as could be expected, increasing the amount of regularization ϵ reduces the variance of $\hat{\mathbf{x}}_{BP}$ but increases its bias. As a sanity check, observe that for $\epsilon \rightarrow 0$ we get that (21) coincides with (13) (recall $m = n$). Since in this case the performance of $\hat{\mathbf{x}}_{BP}$ depends on the couple (β, ϵ) , we cannot obtain clear properties like the observations in Section II-B that hold uniformly for any parameter setting. Yet, as demonstrated below and in the sequel, we have empirically observed that it is possible to find settings of (β, ϵ) that balance the bias and variance of $\hat{\mathbf{x}}_{BP}$ and therefore lead to very good results despite the observed noise.

We verify (21) for the uniform blur kernel mentioned above, and two levels of Gaussian noise: $\sigma_e = \sqrt{0.3}$ and $\sigma_e = \sqrt{2}$. The experiments are performed on the 64×64 version of *cameraman* image, and we use again the ℓ_2 prior. The results are presented in Fig. 6. They show that $\hat{\mathbf{x}}_{BP}$ with good tuning of (β, ϵ) can outperform $\hat{\mathbf{x}}_{LS}$, especially when the noise level is low. This implies that "well-tuned" $\hat{\mathbf{x}}_{BP}$ handles the badly conditioned \mathbf{A} (Fig. 1b) better than $\hat{\mathbf{x}}_{LS}$.

D. The effect of the joint right singular vectors assumption

In this section, we compare the empirical MSE and the analytical formulas in (11), (15) and (21) in cases where the condition $D^T D = V\Gamma^2 V^T$ is violated (recall that the columns of \mathbf{V} are the right singular vectors of \mathbf{A} and Γ^2 is a diagonal matrix, as defined in Section II-A). Since our formulas require the diagonal of Γ^2 (i.e. $\{\gamma_i^2\}$), we compute it as the diagonal of $\mathbf{V}^T D^T D \mathbf{V}$, which is exact under the analysis assumption and can be regarded as an approximation when $D^T D \neq V\Gamma^2 V^T$.

We start with examining the case of $D^T D = \Omega_{DIF}^T \Omega_{DIF} + 0.01 \mathbf{I}_n$, where Ω_{DIF} is the 2D finite difference operator and the diagonal loading is required to make $D^T D \succ 0$. Note that for the deblurring task we have that both \mathbf{A} and $D^T D$ are circulant matrices that can be diagonalized by the

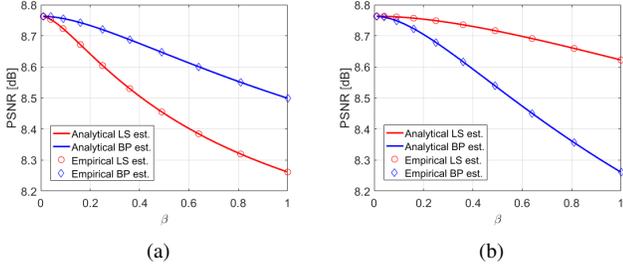


Fig. 9: Compressed sensing with $m = 0.5n$ Gaussian measurements and Haar basis, using ℓ_2 prior. PSNR (for \mathbf{x} that is the projection of *cameraman* onto \mathcal{W}^\perp) vs. β (regularization parameter), for (a) \mathcal{W} that is spanned by the columns of $\mathbf{V}(:, 1:m-500)$, and (b) \mathcal{W} that is spanned by the columns of $\mathbf{V}(:, 1000:m)$.

DFT matrix. Therefore, the condition $\mathbf{D}^T \mathbf{D} = \mathbf{V} \Gamma^2 \mathbf{V}^T$ holds for \mathbf{V} that equals the (inverse) DFT matrix. However, for the SR task \mathbf{A} cannot be singularly decomposed by a Fourier basis. Therefore, the condition cannot be satisfied.

We repeat previous deblurring and SR experiments with the examined $\mathbf{D}^T \mathbf{D}$. The results are presented in Figs. 7a and 8a. For deblurring we see perfect agreement between the empirical results and the analytical formulas (as expected). For SR we see that violating the condition has led to a small gap between the empirical results and the formulas.

Now, we further increase the violation of the condition by breaking the circularity property of $\mathbf{D}^T \mathbf{D}$. We do it by replacing Ω_{DIF} with a non-circulant operator $\tilde{\Omega}$ that performs finite difference only on every 8th pixel (and identity on the rest). We repeat the previous deblurring and SR experiments and present the results in Figs. 7b and 8b. It is easy to see that the deviation of the empirical results from the formulas further grows for both tasks.

The experiments demonstrate that the deviation between the empirical MSE and the analytical expressions is proportional to how much the condition on $\mathbf{D}^T \mathbf{D}$ is violated. Yet, the overall trend in the curves still shares similarity with the analytical results, which motivates considering the observations obtained by the analytical analysis for practical sophisticated priors.

E. Incorporating prior knowledge on \mathbf{x} with the results

The analytical MSE formulas in (11) and (15) are conditioned on the latent image \mathbf{x} , as the expectations are taken only with respect to the noise \mathbf{e} . These expressions have led to observations in Section II that depend only the singular values of \mathbf{A} (i.e. $\{\lambda_i\}$) and do not require prior knowledge on \mathbf{x} (recall that accurately modeling natural images is difficult). The usefulness of these observations for preferring one fidelity term over the other for sophisticated priors is demonstrated in Section IV.

However, a natural question arises: How can one leverage prior knowledge on \mathbf{x} to improve the criterion for choosing the fidelity term?

In this section we briefly demonstrate, using a controlled experiment, how the observations in Section II can be polished given a constraint that \mathbf{x} resides in \mathcal{W}^\perp the orthogonal

complement of a *known* subspace \mathcal{W} . Note that for low-dimensional \mathcal{W} such that $m < n - \dim(\mathcal{W})$, we still have an ill-posed linear inverse problem.

We consider the compressed sensing scenario from Section III-B where $n = 64^2$, $m/n = 0.5$ and $\sigma_e = 0$. In this case \mathbf{A} has (more than 1000) singular values that are larger than 1 and (slightly more than 500) singular values that are smaller than 1 (see Fig. 1d). Therefore, the events in the ‘in particular’-part in Observation 2 do not occur. Indeed, observe that in Fig. 4b none of the estimators is consistently (i.e. for any β) better than the other when \mathbf{x} is the *cameraman* image.

Now, let us use the notation from Section II, where the columns of \mathbf{V} , that is, the right singular vectors of \mathbf{A} , are ordered according to a descending order of the singular values (from 1 to m), and the last $n - m$ columns span the null space of \mathbf{A} . Suppose that \mathcal{W} is the subspace spanned by the columns of $\mathbf{V}(:, 1:m-500)$, where we use Matlab notation, and that $\mathbf{x} \in \mathcal{W}^\perp$. Due to the orthogonality of \mathbf{V} , we have that $[\mathbf{V}^T \mathbf{x}]_i = 0$ for any $1 \leq i \leq m-500$. Substituting this property in (10) and (14), we get

$$\begin{aligned} bias_{LS}^2 &= \sum_{i=m-499}^m bias_{LS}^{2(i)} + \sum_{i=m+1}^n [\mathbf{V}^T \mathbf{x}]_i^2, \\ bias_{BP}^2 &= \sum_{i=m-499}^m bias_{BP}^{2(i)} + \sum_{i=m+1}^n [\mathbf{V}^T \mathbf{x}]_i^2. \end{aligned} \quad (22)$$

Therefore, for the considered CS scenario, we have that $bias_{BP}^2 < bias_{LS}^2$ for any β (because $\lambda_i < 1$ for all $m-499 \leq i \leq m$). Since $\sigma_e = 0$, this implies that $MSE_{BP} < MSE_{LS}$ for any β .

Note that for \mathcal{W} that is the subspace spanned by the columns of $\mathbf{V}(:, 1000:m)$ and $\mathbf{x} \in \mathcal{W}^\perp$ (i.e. \mathbf{x} in a subspace spanned by columns of \mathbf{V} that are either associated with singular values that are greater than 1 or with the null space of \mathbf{A}), similar arguments lead to $MSE_{BP} > MSE_{LS}$ for any β . Fig. 9 verifies both results for a test image \mathbf{x} that is the projection of the *cameraman* image onto \mathcal{W}^\perp (i.e. $\mathbf{x} = \mathbf{P}_{\mathcal{W}^\perp} \mathbf{x}_0$, where \mathbf{x}_0 is the *cameraman* image).

Note that the behavior in Fig. 9 cannot be predicted by the ‘in particular’-part in Observation 2 that considers *all* the singular values of \mathbf{A} , *regardless* of \mathbf{x} . We believe that a detailed study with constraints on \mathbf{x} that better fit images is an interesting direction for future research.

IV. EXPERIMENTS WITH SOPHISTICATED PRIORS

In this section we empirically demonstrate that the behavior of $\hat{\mathbf{x}}_{BP}$ and $\hat{\mathbf{x}}_{LS}$ (the minimizers of $f_{BP}(\tilde{\mathbf{x}})$ and $f_{LS}(\tilde{\mathbf{x}})$) for sophisticated convex and non-convex priors (for whom mathematical analysis is hard or even intractable) strongly correlates with properties for which we have established concrete mathematical reasoning in the case of ℓ_2 priors. Specifically, for super-resolution and deblurring tasks (where the condition number of $\mathbf{A}\mathbf{A}^T$ is very large) BP cost function can lead to significantly improved results compared to the LS cost function, yet, there is inverse proportion between the performance gap and the noise level (since the singular

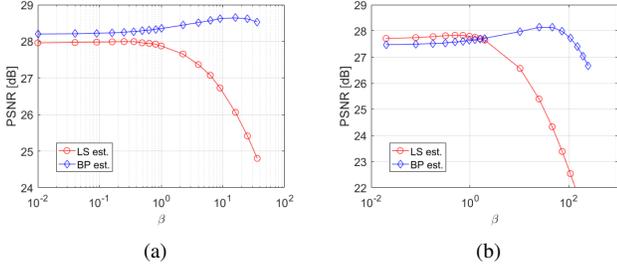


Fig. 10: Super-resolution with Gaussian filter and scale factor of 3, using TV prior and 100 iterations of FISTA. PSNR (averaged over 8 test images) vs. β (regularization parameter), for (a) $\sigma_e = 0$, and (b) $\sigma_e = \sqrt{2}$.

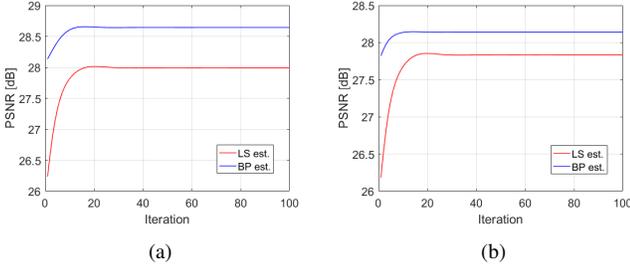


Fig. 11: Super-resolution with Gaussian filter and scale factor of 3, using TV prior. PSNR (for best uniform setting of β , averaged over 8 test images) vs. FISTA iteration number, for (a) $\sigma_e = 0$, and (b) $\sigma_e = \sqrt{2}$.

values of \mathbf{A} are small in these tasks). For Gaussian compressed sensing with low m/n ratio (where the condition number is small and the singular values are greater than 1) $\hat{\mathbf{x}}_{BP}$ is not significantly better than $\hat{\mathbf{x}}_{LS}$, but it is quite robust to noise. However, when the m/n ratio increases (then the condition number increases and some singular values are smaller than 1) the advantage of BP is more significant, but inversely proportional to the noise level.

A. TV prior

We start with the widely-used (isotropic) total-variation (TV) prior [1], which is given by

$$s(\tilde{\mathbf{x}}) = 0.1 \sum_{i,j} \sqrt{|\tilde{x}_{i+1,j} - \tilde{x}_{i,j}|^2 + |\tilde{x}_{i,j+1} - \tilde{x}_{i,j}|^2} \quad (23)$$

for a two-dimensional signal $\tilde{\mathbf{x}}$. The factor 0.1 is used to achieve good performance for $\beta = \sigma_e^2$ in case of denoising ($\mathbf{A} = \mathbf{I}_n$). Obviously, it does not affect the comparison between the methods, since $s(\tilde{\mathbf{x}})$ is multiplied by β that can be set arbitrarily. Note that $s(\tilde{\mathbf{x}})$ is convex, and thus $f_{LS}(\tilde{\mathbf{x}})$ and $f_{BP}(\tilde{\mathbf{x}})$ are also convex functions. We choose to minimize them by the same method: 100 iterations of FISTA [20], which is basically a variant of ISTA (see (32) in the appendix) that is incorporated with Nesterov's accelerated gradient [32]. The step size μ is the typical 1 over the Lipschitz constant of $\nabla \ell(\tilde{\mathbf{x}})$, which in our case can be computed as 1 over the spectral norm of the constant Hessian matrix $\nabla^2 \ell$, i.e. $\mu = 1/\|\mathbf{P}_A\| = 1$ for BP recovery and $\mu = 1/\|\mathbf{A}^T \mathbf{A}\|$ (computed by the power method) for LS recovery. This common choice of step size is known to ensure convergence

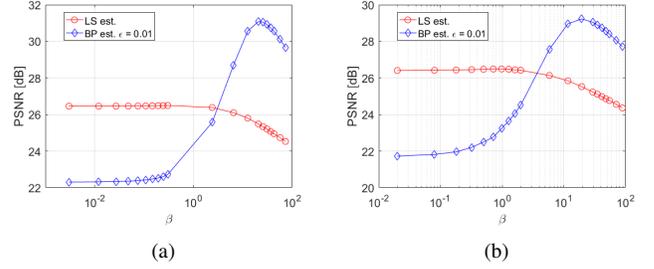


Fig. 12: Deblurring with uniform 9×9 blur kernel, using TV prior and 100 iterations of FISTA. PSNR (averaged over 8 test images) vs. β (regularization parameter), for (a) $\sigma_e = \sqrt{0.3}$, and (b) $\sigma_e = \sqrt{2}$.

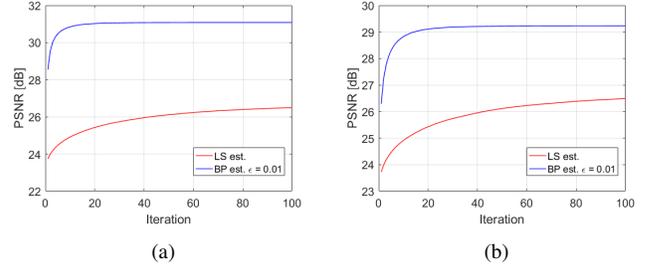


Fig. 13: Deblurring with uniform 9×9 blur kernel, using TV prior. PSNR (for best uniform setting of β , averaged over 8 test images) vs. FISTA iteration number, for (a) $\sigma_e = \sqrt{0.3}$, and (b) $\sigma_e = \sqrt{2}$.

in the convex setting [20]. Several methods for performing proximal mapping of $s(\tilde{\mathbf{x}})$ (i.e. Gaussian denoising associated with the TV prior) exist [33], [34]. Here, we choose to apply split Bregman method [33]. The experiments are performed on the following eight classical test images: *cameraman*, *house*, *peppers*, *Lena*, *Barbara*, *boat*, *hill* and *couple*.

1) *Super-resolution*: We compare the performance of $\hat{\mathbf{x}}_{LS}$ and $\hat{\mathbf{x}}_{BP}$ for SR with Gaussian anti-aliasing kernel (defined in Section III-A) and scale factor of 3. We consider the noiseless case $\sigma_e = 0$, as well as the case of Gaussian noise with $\sigma_e = \sqrt{2}$. For both estimators we initialize FISTA with the bicubic upsampling of \mathbf{y} . For BP, the operator \mathbf{A}^\dagger has fast implementation using the conjugate gradient method [35]. Fig. 10 shows the PSNR of the reconstructions, averaged over all images, for different values of the regularization parameter β . Fig. 11 shows the average PSNR as a function of the iteration number, where for each estimator we use the value of β which has led to its best results in Fig. 10 (0.25 for LS and 16 for BP in Fig. 11a; 0.5 for LS and 46 for BP in Fig. 11b). It can be seen that $\hat{\mathbf{x}}_{BP}$ converges somewhat faster than $\hat{\mathbf{x}}_{LS}$. In Figs. 21c and 21d we also display the results for *cameraman* image in the noiseless case.

Note the agreement of the obtained results with the observations from Section II, even though they have been established for a much simpler convex prior. In the noiseless case, $\hat{\mathbf{x}}_{BP}$ outperforms $\hat{\mathbf{x}}_{LS}$ for any value of β , while in the noisy scenario, this does not hold. However, even in the latter case, $\hat{\mathbf{x}}_{BP}$ (with good tuning of β) outperforms $\hat{\mathbf{x}}_{LS}$ (with good tuning of β). Yet, the gap between them (for optimal tuning) is smaller than in the noiseless case.

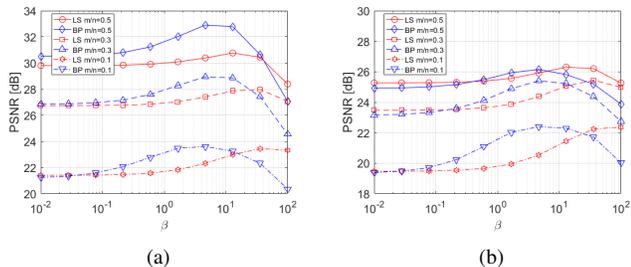


Fig. 14: Compressed sensing with Gaussian measurements, using TV prior and 500 iterations of FISTA. PSNR (averaged over 8 test images) vs. β (regularization parameter), for (a) $\sigma_e = 0$, and (b) SNR of 20 dB.

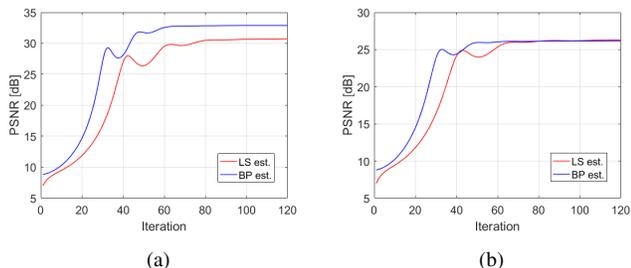


Fig. 15: Compressed sensing with $m = 0.5n$ Gaussian measurements, using TV prior. PSNR (for best uniform setting of β , averaged over 8 test images) vs. FISTA iteration number, for (a) $\sigma_e = 0$, and (b) SNR of 20 dB.

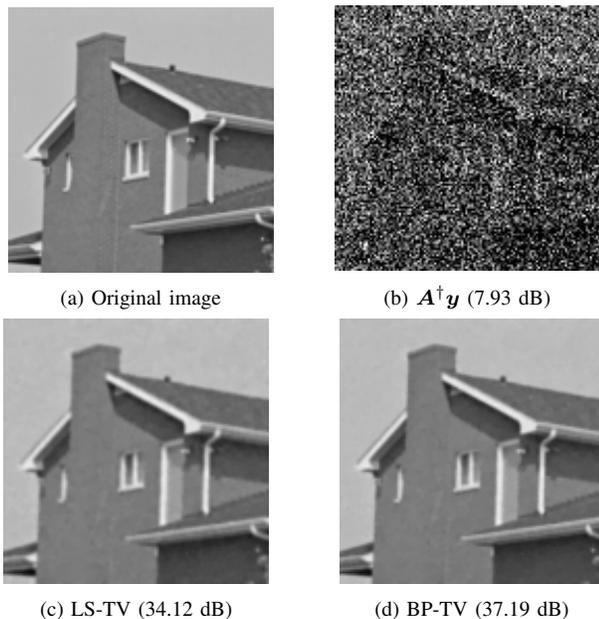


Fig. 16: Compressed sensing with $m = 0.5n$ Gaussian measurements and $\sigma_e = 0$ for *house* image. From left to right and from top to bottom: original image, naive $\mathbf{A}^\dagger \mathbf{y}$, reconstruction of LS fidelity with TV prior, and reconstruction of BP fidelity with TV prior.

2) *Deblurring*: We compare the two estimators for the widely examined 9×9 uniform blur kernel mentioned in Section III-C. We make the common assumption of circular shift-invariant blur operator, which allows very fast implementation of the gradient steps in the optimization of both cost functions using Fast Fourier Transform (FFT). We consider two levels of

Gaussian noise: $\sigma_e = \sqrt{0.3}$ and $\sigma_e = \sqrt{2}$. For both estimators we initialize FISTA with \mathbf{y} , and for $\hat{\mathbf{x}}_{BP}$ we use $\epsilon = 0.01\sigma_e^2$. Fig. 12 shows the average PSNR for different values of β , and Fig. 13 shows the average PSNR as a function of the iteration number, where each estimator uses the best β from Fig. 12 (0.3 for LS and 20.5 for BP in Fig. 13a; 0.98 for LS and 19.5 for BP in Fig. 13b). Note that $\hat{\mathbf{x}}_{BP}$ converges much faster than $\hat{\mathbf{x}}_{LS}$. The difference here for deblurring is more significant than for SR. Visual results for *couple* image in the case of $\sigma_e = \sqrt{2}$ are presented in Figs. 22c and 22d.

The obtained results agree with the observations in Section III-C, in the sense that there exist settings of (β, ϵ) for which $\hat{\mathbf{x}}_{BP}$ outperforms $\hat{\mathbf{x}}_{LS}$. Presumably, even for the more complex TV prior, this is due to a better handling of \mathbf{A} whose condition number is very large. As expected, the performance gap between the estimators (for optimal tuning) decreases when the noise level is higher. However, it is still highly in favor of $\hat{\mathbf{x}}_{BP}$.

3) *Compressed sensing*: We compare the performance of $\hat{\mathbf{x}}_{LS}$ and $\hat{\mathbf{x}}_{BP}$ for CS with Gaussian measurement matrix (i.e. $A_{ij} \sim \mathcal{N}(0, 1/m)$), for which the two cost functions differ (see the discussion in Section III-B). In these CS experiments (only) we decrease the size of the test images to 128×128 pixels, as there is no efficient way to implement the operators \mathbf{A} and \mathbf{A}^T for large dimensions. We consider compression ratios of $m/n = 0.1$, $m/n = 0.3$, and $m/n = 0.5$. For each of them we examine the noiseless case and the case of Gaussian noise with signal-to-noise ratio (SNR) of 20 dB. For both estimators we initialize FISTA with zero and use 500 iterations. As we compute \mathbf{A}^\dagger in advance, both estimators have similar computational cost per iteration. Fig. 14 shows the average PSNR for different values of β . For $m/n = 0.5$ we show in Fig. 15 the average PSNR vs. the iteration number, where each estimator uses the best β from Fig. 14. Again, note that $\hat{\mathbf{x}}_{BP}$ requires less iterations than $\hat{\mathbf{x}}_{LS}$. Visual results for *house* image in the noiseless case are presented in Fig. 16.

The results show correlation with the observations in Section II. In the noiseless case, when the m/n ratio increases (and thus the condition number of $\mathbf{A}\mathbf{A}^T$ increases, e.g. see Figs. 1c and 1d) the performance gap between BP and LS increases in favor of BP. In the noisy case, when the m/n ratio increases the BP estimator becomes more sensitive to noise (due to the increase in the number of singular values that are smaller than 1, again, see Figs. 1c and 1d).

B. BM3D prior

We turn to compare the performance of the two cost functions for the BM3D prior [4], which is based on sparsifying a three-dimensional transformation applied to groups of nearest-neighbor (i.e. similar) patches. This prior is non-convex. In fact, it is also not clear how to precisely formulate its associated $s(\tilde{\mathbf{x}})$. Yet, when implementing proximal algorithms the proximal mapping of $s(\tilde{\mathbf{x}})$ can be replaced with applying the BM3D denoiser as a “black-box”. We use 200 iterations of FISTA to minimize the cost functions with typical step sizes as explained above, and the same eight classical test images.

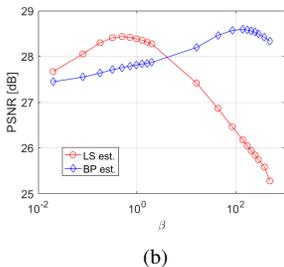
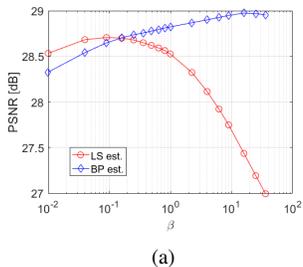


Fig. 17: Super-resolution with Gaussian filter and scale factor of 3, using BM3D prior and 200 iterations of FISTA. PSNR (averaged over 8 test images) vs. β (regularization parameter), for (a) $\sigma_e = 0$, and (b) $\sigma_e = \sqrt{2}$.

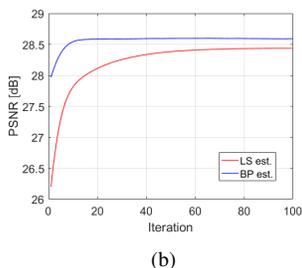
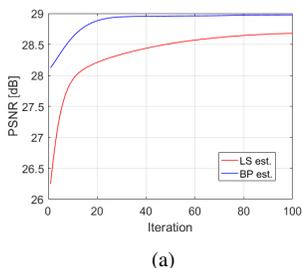


Fig. 18: Super-resolution with Gaussian filter and scale factor of 3, using BM3D prior. PSNR (for best uniform setting of β , averaged over 8 test images) vs. FISTA iteration number, for (a) $\sigma_e = 0$, and (b) $\sigma_e = \sqrt{2}$.

1) *Super-resolution*: We repeat the two SR experiments of Section IV-A1. Fig. 17 shows the average PSNR for different values of β , and Fig. 18 shows the average PSNR as a function of the iteration number, where each estimator uses the best β from Fig. 17 (0.09 for LS and 16 for BP in Fig. 18a; 0.5 for LS and 140 for BP in Fig. 18b). Again, note that \hat{x}_{BP} converges much faster than \hat{x}_{LS} . In Figs. 21e and 21f we display the results for *cameraman* image in the noiseless case.

Note the strong correlation between the obtained results and the observations from Section II, even though the prior is highly non-convex. In the noiseless case, \hat{x}_{BP} outperforms \hat{x}_{LS} for a large range of β . For very small values of β it is inferior to \hat{x}_{LS} , but with only a small gap. From a practitioner point of view, the advantages of using the BP cost here are still clear, since when β is well-tuned (for each of the cost functions) \hat{x}_{BP} is significantly better. Note that in the examined noisy scenario, well-tuned \hat{x}_{BP} is still better than well-tuned \hat{x}_{LS} , but the gap decreases.

2) *Deblurring*: We repeat the two deblurring experiments of Section IV-A2. Fig. 19 shows the average PSNR for different values of β , and Fig. 20 shows the average PSNR as a function of the iteration number, where each estimator uses the best β from Fig. 19 (0.027 for LS and 25.5 for BP in Fig. 20a; 0.5 for LS and 29.5 for BP in Fig. 20b). Figs. 22e and 22f present visual results for *couple* image in the case of $\sigma_e = \sqrt{2}$. The observations that have been made for TV prior stay the same here for the BM3D prior: There exist settings of (β, ϵ) for which \hat{x}_{BP} significantly outperforms \hat{x}_{LS} and converges faster.

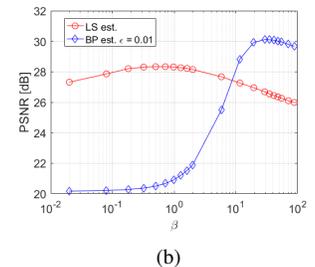
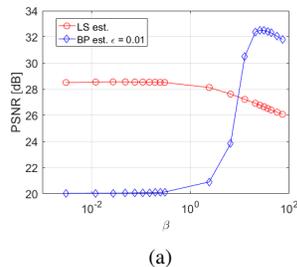


Fig. 19: Deblurring with uniform 9×9 blur kernel, using BM3D prior and 200 iterations of FISTA. PSNR (averaged over 8 test images) vs. β (regularization parameter), for (a) $\sigma_e = \sqrt{0.3}$, and (b) $\sigma_e = \sqrt{2}$.

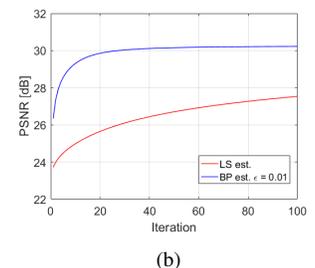
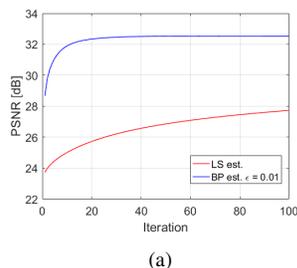


Fig. 20: Deblurring with uniform 9×9 blur kernel, using BM3D prior. PSNR (for best uniform setting of β , averaged over 8 test images) vs. FISTA iteration number, for (a) $\sigma_e = \sqrt{0.3}$, and (b) $\sigma_e = \sqrt{2}$.

C. DCGAN prior

The developments in deep learning [36] in the recent years have led to significant improvement in learning generative models. Methods like variational auto-encoders (VAEs) [37] and generative adversarial networks (GANs) [38] have found success at modeling data distributions. This has naturally led to using pre-trained generative models as priors in imaging inverse problems [17]. Since in popular generative models [37], [38] a generator $\mathcal{G}(\cdot)$ learns a mapping from a low dimensional space $z \in \mathbb{R}^d$ to the signal space $\mathcal{G}(z) \subset \mathbb{R}^n$, one can search for a reconstruction of x only in the range of the generator. This can be formulated by the following non-convex prior

$$s(\tilde{x}) = \begin{cases} 0, & \exists \tilde{z} \in \mathbb{R}^d : \tilde{x} = \mathcal{G}(\tilde{z}) \\ +\infty, & \text{otherwise} \end{cases}. \quad (24)$$

Plugging (24) into the typical cost function (5), we get the objective

$$f_{LS}(\tilde{z}) = \|\mathbf{y} - \mathbf{A}\mathcal{G}(\tilde{z})\|_2^2. \quad (25)$$

Note that for this prior, a regularization parameter β is not required. The recovery of the latent image x is given by $\hat{x}_{LS} = \mathcal{G}(\hat{z}_{LS})$, where \hat{z}_{LS} is a minimizer of (25), which can be obtained by backpropagation and standard gradient based optimizers.

The technique above has been examined recently in [17]. Here, we compare it with the one obtained by a similar approach that uses the BP cost function (6), i.e. we plug (24) into (6), to get the objective

$$f_{BP}(\tilde{z}) = \|\mathbf{A}^\dagger(\mathbf{y} - \mathbf{A}\mathcal{G}(\tilde{z}))\|_2^2, \quad (26)$$



(a) Original image

(b) Bicubic (22.83 dB)



(c) LS-TV (24.37 dB)

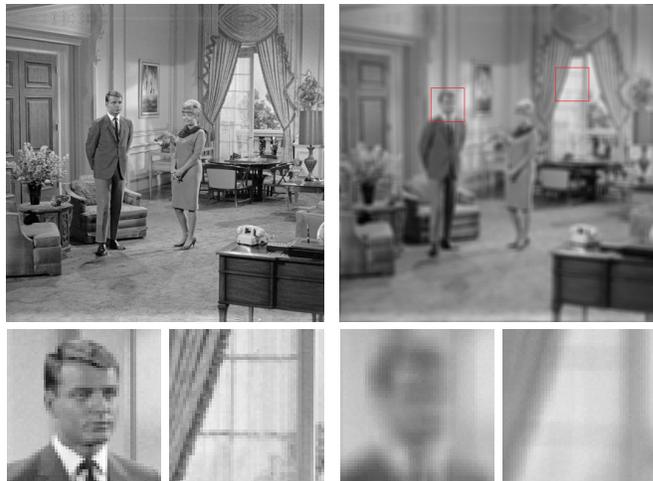
(d) BP-TV (24.94 dB)



(e) LS-BM3D (25.02 dB)

(f) BP-BM3D (25.38 dB)

Fig. 21: Super-resolution with Gaussian filter, scale factor of 3 and $\sigma_e = 0$ for *cameraman* image. From left to right and from top to bottom: original image, bicubic upsampling, reconstruction of LS fidelity with TV prior, reconstruction of BP fidelity with TV prior, reconstruction of LS fidelity with BM3D prior, and reconstruction of BP fidelity with BM3D prior.



(a) Original image

(b) Blurred and noisy image



(c) LS-TV (26.45 dB)

(d) BP-TV (29.28 dB)



(e) LS-BM3D (28.74 dB)

(f) BP-BM3D (30.38 dB)

Fig. 22: Deblurring with uniform 9×9 blur kernel and $\sigma_e = \sqrt{2}$ for *couple* image. From left to right and from top to bottom: original image, blurred and noisy image, reconstruction of LS fidelity with TV prior, reconstruction of BP fidelity with TV prior, reconstruction of LS fidelity with BM3D prior, and reconstruction of BP fidelity with BM3D prior.

TABLE I: Reconstruction PSNR [dB] (averaged over 50 images from CelebA) for super-resolution with Gaussian filter and scale factor of 3, using DCGAN prior and ADAM optimizer.

	Bicubic	LS est.	BP est.
SR $\times 3$	23.04	23.02	23.77

and recover \mathbf{x} by $\hat{\mathbf{x}}_{BP} = \mathcal{G}(\hat{\mathbf{z}}_{BP})$, where $\hat{\mathbf{z}}_{BP}$ is a minimizer of (26).

We use the CelebA dataset [39] and Tensorflow package [40] to train a generator using DCGAN architecture [19] on the cropped version of the images (64×64 pixels), as done in [17]. We use the first 200,000 images (out of 202,599) for training, and the training procedure follows the one in [17], [19]. At test time, all the optimizations with respect to \mathbf{z} are performed using: ADAM [21] with learning rate of 0.1 (as done in [17]), same 10 random initializations of $\tilde{\mathbf{z}}$, and 2000 iterations, which suffice for ensuring that the objectives (25) and (26) stop decreasing. The value of $\tilde{\mathbf{z}}$ that gives the lowest objective is chosen.

1) *Super-resolution*: We compare the performance of $\hat{\mathbf{x}}_{LS}$ and $\hat{\mathbf{x}}_{BP}$ for SR with Gaussian anti-aliasing kernel (defined in Section III-A) and scale factor of 3. Table I shows the PSNR results for the different cost functions, averaged over the last 50 images in CelebA (these images are not included in the training data). Several visual results are shown in Fig. 23.

It can be seen that the BP fidelity yields higher average PSNR and perceptually better recoveries. In fact, in each of the 50 examined images $\hat{\mathbf{x}}_{BP}$ has obtained higher PSNR than $\hat{\mathbf{x}}_{LS}$. This behavior agrees with the previous experiments that demonstrate the advantages of the BP cost for the noiseless SR problem. We also note that even though the results of the simple bicubic upsampling are always perceptually worse than the recoveries that use DCGAN, its PSNR is sometimes higher. This drawback of GAN-based priors is due to the limited representation capabilities of the generators (sometimes referred to as "mode collapse"). A very recent work has suggested to mitigate this deficiency by image-adaptation and back-projections [41].

2) *Compressed sensing*: Due to the small image dimensions, we are able to compare the performance of $\hat{\mathbf{x}}_{BP}$ and $\hat{\mathbf{x}}_{LS}$ for CS with Gaussian measurement matrix (i.e. $A_{ij} \sim \mathcal{N}(0, 1/m)$), for which the two cost functions differ (see the discussion in Section III-B). We use compression ratios of $m/n = 0.1$, $m/n = 0.3$, and $m/n = 0.5$. Table II shows the PSNR results for the different cost functions, averaged over the last 50 images in CelebA. Several visual results are shown in Figs. 24 and 25.

The performance gap between $\hat{\mathbf{x}}_{BP}$ and $\hat{\mathbf{x}}_{LS}$ is negligible for $m/n = 0.1$, and increases in favor of BP when the m/n ratio increases. This behavior correlates with the analysis in Section II (specifically with Observation 3), which explains such behavior for ℓ_2 priors by the fact that when the m/n ratio increases the condition number of $\mathbf{A}\mathbf{A}^T$ increases as well.

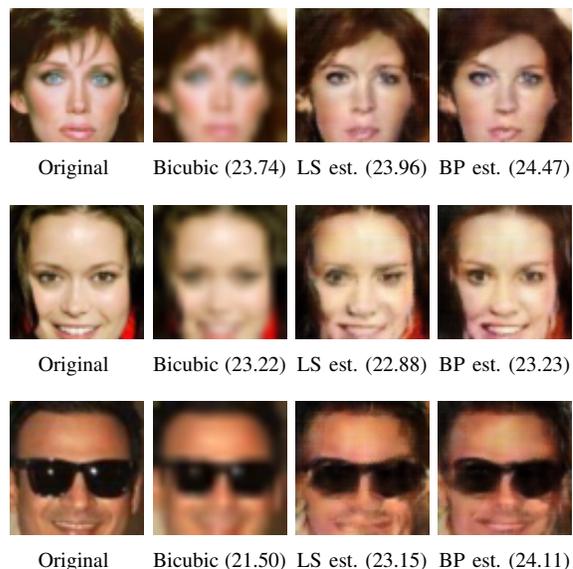


Fig. 23: Super-resolution with Gaussian filter and scale factor of 3, using DCGAN prior.

TABLE II: Reconstruction PSNR [dB] (averaged over 50 images from CelebA) for compressed sensing with Gaussian measurement matrix, using DCGAN prior and ADAM optimizer.

	Naive $\mathbf{A}^\dagger \mathbf{y}$	LS est.	BP est.
CS $m/n = 0.1$	12.07	22.78	22.80
CS $m/n = 0.3$	13.22	23.55	23.62
CS $m/n = 0.5$	14.71	23.67	23.82

V. CONCLUSION

In this work we examined the BP fidelity term for ill-posed linear inverse problems. This term has only been used implicitly by the recently proposed iterative denoising and backward projections (IDBP) framework, and is an alternative to the least squares (LS) term, which is the common choice in most works. We showed that IDBP is essentially a specific optimization scheme, namely the proximal gradient method (known also as ISTA), for minimizing the cost function induced by the BP fidelity term. We analytically compared the two fidelity terms—BP and LS—for the case of ℓ_2 -type prior functions, and obtained mathematically-backed observations in favor of the BP term when the condition number of $\mathbf{A}\mathbf{A}^T$ is large (which is the case in many applications, such as super-resolution and deblurring). Furthermore, we showed that it is possible to leverage prior knowledge on \mathbf{x} to increase the coverage of the observations. Finally, we empirically demonstrated that the behavior for sophisticated priors, such as TV, BM3D and DCGAN, strongly correlates with the theoretically backed properties that we established for ℓ_2 priors. While the mathematical performance analysis in this work is done only for ℓ_2 priors, it provides a good characterization for the advantages of BP and LS compared to each other. Yet, we believe that there are other factors that should be explored with respect to the new fidelity term, such as its behavior with non-convex priors or its effect on the convergence speed

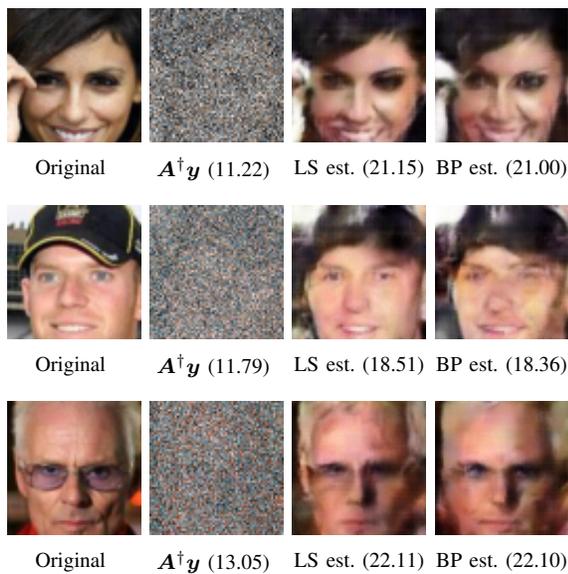


Fig. 24: Compressed sensing with $m = 0.1n$ Gaussian measurements, using DCGAN prior.

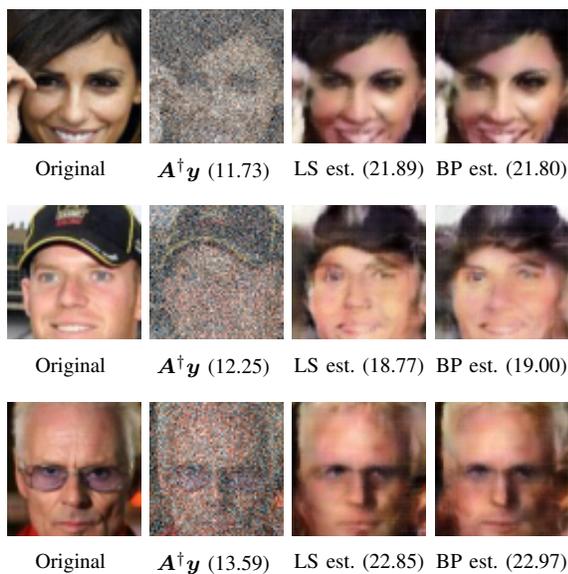


Fig. 25: Compressed sensing with $m = 0.5n$ Gaussian measurements, using DCGAN prior.

of iterative optimization algorithms.

APPENDIX A

THE CONNECTION BETWEEN IDBP [18] AND $f_{BP}(\tilde{\mathbf{x}})$

A. Background

The iterative denoising and backward projections (IDBP) framework [18] is inspired by the plug-and-play priors concept [42], which encourages the usage of existing Gaussian denoisers as “black boxes” to implicitly dictate the prior $s(\tilde{\mathbf{x}})$ when solving inverse problems. Such an approach allows one to use sophisticated denoising methods even when it is not clear how to formulate their associated priors, e.g. convolutional neural network (CNN) denoisers.

Several plug-and-play works have been published [29], [42]–[49]. Most of them consider the typical cost function (5) and directly minimize it using existing iterative optimization schemes, such as FISTA [20], ADMM [50] or quadratic penalty method [51], that include steps in which the proximal mapping of $s(\tilde{\mathbf{x}})$ is used (as explained below, this mapping is equivalent to Gaussian denoising under the prior $s(\tilde{\mathbf{x}})$).

Recently, [18] has suggested, after several manipulations, to solve a different optimization problem

$$\min_{\tilde{\mathbf{x}}, \tilde{\mathbf{z}}} \frac{1}{2(\sigma_e + \delta)^2} \|\tilde{\mathbf{z}} - \tilde{\mathbf{x}}\|_2^2 + s(\tilde{\mathbf{x}}) \quad \text{s.t.} \quad \mathbf{A}\tilde{\mathbf{z}} = \mathbf{y}, \quad (27)$$

where σ_e is the noise level and δ is a design parameter. This work has also proposed an adaptive strategy to set δ , which does not depend on the prior and, contrary to cross-validation, does not require a set of ground truth examples. It has been suggested in [18] to solve (27) using a simple alternating minimization scheme that possesses the plug-and-play property, where the prior term $s(\tilde{\mathbf{x}})$ is handled solely by a Gaussian denoising operation $\mathcal{D}(\cdot; \sigma)$ with noise level $\sigma = \sigma_e + \delta$. In this iterative method, $\tilde{\mathbf{z}}_k$ is obtained by projecting $\tilde{\mathbf{x}}_{k-1}$ onto $\{\mathbf{A}\mathbb{R}^n = \mathbf{y}\}$

$$\begin{aligned} \tilde{\mathbf{z}}_k &= \underset{\tilde{\mathbf{z}}}{\operatorname{argmin}} \|\tilde{\mathbf{z}} - \tilde{\mathbf{x}}_{k-1}\|_2^2 \quad \text{s.t.} \quad \mathbf{A}\tilde{\mathbf{z}} = \mathbf{y} \\ &= \mathbf{A}^\dagger \mathbf{y} + (\mathbf{I}_n - \mathbf{A}^\dagger \mathbf{A}) \tilde{\mathbf{x}}_{k-1} \\ &= \tilde{\mathbf{x}}_{k-1} + \mathbf{A}^\dagger (\mathbf{y} - \mathbf{A} \tilde{\mathbf{x}}_{k-1}). \end{aligned} \quad (28)$$

and $\tilde{\mathbf{x}}_k$ is obtained by

$$\begin{aligned} \tilde{\mathbf{x}}_k &= \underset{\tilde{\mathbf{x}}}{\operatorname{argmin}} \frac{1}{2(\sigma_e + \delta)^2} \|\tilde{\mathbf{z}}_k - \tilde{\mathbf{x}}\|_2^2 + s(\tilde{\mathbf{x}}) \\ &\triangleq \mathcal{D}(\tilde{\mathbf{z}}_k; \sigma_e + \delta). \end{aligned} \quad (29)$$

The two repeating operations lends the method its name: Iterative Denoising and Backward Projections (IDBP). After a stopping criterion is met, the last $\tilde{\mathbf{x}}_k$ is taken as the estimate of the latent \mathbf{x} . Note that in many cases the operation \mathbf{A}^\dagger can be performed efficiently (e.g. the matrix inversion can be avoided using the conjugate gradient method [35]), and thus IDBP is dominated by the complexity of the denoising operation, similarly to other plug-and-play techniques. Using sophisticated denoisers, such as BM3D and CNNs, this algorithm has achieved excellent results for deblurring [18], [22] and super-resolution [23].

B. Obtaining IDBP by applying ISTA on $f_{BP}(\tilde{\mathbf{x}})$

Interestingly, there is another way to develop the exact algorithm, which is different from the way it is developed in [18]. First, note that (27) can be solved directly for $\tilde{\mathbf{z}}$. Similar to (28), we get

$$\tilde{\mathbf{z}}^* = \mathbf{A}^\dagger \mathbf{y} + (\mathbf{I}_n - \mathbf{A}^\dagger \mathbf{A}) \tilde{\mathbf{x}}. \quad (30)$$

Substituting (30) into (27), we reach $\min_{\tilde{\mathbf{x}}} f_{BP}(\tilde{\mathbf{x}})$ with a specific value of the regularization parameter, i.e. $\beta = (\sigma_e + \delta)^2$. Therefore, IDBP is essentially a specific method to minimize the $f_{BP}(\tilde{\mathbf{x}})$ cost function. Let us show that this method coincides with applying the proximal gradient method [20],

[24], popularized under the name ISTA³, on $f_{BP}(\tilde{\mathbf{x}})$. Let us define the proximal mapping, which was introduced by Moreau [53] for convex functions. Here we do not limit this definition to convex functions, though, we emphasize that previous results for proximal mapping of convex functions do not apply to non-convex functions.

Definition 1. *The proximal mapping of a function $s(\cdot)$ at the point $\tilde{\mathbf{z}}$ is defined by*

$$\text{prox}_{s(\cdot)}(\tilde{\mathbf{z}}) \triangleq \underset{\tilde{\mathbf{x}}}{\text{argmin}} \frac{1}{2} \|\tilde{\mathbf{z}} - \tilde{\mathbf{x}}\|_2^2 + s(\tilde{\mathbf{x}}). \quad (31)$$

Clearly, given the same $s(\cdot)$, Gaussian denoising and proximal mapping are tightly connected $\mathcal{D}(\tilde{\mathbf{z}}; \sigma) = \text{prox}_{\sigma^2 s(\cdot)}(\tilde{\mathbf{z}})$.

Assuming a differentiable fidelity term $\ell(\tilde{\mathbf{x}})$ with a Lipschitz continuous gradient $\nabla \ell(\tilde{\mathbf{x}})$, applying ISTA on (2) involves iterations of

$$\tilde{\mathbf{x}}_k = \text{prox}_{\mu\beta s(\cdot)}(\tilde{\mathbf{x}}_{k-1} - \mu \nabla \ell(\tilde{\mathbf{x}}_{k-1})), \quad (32)$$

where μ is a step-size, which ensures convergence for convex $s(\cdot)$ if it is equal to (or smaller than) 1 over the Lipschitz constant of $\nabla \ell(\tilde{\mathbf{x}})$ [20].

Proposition 1. *The IDBP algorithm, given in (28) and (29), coincides with applying ISTA (32) on the cost function $f_{BP}(\tilde{\mathbf{x}})$.*

Proof. Let us compute $\nabla \ell_{BP}(\tilde{\mathbf{x}})$. Using the properties $\mathbf{P}_A \triangleq \mathbf{A}^\dagger \mathbf{A} = \mathbf{P}_A^T = \mathbf{P}_A^2$ and $\mathbf{P}_A \mathbf{A}^\dagger = \mathbf{A}^\dagger$, we get

$$\begin{aligned} \nabla \ell_{BP}(\tilde{\mathbf{x}}) &= -\mathbf{P}_A (\mathbf{A}^\dagger \mathbf{y} - \mathbf{P}_A \tilde{\mathbf{x}}) \\ &= -\mathbf{A}^\dagger (\mathbf{y} - \mathbf{A} \tilde{\mathbf{x}}). \end{aligned} \quad (33)$$

The Lipschitz constant of $\nabla \ell_{BP}(\tilde{\mathbf{x}})$ can be computed here as the spectral norm of the constant Hessian matrix $\nabla^2 \ell_{BP}$. Therefore, μ can be chosen as

$$\mu = \frac{1}{\|\nabla^2 \ell_{BP}(\tilde{\mathbf{x}})\|} = \frac{1}{\|\mathbf{P}_A\|} = 1, \quad (34)$$

where we use the fact that the spectral norm of a non-trivial orthogonal projection is 1. Now, due to the connection $\mathcal{D}(\tilde{\mathbf{z}}; \sigma) = \text{prox}_{\sigma^2 s(\cdot)}(\tilde{\mathbf{z}})$, (32) can be written as

$$\tilde{\mathbf{x}}_k = \mathcal{D}(\tilde{\mathbf{x}}_{k-1} - \mu \nabla \ell(\tilde{\mathbf{x}}_{k-1}); \sqrt{\mu\beta}). \quad (35)$$

Finally, by plugging (33) and (34) into (35) and setting $\beta = (\sigma_e + \delta)^2$, we get the IDBP scheme, which is presented in (28) and (29). \square

The connection between IDBP and ISTA, allows IDBP to adopt the theoretical results of the latter. Yet, note that the powerful global convergence (obtaining the optimal value of the objective) of ISTA holds only for denoisers that are associated with convex prior functions [20]. This limitation is shared also with ADMM-based plug-and-play schemes [43].

³ISTA is the abbreviation of Iterative Shrinkage-Thresholding Algorithm, initially designed for $s(\tilde{\mathbf{x}}) = \|\tilde{\mathbf{x}}\|_1$ [52].

REFERENCES

- [1] L. I. Rudin, S. Osher, and E. Fatemi, "Nonlinear total variation based noise removal algorithms," *Physica D: nonlinear phenomena*, vol. 60, no. 1-4, pp. 259–268, 1992.
- [2] A. Buades, B. Coll, and J.-M. Morel, "A review of image denoising algorithms, with a new one," *Multiscale Modeling & Simulation*, vol. 4, no. 2, pp. 490–530, 2005.
- [3] M. Elad and M. Aharon, "Image denoising via sparse and redundant representations over learned dictionaries," *IEEE Transactions on Image Processing*, vol. 15, no. 12, pp. 3736–3745, 2006.
- [4] K. Dabov, A. Foi, V. Katkovich, and K. Egiazarian, "Image denoising by sparse 3-D transform-domain collaborative filtering," *IEEE Transactions on Image Processing*, vol. 16, no. 8, pp. 2080–2095, 2007.
- [5] M. Bertalmio, G. Sapiro, V. Caselles, and C. Ballester, "Image inpainting," in *Proceedings of the 27th annual conference on Computer graphics and interactive techniques*, pp. 417–424, ACM Press/Addison-Wesley Publishing Co., 2000.
- [6] A. Criminisi, P. Pérez, and K. Toyama, "Region filling and object removal by exemplar-based image inpainting," *IEEE Transactions on Image Processing*, vol. 13, no. 9, pp. 1200–1212, 2004.
- [7] M. Elad, J.-L. Starck, P. Querre, and D. L. Donoho, "Simultaneous cartoon and texture image inpainting using morphological component analysis (MCA)," *Applied and Computational Harmonic Analysis*, vol. 19, no. 3, pp. 340–358, 2005.
- [8] J. A. Guerrero-Colón, L. Mancera, and J. Portilla, "Image restoration using space-variant Gaussian scale mixtures in overcomplete pyramids," *IEEE Transactions on Image Processing*, vol. 17, no. 1, pp. 27–41, 2008.
- [9] A. Danielyan, V. Katkovich, and K. Egiazarian, "BM3D frames and variational image deblurring," *IEEE Transactions on Image Processing*, vol. 21, no. 4, pp. 1715–1728, 2012.
- [10] J. Yang, J. Wright, T. S. Huang, and Y. Ma, "Image super-resolution via sparse representation," *IEEE transactions on image processing*, vol. 19, no. 11, pp. 2861–2873, 2010.
- [11] W. Dong, L. Zhang, G. Shi, and X. Li, "Nonlocally centralized sparse representation for image restoration," *IEEE Transactions on Image Processing*, vol. 22, no. 4, pp. 1620–1630, 2013.
- [12] D. L. Donoho, "Compressed sensing," *IEEE Transactions on information theory*, vol. 52, no. 4, pp. 1289–1306, 2006.
- [13] M. F. Duarte, M. A. Davenport, D. Takhar, J. N. Laska, T. Sun, K. E. Kelly, R. G. Baraniuk, et al., "Single-pixel imaging via compressive sampling," *IEEE Signal Processing Magazine*, vol. 25, no. 2, p. 83, 2008.
- [14] E. J. Candès and M. B. Wakin, "An introduction to compressive sampling," *IEEE signal processing magazine*, vol. 25, no. 2, pp. 21–30, 2008.
- [15] T. Blumensath, "Sampling and reconstructing signals from a union of linear subspaces," *IEEE Transactions on Information Theory*, vol. 57, no. 7, pp. 4660–4671, 2011.
- [16] T. Tirer and R. Giryes, "Generalizing CoSaMP to signals from a union of low dimensional linear subspaces," *Applied and Computational Harmonic Analysis*, 2018.
- [17] A. Bora, A. Jalal, E. Price, and A. G. Dimakis, "Compressed sensing using generative models," in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 537–546, JMLR. org, 2017.
- [18] T. Tirer and R. Giryes, "Image restoration by iterative denoising and backward projections," *IEEE Transactions on Image Processing*, vol. 28, no. 3, pp. 1220–1234, 2019.
- [19] A. Radford, L. Metz, and S. Chintala, "Unsupervised representation learning with deep convolutional generative adversarial networks," *arXiv preprint arXiv:1511.06434*, 2015.
- [20] A. Beck and M. Teboulle, "A fast iterative shrinkage-thresholding algorithm for linear inverse problems," *SIAM journal on imaging sciences*, vol. 2, no. 1, pp. 183–202, 2009.
- [21] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [22] T. Tirer and R. Giryes, "An iterative denoising and backwards projections method and its advantages for blind deblurring," in *2018 25th IEEE International Conference on Image Processing (ICIP)*, pp. 973–977, IEEE, 2018.
- [23] T. Tirer and R. Giryes, "Super-resolution via image-adapted denoising CNNs: Incorporating external and internal learning," *IEEE Signal Processing Letters*, vol. 26, no. 7, pp. 1080–1084, 2019.
- [24] P. L. Combettes and J.-C. Pesquet, "Proximal splitting methods in signal processing," in *Fixed-point algorithms for inverse problems in science and engineering*, pp. 185–212, Springer, 2011.

- [25] R. A. Horn and C. R. Johnson, *Matrix analysis*. Cambridge university press, 2012.
- [26] E. T. Hale, W. Yin, and Y. Zhang, “Fixed-point continuation for ℓ_1 -minimization: Methodology and convergence,” *SIAM Journal on Optimization*, vol. 19, no. 3, pp. 1107–1130, 2008.
- [27] R. Giryes, Y. C. Eldar, A. M. Bronstein, and G. Sapiro, “Tradeoffs between convergence speed and reconstruction accuracy in inverse problems,” *IEEE Transactions on Signal Processing*, vol. 66, no. 7, pp. 1676–1690, 2018.
- [28] C. T. Kelley, *Iterative methods for linear and nonlinear equations*, vol. 16. Siam, 1995.
- [29] K. Zhang, W. Zuo, S. Gu, and L. Zhang, “Learning deep CNN denoiser prior for image restoration,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3929–3938, 2017.
- [30] M. Lustig, D. Donoho, and J. M. Pauly, “Sparse MRI: The application of compressed sensing for rapid mr imaging,” *Magnetic Resonance in Medicine: An Official Journal of the International Society for Magnetic Resonance in Medicine*, vol. 58, no. 6, pp. 1182–1195, 2007.
- [31] J. Kovacevic, A. Chebira, *et al.*, “An introduction to frames,” *Foundations and Trends® in Signal Processing*, vol. 2, no. 1, pp. 1–94, 2008.
- [32] Y. E. Nesterov, “A method for solving the convex programming problem with convergence rate $O(1/k^2)$,” in *Dokl. akad. nauk Sssr*, vol. 269, pp. 543–547, 1983.
- [33] T. Goldstein and S. Osher, “The split Bregman method for L1-regularized problems,” *SIAM journal on imaging sciences*, vol. 2, no. 2, pp. 323–343, 2009.
- [34] A. Beck and M. Teboulle, “Fast gradient-based algorithms for constrained total variation image denoising and deblurring problems,” *IEEE transactions on image processing*, vol. 18, no. 11, pp. 2419–2434, 2009.
- [35] M. R. Hestenes and E. Stiefel, *Methods of conjugate gradients for solving linear systems*, vol. 49. 1952.
- [36] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning*. 2016.
- [37] D. P. Kingma and M. Welling, “Auto-encoding variational bayes,” *arXiv preprint arXiv:1312.6114*, 2013.
- [38] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial nets,” in *Advances in neural information processing systems*, pp. 2672–2680, 2014.
- [39] Z. Liu, P. Luo, X. Wang, and X. Tang, “Deep learning face attributes in the wild,” in *Proceedings of the IEEE international conference on computer vision*, pp. 3730–3738, 2015.
- [40] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, *et al.*, “Tensorflow: A system for large-scale machine learning,” in *12th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 16)*, pp. 265–283, 2016.
- [41] S. Abu Hussein, T. Tirer, and R. Giryes, “Image-adaptive GAN based reconstruction,” *arXiv preprint arXiv:1906.05284*, 2019.
- [42] S. V. Venkatakrisnan, C. A. Bouman, and B. Wohlberg, “Plug-and-play priors for model based reconstruction,” in *Global Conference on Signal and Information Processing (GlobalSIP), 2013 IEEE*, pp. 945–948, IEEE, 2013.
- [43] S. Sreehari, S. V. Venkatakrisnan, B. Wohlberg, G. T. Buzzard, L. F. Drummy, J. P. Simmons, and C. A. Bouman, “Plug-and-play priors for bright field electron tomography and sparse interpolation,” *IEEE Transactions on Computational Imaging*, vol. 2, no. 4, pp. 408–423, 2016.
- [44] Y. Romano, M. Elad, and P. Milanfar, “The little engine that could: Regularization by denoising (RED),” *SIAM Journal on Imaging Sciences*, vol. 10, no. 4, pp. 1804–1844, 2017.
- [45] A. M. Teodoro, J. M. Bioucas-Dias, and M. A. Figueiredo, “Image restoration and reconstruction using variable splitting and class-adapted image priors,” in *Image Processing (ICIP), 2016 IEEE International Conference on*, pp. 3518–3522, IEEE, 2016.
- [46] U. S. Kamilov, H. Mansour, and B. Wohlberg, “A plug-and-play priors approach for solving nonlinear imaging inverse problems,” *IEEE Signal Processing Letters*, vol. 24, no. 12, pp. 1872–1876, 2017.
- [47] S. H. Chan, X. Wang, and O. A. Elgandy, “Plug-and-play ADMM for image restoration: Fixed-point convergence and applications,” *IEEE Transactions on Computational Imaging*, vol. 3, no. 1, pp. 84–98, 2017.
- [48] Y. Sun, B. Wohlberg, and U. S. Kamilov, “An online plug-and-play algorithm for regularized image reconstruction,” *IEEE Transactions on Computational Imaging*, 2019.
- [49] S. Ono, “Primal-dual plug-and-play image restoration,” *IEEE Signal Processing Letters*, vol. 24, no. 8, pp. 1108–1112, 2017.
- [50] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, “Distributed optimization and statistical learning via the alternating direction method of multipliers,” *Foundations and Trends® in Machine Learning*, vol. 3, no. 1, pp. 1–122, 2011.
- [51] J. Nocedal and S. J. Wright, *Numerical Optimization*. Springer, 2006.
- [52] I. Daubechies, M. DeFrise, and C. De Mol, “An iterative thresholding algorithm for linear inverse problems with a sparsity constraint,” *Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences*, vol. 57, no. 11, pp. 1413–1457, 2004.
- [53] J.-J. Moreau, “Proximité et dualité dans un espace hilbertien,” *Bull. Soc. Math. France*, vol. 93, no. 2, pp. 273–299, 1965.