

Sparse image reconstruction for molecular imaging

Michael Ting*, *Member, IEEE*, Raviv Raich, *Member, IEEE*, and Alfred O. Hero III *Fellow, IEEE*

Abstract—The application that motivates this paper is molecular imaging at the atomic level. When discretized at sub-atomic distances, the volume is inherently sparse. Noiseless measurements from an imaging technology can be modeled by convolution of the image with the system point spread function (psf). Such is the case with magnetic resonance force microscopy (MRFM), an emerging technology where imaging of an individual tobacco mosaic virus was recently demonstrated with nanometer resolution. We also consider additive white Gaussian noise (AWGN) in the measurements. Many prior works of sparse estimators have focused on the case when \mathbf{H} has low coherence; however, the system matrix \mathbf{H} in our application is the convolution matrix for the system psf. A typical convolution matrix has high coherence. The paper therefore does not assume a low coherence \mathbf{H} . A discrete-continuous form of the Laplacian and atom at zero (LAZE) p.d.f. used by Johnstone and Silverman is formulated, and two sparse estimators derived by maximizing the joint p.d.f. of the observation and image conditioned on the hyperparameters. A thresholding rule that generalizes the hard and soft thresholding rule appears in the course of the derivation. This so-called hybrid thresholding rule, when used in the iterative thresholding framework, gives rise to the hybrid estimator, a generalization of the lasso. Unbiased estimates of the hyperparameters for the lasso and hybrid estimator are obtained via Stein's unbiased risk estimate (SURE). A numerical study with a Gaussian psf and two sparse images shows that the hybrid estimator outperforms the lasso.

I. INTRODUCTION

The structures of biological molecules like proteins and viri are of interest to the medical community [1]. Existing methods for imaging at the nanometer or even sub-nanometer scale include atomic force microscopy (AFM), electron microscopy (EM), and X-ray crystallography [2], [3]. At the sub-atomic scale, a molecule is naturally a sparse image. That is, the volume imaged consists of mostly space with a few locations occupied by atoms. The application in particular that motivates this paper is MRFM [4], a technology that potentially offers advantages not existent in currently used methods. In particular, MRFM is non-destructive and capable of 3-d imaging. Recently, imaging of a biological sample with nanometer resolution was demonstrated [5]. Given that MRFM and indeed even AFM [6] measures the convolution of the image with a point spread function (psf), a deconvolution must be performed in order to obtain the molecular image. This paper considers the following problem: suppose one observes

a linear transformation of a sparse image corrupted by AWGN. With only knowledge of the linear transformation and noise variance, the goal is to reconstruct the unknown sparse image.

The *system matrix* \mathbf{H} is the linear transformation that, in the case of MRFM, represents convolution with the MRFM psf. Several prior works are only applicable when the system matrix has small pairwise correlation, i.e., low coherence or low collinearity [7]–[10]. Others assume that the columns of \mathbf{H} come from a specific random distribution, e.g., the uniform spherical ensemble (USE), or the uniform random projection ensemble (URPE) [11]. These assumptions are inapplicable when \mathbf{H} represents convolution with the MRFM psf. In general, a convolution matrix for a continuous psf would not have low coherence. Such is the case with MRFM. The coherence of the simulated MRFM psf used in the simulation study section is at least 0.557.

The lasso, the estimator formed by maximizing the penalized likelihood criterion with a l_1 penalty on the image values [12], is known to promote sparsity in the estimate. The Bayesian interpretation of the lasso is the maximum a posteriori (MAP) estimate with an i.i.d. Laplacian p.d.f. on the image values [13]. Consider the following: given M i.i.d. samples of a Laplacian distribution, the expected number of samples equal to 0 is zero. The Laplacian p.d.f. is more convincingly described as a heavy-tailed distribution rather than a sparse distribution. Indeed, when used in a suitable hierarchical model such as in sparse Bayesian learning [14], the Gaussian r.v., not commonly considered as a sparse distribution, results in a sparse estimator. While using a sparse prior is clearly not a necessary condition for formulating a sparse estimator, one wonders if a better sparse estimator can be formed if a sparse prior is used instead.

In [15], the mixture of a Dirac delta and a symmetric, unimodal density with heavy tails is considered; a sparse denoising estimator is then obtained via marginal maximum likelihood (MML). The LAZE distribution is a specific member of the mixture family. Going through the same thought experiment previously mentioned with the LAZE distribution, one obtains an intuitive result: Mw samples equal 0, where w is the weight placed on the Dirac delta. Unlike the Laplacian p.d.f., the LAZE p.d.f. is both heavy-tailed and sparse. Under certain conditions, the estimator achieves the asymptotic minimax risk to within a constant factor [15, Thm. 1]. The lasso estimator can be implemented in an iterative thresholding framework using the soft thresholding rule [16], [17]. Use of a thresholding rule based on the LAZE prior in the iterative thresholding framework can potentially result in better performance.

This paper develops several methods to enable Bayes-optimal nanoscale molecular imaging. In particular, advances are made in these three areas.

This research was partially supported by the Army Research Office under contract W911NF-05-1-0403.

M. Ting is with Seagate Technology, Pittsburgh, PA, 15222. Email: m_ting@ieee.org; Phone: +1 412 519 6946; Fax: +1 412 918 7010.

R. Raich is with the Oregon State University, Corvallis, OR, 97331. Email: raich@eecs.oregonstate.edu; Phone: +1 541 737 9862; Fax: +1 541 737 1300.

A. O. Hero III is with the University of Michigan, Ann Arbor, MI, 48109. Email: hero@eecs.umich.edu; Phone: +1 734 763 0564; Fax: +1 734 763 8041.

- 1) First, we introduce a mixed discrete-continuous LAZE prior for use in the MAP/maximum likelihood (ML) framework. Knowing only that the image is sparse, but lacking any precise information on the sparsity level, selection of the hyperparameters or regularization parameters has to be empirical or data-driven. The sparse image and hyperparameters are jointly estimated by maximizing the joint p.d.f. of the observation and unknown sparse image conditioned on the hyperparameters. Two sparse Bernoulli-Laplacian MAP/ML estimators based on the discrete-continuous LAZE p.d.f. are introduced: MAP1 and MAP2.
- 2) The second contribution of the paper is the introduction of the hybrid estimator, which is formed by exclusively using the *hybrid* thresholding rule in the iterative thresholding framework. The hybrid thresholding rule is a generalization of the soft and hard thresholding rules. In order to apply this to the molecular imaging problem, it is necessary to estimate the hyperparameters in a data-driven fashion.
- 3) Thirdly, SURE is applied to estimate the hyperparameter of lasso and of the hybrid estimator proposed above. The SURE-equipped versions of lasso and hybrid estimator are referred to as lasso-SURE and H-SURE. Our lasso-SURE result is a generalization of the results in [18], [19]. Alternative lasso hyperparameter selection methods exist, e.g., [20]. In [20], however, a prior is placed on the support of the image values that discourages the selection of high correlated columns of \mathbf{H} . Since the \mathbf{H} we consider has columns that are highly correlated, this predisposes a certain amount of separation between the support of the estimated image values, i.e., the sparse image estimate will be resolution limited. A number of other general-purpose techniques exist as well, e.g., cross validation (CV), generalized CV (GCV), MML [21]. Some are, however, more tractable than others. For example, a closed form expression of the marginal likelihood cannot be obtained for the Laplacian prior: approximations have to be made [13].

A simulation study is performed. In the first part, LS, oracular LS, SBL, stagewise orthogonal matching pursuit (StOMP), and the four proposed sparse estimators, are compared. Two image types (one binary-valued and another based on the LAZE p.d.f.) are studied under two signal-to-noise ratio (SNR) conditions (low and high). MAP2 has the best performance in the two low SNR cases. In one of the high SNR cases, H-SURE has the best performance, while in the other, SBL is arguably the best performing method. When the hyperparameters are estimated via SURE, H-SURE is sparser than lasso-SURE and achieves lower l_p error for $p = 0, 1, 2$ as well as lower detection error E_d . In the second part of the numerical study, the performance of the proposed sparse estimators is studied across the range of SNRs between the low and high values considered in the first part. A 3-d reconstruction example is given in the third part, where the LS and lasso-SURE estimator are compared. This serves to demonstrate the applicability of lasso-SURE on a relatively large problem.

The paper is organized into the following sections. First, the sparse image deconvolution problem is formulated in Section II. The algorithms are discussed in Section III: there are three parts to this section. The two MAP/ML estimators based on the discrete-continuous LAZE prior are derived in Section III-A. This is followed by the introduction of the hybrid estimator in Section III-B. Stein's unbiased risk estimate is applied in Section III-C to derive lasso-SURE and H-SURE. Section IV contains a numerical study comparing the proposed algorithms with several existing sparse reconstruction methods. A summary of the work and future directions in Section V concludes the paper.

II. PROBLEM FORMULATION

Consider a 2-d or 3-d image, and denote its vector version by $\underline{\theta} \in \mathbb{R}^M$. In this paper, $\underline{\theta}$ is assumed to be *sparse*, viz., the percentage of non-zero θ_i is small. Suppose that the measurement $\underline{y} \in \mathbb{R}^N$ is given by

$$\underline{y} = \mathbf{H}\underline{\theta} + \underline{w}, \text{ where } \underline{w} \sim \mathcal{N}(\underline{0}, \sigma^2 \mathbf{I}), \quad (1)$$

where $\mathbf{H} \in \mathbb{R}^{N \times M}$ is termed the *system matrix*, and \underline{w} is AWGN. The problem considered can be stated as: given \underline{y} , \mathbf{H} , and $\sigma > 0$, estimate $\underline{\theta}$ knowing that it is sparse. Without loss of generality, one can assume that the columns of \mathbf{H} have unit l_2 norm. In the problem formulation, note that knowledge of the sparseness of $\underline{\theta}$, viz., $\|\underline{\theta}\|_0$, is *not* known a priori.

It should be noted that, while the sparsity considered in (1) is in the natural basis of $\underline{\theta}$, a wavelet basis has been considered in other works, e.g. [19]. It may be possible to re-formulate (1) using some other basis so that the corresponding system matrix has low coherence. This question is beyond the scope of the paper. The emphasis here is on (1) and on sparsity in the natural basis. If \mathbf{H} had full column rank, an equivalent problem formulation is available. Since $(\mathbf{H}'\mathbf{H})$ is invertible, (1) can be re-written as

$$\tilde{\underline{y}} = \underline{\theta} + \tilde{\underline{w}}, \text{ where } \tilde{\underline{w}} \sim \mathcal{N}(\underline{0}, \sigma^2 \mathbf{H}'(\mathbf{H}')^{-1}) \quad (2)$$

where $\tilde{\underline{y}} \triangleq \mathbf{H}'\underline{y}$; $\mathbf{H}' \triangleq (\mathbf{H}'\mathbf{H})^{-1}\mathbf{H}'$ is the pseudoinverse of \mathbf{H} ; and $\tilde{\underline{w}} \triangleq \mathbf{H}'\underline{w}$ is colored Gaussian noise. Deconvolution of $\underline{\theta}$ from \underline{y} in AWGN is therefore equivalent to denoising of $\underline{\theta}$ in colored Gaussian noise. In the special case that \mathbf{H} is orthonormal, $\tilde{\underline{w}}$ is also AWGN.

III. ALGORITHMS

A. Bernoulli-Laplacian MAP/ML sparse estimators

This section considers the case when the discrete-continuous i.i.d. LAZE prior is used for $p(\underline{\theta}|\underline{\zeta})$, with $\underline{\theta}$ and $\underline{\zeta}$ simultaneously estimated via MAP/ML. For the continuous distribution, $\underline{\theta}$, $\underline{\zeta}$ are obtained as the maximizers of the conditional density $p(\underline{y}, \underline{\theta}|\underline{\zeta})$, viz.,

$$\hat{\underline{\theta}}, \hat{\underline{\zeta}} = \underset{\underline{\theta}, \underline{\zeta}}{\operatorname{argmax}} \log p(\underline{y}, \underline{\theta}|\underline{\zeta}) \quad (3)$$

If $\underline{\zeta}$ were constant, $\hat{\underline{\theta}}$ obtained from (3) would be the MAP estimate. If $\underline{\theta}$ were constant, the resulting $\hat{\underline{\zeta}}$ would be the ML estimate. Since these two principles are at work, it cannot be

said that the estimates obtained via (3) are strictly MAP or ML.

Recall that the LAZE p.d.f. is given by

$$p(\theta_i) = (1 - w)\delta(\theta_i) + w\gamma(\theta_i; a), \quad (4)$$

where $\gamma(x; a) = (1/2)ae^{-a|x|}$ is the Laplacian p.d.f. The Dirac delta function is difficult to work with in the context of maximizing the conditional p.d.f. in (3). Consider then a mixed *discrete-continuous* version of (4). Define the random variables $\tilde{\theta}_i$ and I_i such that $\theta_i = I_i\tilde{\theta}_i$, $1 \leq i \leq M$. The r.v.s $\tilde{\theta}_i, I_i$ have the following density:

$$I_i = \begin{cases} 0 & \text{with probability } (1 - w) \\ 1 & \text{with probability } w \end{cases} \quad (5)$$

$$p(\tilde{\theta}_i | I_i) = \begin{cases} g(\tilde{\theta}_i) & I_i = 0 \\ \gamma(\tilde{\theta}_i; a) & I_i = 1 \end{cases}, \quad (6)$$

where $g(\cdot)$ is some p.d.f. that will be specified later on. It is assumed that $(\tilde{\theta}_i, I_i)$ are i.i.d. I_i assumes the role of the Dirac delta: its introduction necessitates use of the auxiliary density g in (6). Instead of (3), consider the optimality criterion

$$\hat{\underline{\theta}}, \hat{\underline{I}}, \hat{\underline{\zeta}} = \underset{\underline{\theta}, \underline{I}, \underline{\zeta}}{\operatorname{argmax}} \log p(\underline{\theta}, \underline{I} | \underline{y}, \underline{\zeta}) \quad (7)$$

Let $\mathcal{I}_1 \triangleq \{i : I_i = 1\}$ and $\mathcal{I}_0 \triangleq \bar{\mathcal{I}}_1 = \{i : I_i = 0\}$. The maximization of (7) is equivalent to the maximization of

$$\begin{aligned} \Psi_{\text{map}} \triangleq & -\frac{\|\mathbf{H}\hat{\underline{\theta}} - \underline{y}\|^2}{2\sigma^2} + (M - \|\hat{\underline{I}}\|_0) \log(1 - w) \\ & + \|\hat{\underline{I}}\|_0 \log w + \sum_{i \in \mathcal{I}_1} \log \left(\frac{1}{2} a e^{-a|\hat{\theta}_i|} \right) + \sum_{i \in \mathcal{I}_0} \log g(\hat{\theta}_i) \end{aligned} \quad (8)$$

We propose to maximize (8) in a block coordinate-wise fashion [22] via Algorithm 1. Note that $\hat{\theta}_i = \hat{\theta}_i \hat{I}_i$. A superscript “(n)” attached to a variable indicates its value in the n th iteration.

Algorithm 1 Block coordinate maximization of MAP criterion Ψ_{map} .

Require: $\hat{\underline{\theta}}^{(0)}, \hat{\underline{I}}^{(0)}, \epsilon > 0$

1: $n \leftarrow 0$

2: **repeat**

3: $n \leftarrow n + 1$

4: $\hat{\underline{\zeta}}^{(n)} \leftarrow \underset{\underline{\zeta}}{\operatorname{argmax}} \Psi_{\text{map}}(\hat{\underline{\theta}}^{(n-1)}, \hat{\underline{I}}^{(n-1)}, \underline{\zeta})$

5: $\hat{\underline{\theta}}^{(n)}, \hat{\underline{I}}^{(n)} \leftarrow \underset{\underline{\theta}, \underline{I}}{\operatorname{argmax}} \Psi_{\text{map}}(\underline{\theta}, \underline{I}, \hat{\underline{\zeta}}^{(n)})$

6: **until** $\|\hat{\underline{\theta}}^{(n)} - \hat{\underline{\theta}}^{(n-1)}\| < \epsilon$

The p.d.f. g arises as an extra degree of freedom due to the introduction of the indicator variables I_i . Consider two cases: first, let $g(x) = \gamma(x; a)$ in (8). This will give rise to the algorithm MAP1. Second, let $g(x)$ be an arbitrary p.d.f. such that: (1) $|g(x)| < \infty$ for all $x \in \mathbb{R}$; (2) $\sup g(x)$ is attained for some $x \in \mathbb{R}$; and (3) $g(x)$ is *independent* of a, w . By selecting g that satisfies these three properties, the algorithm MAP2 is thus obtained.

1) *MAP1*: Let $\Psi_{\text{map1}}(\hat{\underline{\theta}}, \hat{\underline{I}}, \hat{\underline{\zeta}})$ denote the function obtained by setting $g(x) = \gamma(x; a)$. Step (4) of Algorithm 1 is determined by the solution to $\nabla_{\underline{\zeta}} \Psi_{\text{map1}} = 0$. This is solved as

$$\hat{a} = \frac{M}{\|\hat{\underline{\zeta}}\|_1} \quad \text{and} \quad \hat{w} = \frac{\|\hat{\underline{\zeta}}\|_0}{M}. \quad (9)$$

It can be verified that the Hessian $\nabla_{\underline{\zeta}} \nabla_{\underline{\zeta}}^T \Psi_{\text{map1}}$ is negative definite for all $a > 0$ and $0 < w < 1$. Given n samples x_1, \dots, x_n drawn from a Laplacian p.d.f. $\gamma(\cdot; a)$, the ML estimate of a is $\hat{a}_{\text{ML}} = n(\sum_{i=1}^n |x_i|)^{-1}$. The estimate \hat{a} in (9) is therefore the ML estimate of a where all of the $\hat{\theta}_i$ s are used.

The maximization in step (5) of Algorithm 1 can be obtained by applying the EM algorithm [16]. Recall that EM can be applied using $\underline{z} = \underline{\theta} + \alpha \underline{w}_1$ as the *complete data*, where $\underline{w}_1 \sim \mathcal{N}(0, \alpha^2 \mathbf{I})$ and $\alpha \leq \sigma / \|\mathbf{H}\|_2$. Denote by $\hat{\underline{\theta}}^{(n)}, \hat{\underline{z}}^{(n)}$ the estimates in the n th EM iteration. The E-step is the Landweber iteration

$$\hat{\underline{z}}^{(n)} = \hat{\underline{\theta}}^{(n-1)} + \left(\frac{\alpha}{\sigma} \right)^2 \mathbf{H}^T (\underline{y} - \mathbf{H} \hat{\underline{\theta}}^{(n-1)}). \quad (10)$$

Define the *hybrid thresholding rule* as

$$T_{\text{hy}}(x; t_1, t_2) \triangleq (x - \operatorname{sgn}(x)t_2)I(|x| > t_1), \quad (11)$$

where t_1 and t_2 are restricted to $0 \leq t_2 \leq t_1$. See Fig. 1. This is a generalization of the soft and hard thresholding rules. The soft thresholding rule $T_s(x; t) = T_{\text{hy}}(x; t, t)$, and the hard thresholding rule $T_h(x; t) = xI(|x| > t) = T_{\text{hy}}(x; t, 0)$. The

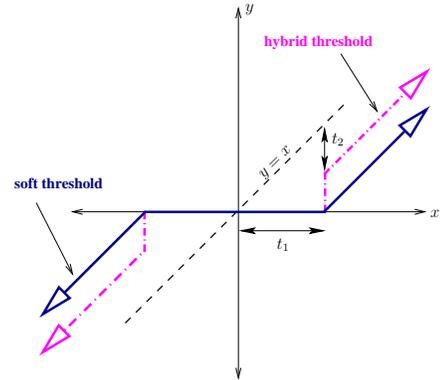


Fig. 1. Hybrid thresholding rule.

M-step of the EM algorithm is given by

$$\theta_i = \begin{cases} T_{\text{hy}}(\hat{z}_i^{(n)}; \alpha a^2 + \kappa_1(\alpha, w), \alpha a^2) & 0 < w \leq \frac{1}{2} \\ T_s(\hat{z}_i^{(n)}; \alpha a^2) & \frac{1}{2} < w \leq 1 \end{cases} \quad (12)$$

where $\kappa_1(\alpha, w) \triangleq \sqrt{2\alpha^2 \log((1-w)/w)}$. Recall that $\theta_i = \tilde{\theta}_i I_i$. If $w > 1/2$, the soft-thresholding rule is applied in the Q -step of the EM iterations of MAP1. These iterations produce the lasso estimate with hyperparameter $\zeta = 2\alpha a^2$. However, if $0 < w \leq 1/2$, a larger thresholding value is used that increases the smaller w becomes.

2) MAP2: From (6) and the assumptions on g , $\tilde{\theta}_i \neq 0$ w.p. 1. Consequently, the set

$$\mathcal{I}_1 = \{i : I_i = 1\} = \{i : \theta_i \neq 0\} \text{ w.p. 1.} \quad (13)$$

This implies $\|\underline{I}\|_0 = \|\underline{\theta}\|_0$ w.p. 1. Apply (13) to the criterion to maximize, viz., (8), and denote the result by $\Psi_{\text{map2}}(\hat{\underline{\theta}}, \underline{I}, \underline{\zeta})$. One gets

$$\begin{aligned} \Psi_{\text{map2}} = & -\frac{1}{2\sigma^2} \|\mathbf{H}\hat{\underline{\theta}} - \underline{y}\|^2 + (M - \|\underline{I}\|_0) \log(1 - w) \\ & + \|\underline{I}\|_0 \log w + \|\underline{\theta}\|_0 \log \frac{a}{2} - a \|\underline{\theta}\|_1 + \sum_{\{i: I_i=0\}} \log g(\tilde{\theta}_i). \end{aligned} \quad (14)$$

The maximization in step (i) is obtained by solving for $\nabla_{\underline{\zeta}} \Psi_{\text{map2}} = 0$, which produces

$$\hat{a} = \frac{\|\hat{\underline{\theta}}\|_0}{\|\hat{\underline{\theta}}\|_1} \text{ and } \hat{w} = \frac{\|\hat{\underline{\theta}}\|_0}{M}. \quad (15)$$

As before, one can verify that the Hessian $\nabla_{\underline{\phi}} \nabla_{\underline{\phi}}^T \Psi_2$ is negative definite for all $a > 0$ and $0 < w < 1$. It is instructive to compare the hyperparameter estimates of MAP1 vs. MAP2, i.e., (9) vs. (15). The main difference lies in the estimation of a . Assuming that the estimates $\hat{\underline{I}}$ and $\hat{\underline{\theta}}$ obey (13), one can rewrite the MAP2 estimate $\hat{a} = |\mathcal{I}_1| / \sum_{i \in \mathcal{I}_1} |\hat{\theta}_i|$. This is the ML estimate using only the $\hat{\theta}_i, i \in \mathcal{I}_1$, i.e., the non-zero voxels. On the other hand, the MAP1 estimate of a can be written as

$$\hat{a} = \frac{|\mathcal{I}_1| + |\mathcal{I}_0|}{\sum_{i \in \mathcal{I}_1} |\hat{\theta}_i| + \sum_{i \in \mathcal{I}_0} |\hat{\theta}_i|}. \quad (16)$$

As with MAP1, the maximization in step (5) of Algorithm 1 can be obtained by applying the EM algorithm with the complete data $\underline{z} = \underline{\theta} + \alpha \underline{w}_1$. The E-step is given by (10), which is the same as MAP1's E-step. Define

$$g^* \triangleq \sup_x g(x) \text{ and } r \triangleq \frac{g^*}{a/2} \frac{1-w}{w}. \quad (17)$$

The resulting $\underline{\theta}$ in the M-step is given by the following thresholding rule

$$\theta_i = \begin{cases} T_{\text{hy}}(\hat{z}_i^{(n)}; a\alpha^2 + \kappa_2(\alpha, r), a\alpha^2) & r \geq 1 \\ T_s(\hat{z}_i^{(n)}; a\alpha^2) & 0 \leq r < 1 \end{cases} \quad (18)$$

where $\kappa_2(\alpha, r) \triangleq \sqrt{2\alpha^2 \log r}$, which is similar to the M-step of MAP1. Indeed, the M-step of MAP1 can be obtained by setting $g^* = a/2$. Just like in MAP1, the EM iterations of MAP2 produce a larger threshold the sparser the hyperparameter w is. As well, if a is smaller, r increases. Since the variance of the Laplacian $\gamma(\cdot; a)$ is $2/a^2$, a smaller a implies a larger variance of the Laplacian. Use of a larger threshold is therefore appropriate.

The tuning parameter g^* can be regarded as an extra degree of freedom that arises due to g being independent of a, w . The MAP2 M-step is a function of g^* , and a suitable value has to be selected. In contrast, MAP1 has no free tuning parameter(s).

B. Hybrid thresholding rule in the iterative framework

Define the *hybrid estimator* to be the estimator formed by using the hybrid thresholding rule (11) in the iterative framework [16, (24)], viz.,

$$\hat{\underline{\theta}}^{(n+1)} = S_{T_{\text{hy}}, \underline{\zeta}} \left(\hat{\underline{\theta}}^{(n)} + (\alpha/\sigma)^2 \mathbf{H}'(\underline{y} - \mathbf{H}\hat{\underline{\theta}}^{(n)}) \right). \quad (19)$$

where $S_{T_{\text{hy}}, \underline{\zeta}}(\underline{x}) = \sum_i T_{\text{hy}}(x_i; \underline{\zeta}) e_i$ and $e_i \in \mathbb{R}^M, i = 1, \dots, M$ are the standard unit vectors. Due to the hybrid thresholding rule being a generalization of the soft thresholding rule, the hybrid estimator potentially offers better performance than lasso. The cost function of the hybrid estimator is given in Prop. 1.

Proposition 1 Consider the iterations (19) when $\|\mathbf{H}\|_2 < 1$ and $\alpha = \sigma$. The iterations minimize the cost function

$$\Psi_{\underline{\zeta}, \text{hy}}(\underline{\theta}) = \|\mathbf{H}\underline{\theta} - \underline{y}\|_2^2 + \sum_i J_1(\theta_i)$$

$$\begin{aligned} \text{where: } J_1(x) = & I(|x| < \zeta_1 - \zeta_2) [-(x - \text{sgn}(x)\zeta_1)^2 + 2\zeta_1\zeta_2] \\ & + I(|x| \geq \zeta_1 - \zeta_2) (2\zeta_2|x| + \zeta_2^2) \end{aligned} \quad (20)$$

Proof. This is an application of Thm. 3 in Appendix I. When $\zeta_1 = \zeta_2 = \zeta$, $J_1(x) = 2\zeta|x| + \zeta^2$, which gives rise to the lasso estimator, as expected. ■

C. Using SURE to empirically estimate the hyperparameters

In this section, SURE is applied to estimate the regularization parameter of lasso and the hybrid estimator. Consider the l_2 risk measure

$$R(\underline{\theta}, \underline{\zeta}) = \frac{1}{N} E_{\underline{Y}} \|\hat{\underline{\theta}} - \underline{\theta}\|_2^2 \quad (21)$$

for lasso. Since $\underline{\theta}$ is not known, this risk cannot be computed; however, one can compute an *unbiased* estimate of the risk [23]. Denote the unbiased estimate by $\hat{R}(\underline{\zeta})$: $\underline{\zeta}$ can then be estimated as $\hat{\underline{\zeta}} = \text{argmin}_{\underline{\zeta} \in \Omega} \hat{R}(\underline{\zeta})$, where Ω is the set of valid values. When $\mathbf{H} = \mathbf{I}$, an expression for $\hat{R}(\underline{\zeta})$ is derived in [18, (11)]. When $\mathbf{H} \neq \mathbf{I}$, however, Stein's unbiased estimate [23] cannot be applied to evaluate (21). In [19], the alternative l_2 risk

$$R(\underline{\theta}, \underline{\zeta}) = \frac{1}{N} E_{\underline{Y}} \|\mathbf{H}(\hat{\underline{\theta}} - \underline{\theta})\|_2^2 \quad (22)$$

is proposed instead. Equation (22) was evaluated for a diagonal \mathbf{H} in [19].

The first theorem in this section generalizes the result of [19] by developing $\hat{R}(\underline{\zeta})$ for arbitrary full column rank \mathbf{H} . The second theorem in this section derives (22) when $\hat{\underline{\theta}}$ is the hybrid estimator. For this result, \mathbf{H} is also an arbitrary full column matrix. If the convolution matrix can be approximated by 2d or 3d circular convolution, the full column rank assumption is equivalent to the 2d or 3d DFT of the psf having no spectral nulls. The proofs of the two theorems are given in Appendix II.

1) SURE for lasso:

Theorem 1 Assume that the columns of \mathbf{H} are linearly independent, and $\hat{\underline{\theta}}$ is the lasso estimator. The unbiased risk estimate (22) is

$$\hat{R}(\zeta) = \sigma^2 + \frac{1}{N} \|\underline{e}\|_2^2 + \frac{2\sigma^2}{N} \|\hat{\underline{\theta}}\|_0 \quad (23)$$

where $\underline{e} = \underline{y} - \mathbf{H}\hat{\underline{\theta}}$ is the reconstruction error.

Since the hyperparameter $\zeta \geq 0$, it can be estimated via

$$\hat{\zeta} = \operatorname{argmin}_{\zeta \geq 0} \hat{R}(\zeta) \quad (24)$$

where $\hat{R}(\zeta)$ is given in (23). LARS can be used to compute (24). Note that LARS requires the linear independence of the columns of \mathbf{H} . The estimator $\hat{\underline{\theta}}_i(\hat{\zeta})$ with $\hat{\zeta}$ obtained via (24) will be referred to as lasso-SURE.

2) SURE for the hybrid estimator: Several definitions are in order first.

Definition 1 Suppose that $\hat{\underline{\theta}} \in \mathbb{R}^M$ has $\|\hat{\underline{\theta}}\|_0 = M - r$. Denote the non-zero components of $\hat{\underline{\theta}}$ by x_i , $1 \leq i \leq M - r$. The permutation matrix $\mathbf{P}(\hat{\underline{\theta}}) \in \mathbf{R}^{M \times M}$ is said to order the zero and non-zero components of $\hat{\underline{\theta}}$ if $\mathbf{P} \operatorname{diag}(\hat{\underline{\theta}}) \mathbf{P}' = \operatorname{diag}(0, \dots, 0, x_1, \dots, x_{M-r})$.

Note that \mathbf{P} in the above definition is not unique. As \mathbf{P} is a permutation matrix, it is orthogonal.

Definition 2 For a matrix $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{q \times r}$, let c_n be a non-zero sequence of length at most q s.t. $1 \leq c_n \leq q$. Similarly, let d_n be non-zero sequence of length at most r s.t. $1 \leq d_n \leq r$. The submatrix $\mathbf{A}[c_n, d_n] = (\alpha_{ij})$ is such that $\alpha_{ij} = a_{c_i, d_j}$.

Define $\Delta_\zeta \triangleq \zeta_1 - \zeta_2$ and

$$\mathbf{U}(\underline{\theta}) \triangleq \begin{cases} \operatorname{diag}[\operatorname{rect}(\frac{\theta_1}{2\Delta_\zeta}), \dots, \operatorname{rect}(\frac{\theta_M}{2\Delta_\zeta})] & \Delta_\zeta > 0 \\ \mathbf{0} & \Delta_\zeta = 0 \end{cases} \quad (25)$$

where $\operatorname{rect}(x) = 1, |x| \leq 1/2$ and 0 otherwise. Recall that $0 \leq \zeta_2 \leq \zeta_1$ by assumption, so $\Delta_\zeta \geq 0$. Let $\mathbf{G}(\mathbf{H}) \triangleq \mathbf{H}'\mathbf{H}$ denote the Gram matrix of \mathbf{H} . For a given $\hat{\underline{\theta}}$, set

$$\mathbf{C}_1(\hat{\underline{\theta}}) \triangleq (\mathbf{P}\mathbf{G}(\mathbf{H})\mathbf{P}') [r+1 : M, r+1 : M] \quad \text{and} \quad (26)$$

$$\mathbf{C}_2(\hat{\underline{\theta}}) \triangleq -\frac{1}{2} (\mathbf{P}\mathbf{U}(\hat{\underline{\theta}})\mathbf{P}') [r+1 : M, r+1 : M], \quad (27)$$

where \mathbf{P} is a matrix that orders the zero and non-zero components of $\hat{\underline{\theta}}$.

Theorem 2 Suppose that the columns of \mathbf{H} are linearly independent and that $\mathbf{G}(\mathbf{H})$ does not have an eigenvalue of $1/2$. With $\hat{\underline{\theta}}$ denoting the hybrid estimator, the unbiased risk estimate (22) is

$$\hat{R}(\underline{\zeta}) = \sigma^2 + \frac{1}{N} \|\underline{e}\|_2^2 + \frac{2\sigma^2}{N} \operatorname{tr}(\mathbf{C}_1[\mathbf{C}_1 + \mathbf{C}_2]^{-1}) \quad (28)$$

where $\underline{e} = \underline{y} - \mathbf{H}\hat{\underline{\theta}}$.

To evaluate (28) for a particular $\hat{\underline{\theta}}$, one would have to construct the matrix \mathbf{P} ; then, invert the $(M - r) \times (M - r)$

matrix $(\mathbf{C}_1 + \mathbf{C}_2)$. If $\hat{\underline{\theta}}$ is sparse, $(M - r)$ is small, and the inversion would not be computationally demanding. The optimum $\underline{\zeta}$ is the $\underline{\zeta} \in \{(\zeta_1, \zeta_2) : \zeta_1 \geq 0, 0 \leq \zeta_2 \leq \zeta_1\}$ that minimizes $\hat{R}(\underline{\zeta})$. The corresponding $\hat{\underline{\theta}}_{\text{hy}}(\underline{\zeta})$ would be the output. This method will be referred to as Hybrid-SURE, or for short, H-SURE.

IV. SIMULATION STUDY

In Section IV-B, the following classes of methods are compared: (i) least-squares (LS) and oracular LS; (ii) the proposed sparse reconstruction methods; and (iii) other existent sparse methods, viz., SBL and StOMP.

The LS solution is implemented via the Landweber algorithm [24]. It provides a ‘‘worst-case’’ bound for the l_2 error, i.e., $\|\underline{e}\|_2$. Since the LS estimate does not take into account the sparsity of $\underline{\theta}$, one would expect it to have worse performance than estimates that do. In the oracular LS method, on the other hand, one knows the support of $\underline{\theta}$, and regresses the measurement \underline{y} on the corresponding columns of \mathbf{H} [25]. The oracular LS estimate consequently provides a ‘‘best-case’’ bound for the l_2 error; however, the oracular LS estimate is unimplementable in reality, as it requires prior knowledge of the support of $\underline{\theta}$. The second class of methods includes the two MAP/ML variants, MAP1 and MAP2; in addition, lasso-SURE and H-SURE are also tested. Finally, in order to benchmark the proposed methods to other sparse methods, SBL and StOMP are included in the simulation study. The Sparselab toolbox is used to obtain the StOMP estimate. The CFAR and CFDR approaches to threshold selection are applied [11]. For CFAR selection, the per-iteration false alarm rate of $1/50$ is used. For CFDR selection, the discovery rate is set to 0.5. Although a multitude of other sparse reconstruction methods exist, they are not included in the simulation study due to a lack of space.

Two sparse images $\underline{\theta}$ are investigated in Section IV-B: a binary-valued image, and an image based on the LAZE prior (4). The binary-valued image has 12 pixels set to one, and the rest are zero. The LAZE image, i.e., the image based on the LAZE prior, can be regarded as a realization of the LAZE prior with $a = 1$ and $w = 0.026$. They are depicted in Fig. 2a,b respectively. The two images are of size 32×32 , as is \underline{y} : so, $M = N = 1024$. The matrix \mathbf{H} , of size 1024×1024 , is the convolution matrix for the Gaussian blur point spread function (psf). In order to satisfy the requirements of Thm. 1 and 2, the columns of \mathbf{H} are linearly independent and $\mathbf{G}(\mathbf{H})$ does not have an eigenvalue of $1/2$. The Gaussian blur is illustrated in Fig. 2c.

The Gaussian blur convolution matrix has columns that are highly correlated: the coherence $\mu = 0.86089$. Let $\Lambda(\underline{\theta}) \triangleq \{i : \theta_i \neq 0\}$. The stability and support results of lasso all require that

$$\mu |\Lambda(\underline{\theta})| \lesssim c \quad (29)$$

where $c = 1/2$ or $1/4$ in order that some statement of recoverability holds [8]–[10], [25]. For a given \mathbf{H} , (29) places an upper bound on $|\Lambda(\underline{\theta})|$ for which recoverability of $\underline{\theta}$ is assured in some fashion. With the Gaussian blur \mathbf{H} , $|\Lambda(\underline{\theta})| \lesssim c/0.86089 < 1$ for both $c = 1/4$ and $1/2$. Since $\|\underline{\theta}\|_0 = 12$,

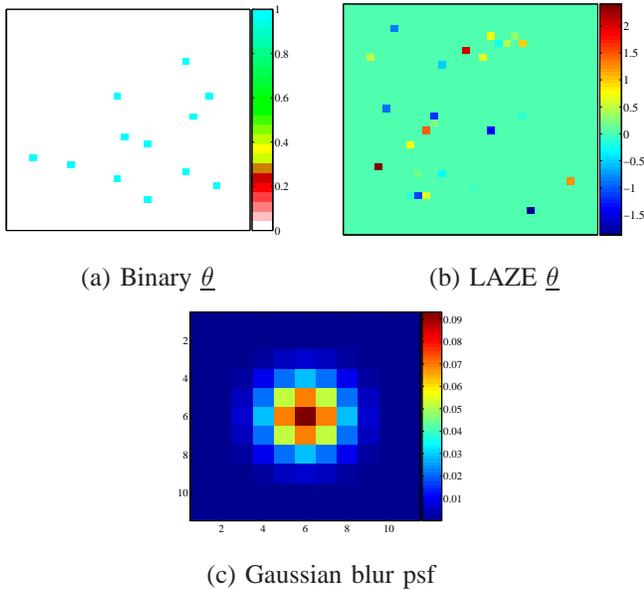


Fig. 2. Illustration of the two types of $\underline{\theta}$ used in the simulations; as well, the Gaussian blur psf is shown.

the simulation study is outside of the coverage of existing recoverability theorems.

In Section IV-C, the performance of the proposed sparse methods over a range of SNRs is investigated. The binary-valued image and Gaussian blur psf are considered in this section. In addition to the proposed sparse methods, the LS estimate is included as a point of reference. Lastly, a 3d MRFM example of dimension $128 \times 128 \times 32$ is given in Section IV-D comparing the LS estimate and lasso-SURE. This serves to illustrate the computational feasibility of lasso-SURE for a relatively large problem.

The proposed algorithms are implemented as previously outlined. The tuning parameter g^* of MAP2 is set to $1/\sqrt{2}$ in Section IV-B and IV-C. LARS is used to compute the lasso-SURE estimator. H-SURE is suboptimally implemented: the minimizing $\underline{\zeta} = (\zeta_1, \zeta_2)'$ is obtained via two line searches. The first, along the $(1/\sqrt{2}, 1/\sqrt{2})$ direction in the (ζ_1, ζ_2) plane, is done using lasso-SURE. A subsequent line search in the $(1, 0)$ direction is performed, i.e., ζ_2 is kept constant and ζ_1 is increased. Define the SNR as $\text{SNR} \triangleq (N^{-1}\|\mathbf{H}\underline{\theta}\|^2)/\sigma^2$, and the SNR in dB as $\text{SNR}_{\text{dB}} \triangleq 10 \log_{10} \text{SNR}$.

A. Error criteria

Recall that the reconstruction error $\underline{e} = \underline{\theta} - \hat{\underline{\theta}}$. Several error criteria are considered in the performance assessment of a sparse estimator.

- $\|\underline{e}\|_p$ for $p = 0, 1, 2$.
- The *detection error* criterion defined by

$$E_d(\underline{\theta}, \hat{\underline{\theta}}; \delta) \triangleq \sum_{i=1}^M |I(\theta_i = 0) - I(|\hat{\theta}_i| < \delta)| \quad (30)$$

Values of $\hat{\theta}_i$ such that $|\hat{\theta}_i| < \delta$ are considered equivalent to 0. This is used to handle the effect of finite-precision computing. More importantly, it addresses the fact that,

to the human observer, small non-zero values are not discernible from zero values. In the study, $\delta = 10^{-2}\|\underline{\theta}\|_\infty$ is selected. This error criterion is effectively a 0-1 penalty on the support of $\underline{\theta}$. Accurately determining the support of a sparse $\underline{\theta}$ is more critical than its actual values [7], [26].

- The number of non-zero values of $\hat{\underline{\theta}}$, i.e., $\|\hat{\underline{\theta}}\|_0$. One would like $\|\hat{\underline{\theta}}\|_0 \approx \|\underline{\theta}\|_0$, which is small if $\underline{\theta}$ is indeed sparse.

B. Performance under low and high SNR

The performance of the estimators is given in Table I for the binary-valued $\underline{\theta}$ with the SNR equal to 1.76 dB (low SNR) and 20 dB (high SNR). The number reported in Table I is the mean over the simulation runs. For each performance criterion, the best mean number is underlined. The oracular LS estimate is excluded from this assessment, as it cannot be implemented without prior knowledge. In terms of $\|\hat{\underline{\theta}}\|_0$, the best number is the value closest to $\|\underline{\theta}\|_0$. Recall that for the binary-valued image $\underline{\theta}$, $\|\underline{\theta}\|_0 = 12$. The best number for the other performance criterion is the value closest to 0.

TABLE I
PERFORMANCE OF THE RECONSTRUCTION METHODS FOR THE BINARY-VALUED $\underline{\theta}$.

Method	Error criterion				
	$\ \underline{e}\ _0$	$\ \underline{e}\ _1$	$\ \underline{e}\ _2$	$E_d(\underline{\theta}, \hat{\underline{\theta}})$	$\ \hat{\underline{\theta}}\ _0$
SNR = 1.76 dB					
Oracular LS	12	0.880	0.309	0	12
LS	1024	579	22.6	1.00×10^3	1024
SBL	1024	13.8	2.35	58.7	1024
StOMP ^{CFAR}	335	1.46×10^4	1.50×10^3	322	335
StOMP ^{CFDR}	454	7.95×10^4	7.17×10^3	442	454
MAP1	<u>12</u>	12	3.46	12	0
MAP2	15.5	<u>2.72</u>	<u>0.912</u>	<u>3.68</u>	<u>15.3</u>
lasso-SURE	60	7.83	1.51	44.2	60.6
H-SURE	39.3	7.25	1.51	27.0	39.3
SNR = 20 dB					
Oracular LS	12	0.112	0.0394	0	12
LS	1024	86.1	3.67	929	1024
SBL	1024	1.19	0.184	32.2	1024
StOMP ^{CFAR}	377	4.33×10^3	457	361	377
StOMP ^{CFDR}	459	1.36×10^4	1.19×10^3	446	459
MAP1	43.9	1.07	0.209	22.9	43.9
MAP2	230	3.82	0.380	114	230
lasso-SURE	61.2	0.923	0.176	15.7	61.8
H-SURE	<u>22.0</u>	<u>0.584</u>	<u>0.152</u>	<u>7.5</u>	<u>22.0</u>

In the low SNR case, MAP2 has the best performance. MAP1 consistently produces the *trivial* estimate of all zeros, as evidenced by the mean value of $\|\hat{\underline{\theta}}\|_0$ being equal to 0. The trivial all-zero estimate results in $\|\underline{e}\|_p = \|\underline{\theta}\|_p$ for $p = 0, 1, 2$. For a sparse $\underline{\theta}$, a small $\|\underline{e}\|_0$ therefore is not necessarily an indicator of good performance. A second comment regarding $\|\hat{\underline{\theta}}\|_0$ is that it does not always give an accurate assessment of the *perceived* sparsity of the reconstruction. In Table I, SBL never produces a strictly sparse estimate, as the mean

$\|\hat{\theta}\|_0$ equals the maximal value of 1024. However, consider Fig. 3a, where the SBL estimate for one noise realization at an SNR of 1.76 dB is depicted. The $\hat{\theta}$ looks sparser than would be suggested by $\|\hat{\theta}\|_0 = 1024$. This is because many of the non-zero pixel values have a small magnitude, and are visually indistinguishable from zero. The SBL estimate has many spurious non-zero pixels, in addition to blurring around several non-zero pixel locations. Negative values are present in the reconstruction, although the binary θ is non-negative.

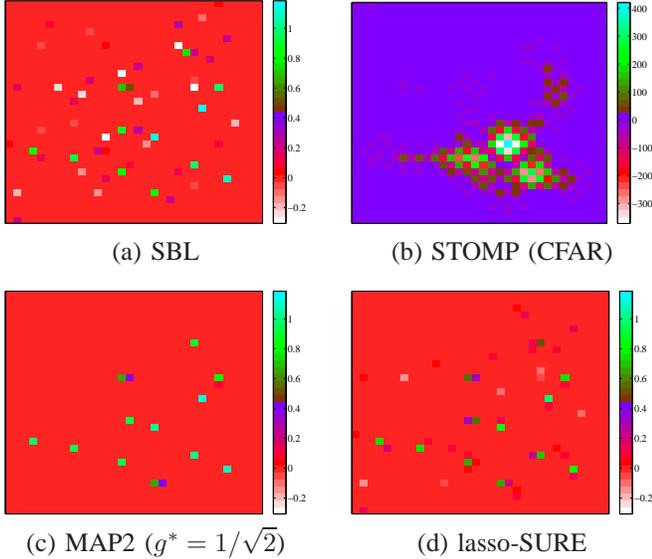


Fig. 3. Reconstructed images for the binary-valued θ under an SNR of 1.76 dB for SBL, StOMP (CFAR), MAP2 ($g^* = 1/\sqrt{2}$), and lasso-SURE.

The StOMP^{CFAR}, MAP2, and lasso-SURE estimate are illustrated in Figs. 3b–d respectively. The StOMP^{CFAR} $\hat{\theta}$ has large positive and negative values. It does not seem like a sufficient number of stages have been taken. While blurring around several non-zero voxels are evident in the MAP2 estimate, $\hat{\theta}$ closely resembles θ , cf. Fig. 2a. None of the estimators considered here take into account positivity. From Fig. 3b, however, one sees that the MAP2 estimate has no negative values. Qualitatively, the lasso-SURE estimate looks better than SBL, but worse than MAP2. This is reflected in the quantitative performance criteria in Table I.

In the high SNR case, H-SURE has the best performance. The mean values of all the performance criteria decrease as compared to lasso-SURE. The greatest decreases are in $\|\underline{e}\|_0$, E_d , and $\|\hat{\theta}\|_0$. They indicate that the H-SURE estimator is properly zeroing out spurious non-zero values and producing a sparser estimate than lasso-SURE. However, this comes at a price of higher computational complexity.

Examine next the performance of the reconstruction methods with the LAZE image. One expects MAP1 and MAP2 to have better performance than the other methods, as the image θ is generated using the LAZE prior. The numbers for the performance criteria are given in Table II. Again, the reconstruction method with the best number for each criterion is underlined. For the LAZE θ , $\|\theta\|_0 = 27$.

In the low SNR case, MAP2 has the advantage. MAP1 produces the trivial estimate of all zeros, just as in the

TABLE II
PERFORMANCE OF THE RECONSTRUCTION METHODS FOR THE LAZE θ .

Method	Error criterion				
	$\ \underline{e}\ _0$	$\ \underline{e}\ _1$	$\ \underline{e}\ _2$	$E_d(\theta, \hat{\theta})$	$\ \hat{\theta}\ _0$
SNR = 1.76 dB					
Oracular LS	27	5.71	1.55	0.56	27
LS	1024	807	31.6	977	1024
SBL	1024	28.1	3.99	72.6	1024
StOMP ^{CFAR}	264	4.37×10^3	558	244	257
StOMP ^{CFDR}	409	1.62×10^4	1.65×10^3	386	405
MAP1	<u>27</u>	21.2	5.21	27	0
MAP2	30.9	<u>17.5</u>	3.98	<u>25.1</u>	<u>9.77</u>
LASSO-SURE	92.6	20.3	3.15	69.3	81.9
H-SURE	67.2	19.1	<u>3.14</u>	51.1	54.7
SNR = 20 dB					
Oracular LS	27	0.686	0.190	0.6	27
LS	1024	122	5.34	856	1024
SBL	1024	<u>4.32</u>	<u>0.814</u>	33.7	1024
StOMP ^{CFAR}	336	1.00×10^4	110	305	330
StOMP ^{CFDR}	438	2.67×10^4	250	408	435
MAP1	<u>69.7</u>	6.53	1.34	31.9	<u>63.8</u>
MAP2	216	10.8	1.44	86.6	212
LASSO-SURE	119	6.63	1.32	<u>31.1</u>	116
H-SURE	84.4	6.73	1.35	33.0	78.7

case of the binary-valued θ . The high SNR case has mixed results. While SBL has the best mean $\|\underline{e}\|_1$ and $\|\underline{e}\|_2$, the best result for the other three criteria each occur at a different method. The fact that MAP1 and MAP2 did not produce superior performance over the other methods in the case of the LAZE image is unintuitive. As the SNR increases, however, the hyperparameter estimates become biased [27]. The other unintuitive result is that $E_d(\theta, \hat{\theta}; \delta)$ for the oracular LS estimate is not zero. This arises because of the choice of δ . Since $\delta = 10^{-2} \|\theta\|_\infty$, the values of $\hat{\theta}_i$ that are smaller than δ in absolute value are thresholded to zero. This results in a non-zero E_d in some cases.

C. Performance vs. SNR of the proposed reconstruction methods

The performance of the proposed reconstruction methods when applied to the binary-valued θ is examined with respect to SNR. The intent in this subsection is to study the behavior of the proposed methods at SNR values in between the low and high values of 1.76 dB and 20 dB respectively. As with the previous section, the MAP2 estimator is used with $g^* = 1/\sqrt{2}$. For each estimator, the mean is plotted along with error bars of one standard deviation. The error plots are given in Fig. 4. Note that in Fig. 4e, the MAP1 curve is missing the first several SNR values because $\|\hat{\theta}\|_0 = 0$ and the y-axis is in a log scale.

First, consider the $\|\underline{e}\|_0$, $\|\underline{e}\|_1$, and $\|\underline{e}\|_2$ error criteria. MAP1 is unable to distinguish the location of the non-zero pixels in low SNR. Under high SNR conditions, it has performance that is comparable to lasso-SURE and H-SURE in terms of the $\|\underline{e}\|_1$ and $\|\underline{e}\|_2$ errors. The value of $\|\underline{e}\|_0$ increases with

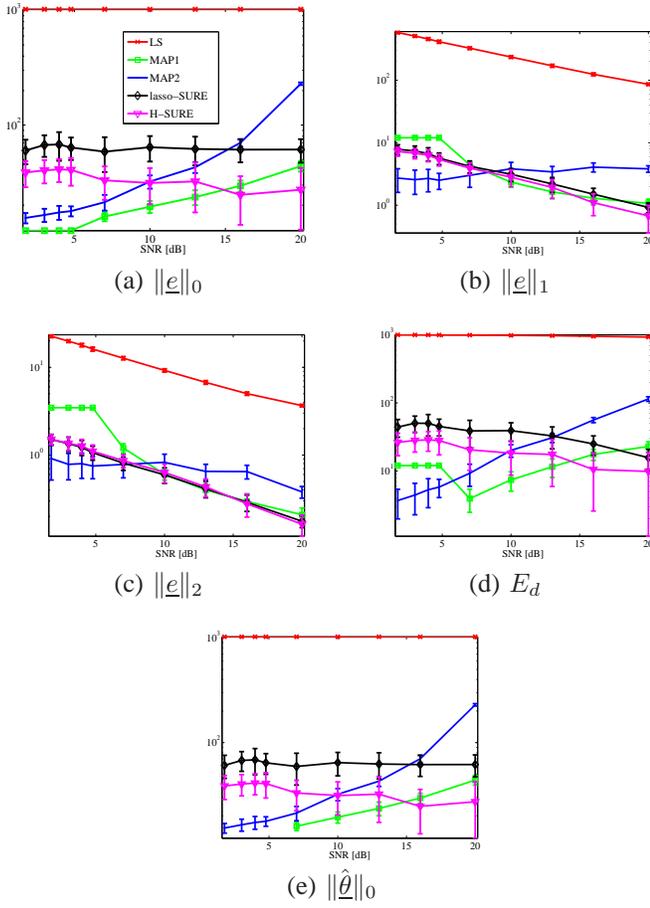


Fig. 4. Performance vs. SNR for Landweber iterations, MAP1, MAP2, lasso-SURE, and H-SURE when applied to the binary-valued $\underline{\theta}$.

respect to increasing SNR for MAP1. Taken together with the $\|\underline{e}\|_1$ and $\|\underline{e}\|_2$ curves, the trend is indicative of small non-zero coefficients appearing in $\hat{\underline{\theta}}$ that are spurious. MAP2 also has the same behavior with respect to $\|\underline{e}\|_0$; however, a performance gap under high SNR exists in its $\|\underline{e}\|_1$ and $\|\underline{e}\|_2$ curves as compared to MAP1, lasso-SURE, and H-SURE. The lasso-SURE and H-SURE estimates have curves that decrease as the SNR increases. H-SURE’s error curve is lower than lasso-SURE’s for $\|\underline{e}\|_0$ and $\|\underline{e}\|_1$, and it is almost identical for $\|\underline{e}\|_2$.

Consider next the E_d and $\|\hat{\underline{\theta}}\|_0$ error criterion. The lasso-SURE curve for $\|\hat{\underline{\theta}}\|_0$ is relatively flat, and its E_d curve decreases for high SNR. This indicates that, while the number of non-zero coefficients in $\hat{\underline{\theta}}$ remains the same, the amplitude at the spurious locations are decreasing. With MAP1 and MAP2, the opposite trend is true. For low SNR, the number of non-zero coefficients in $\hat{\underline{\theta}}$ is small, but increases with higher SNR. A similar increase can be seen in the E_d curves. One can conclude that the number of spurious non-zero locations is increasing. This phenomenon is due to the bias of the hyperparameter estimates [27]. With H-SURE, both the E_d and $\|\hat{\underline{\theta}}\|_0$ curves decrease as the SNR increases. This behavior is intuitive, as higher SNR should result in better performance. We note that H-SURE’s E_d curve is lower than lasso-SURE’s; moreover, H-SURE’s $\|\hat{\underline{\theta}}\|_0$ curve is closer to $\|\underline{\theta}\|_0 = 12$ than

lasso-SURE’s.

D. MRFM reconstruction example

A three dimensional example using the hydrogen atom locations of the DNA molecule (PDB ID: 103D) [28] as $\underline{\theta}$ and the 3d MRFM psf is carried out in this subsection. Both $\underline{\theta}$ and \underline{y} have dimension $128 \times 128 \times 32$, and the SNR is 4.77 dB. Each hydrogen location in $\underline{\theta}$ is set to 1, and the rest of the locations set to 0. The resulting image $\underline{\theta}$ has a helical structure: see Fig. 5a. The image represented by $\mathbf{H}\underline{\theta}$ is illustrated in Fig. 5b. The LS and lasso-SURE estimates

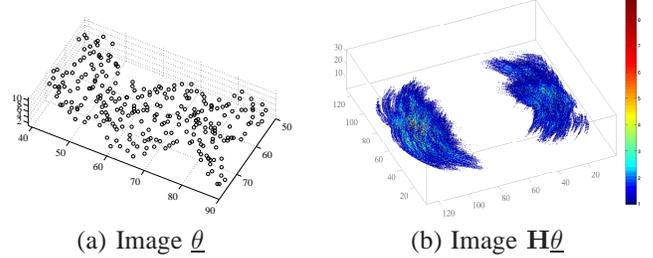


Fig. 5. Image $\underline{\theta}$ and noiseless projection $\mathbf{H}\underline{\theta}$ used in the MRFM reconstruction example.

are given in Fig. 6 and 7 respectively. The 3d figures plot

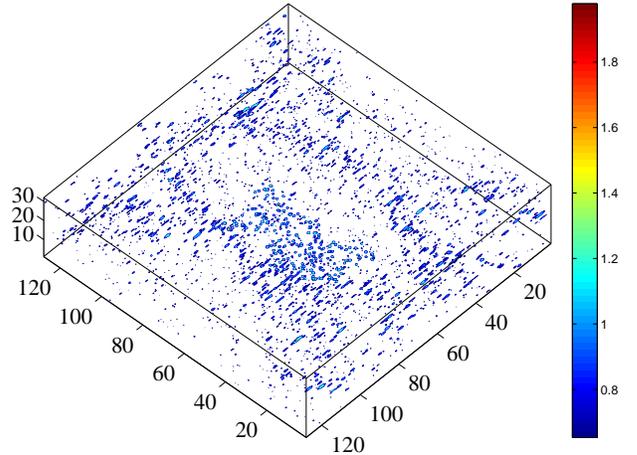


Fig. 6. The LS estimate of the MRFM example under a SNR of 4.77 dB.

contours for several values. The white volume in Fig. 6 does not indicate $\hat{\theta}_i = 0$; rather, the $\hat{\theta}_i$ are at a value smaller than the lowest color bar value. On the other hand, the white volume of the lasso-SURE estimate is mostly $\hat{\theta}_i = 0$. The histogram of $\hat{\theta}_i$ for the LS and lasso-SURE estimator given in Fig. 8a,b respectively illustrate this point. The sharp peak at 0 in the lasso-SURE histogram suggests that the lasso estimator incorporates a thresholding rule, which it does. The $\hat{\theta}_i$ values are separated into two distinct sets: the sparse image centered around 0.95 and the background around 0. In contrast, the histogram of $\hat{\theta}_i$ for Landweber is not separated in this fashion, nor does it have a sharp peak at 0.

V. SUMMARY AND FUTURE DIRECTIONS

Use of a mixed discrete-continuous LAZE prior and jointly estimating $(\underline{\theta}, \underline{\zeta})$ as the maximizer of $p(\underline{y}, \underline{\theta} | \underline{\zeta})$ gives rise to the

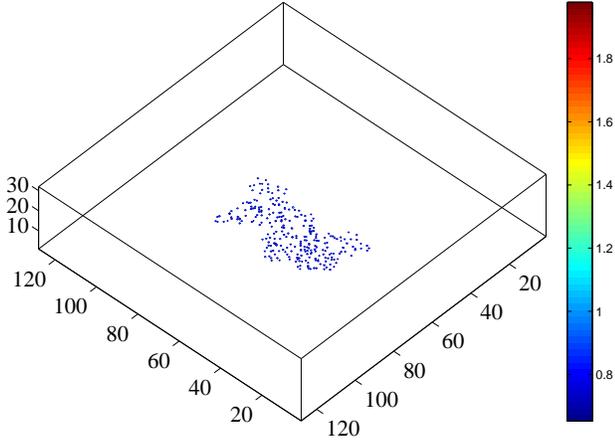


Fig. 7. The lasso-SURE estimate of the MRFM example under a SNR of 4.77 dB.

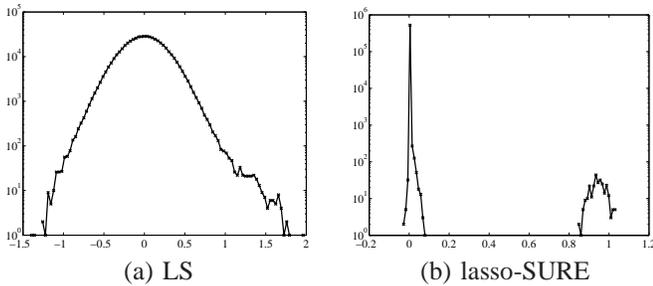


Fig. 8. Histogram of $\hat{\theta}_i$ for the LS and lasso-SURE estimator.

Bernoulli-Laplacian sparse estimators MAP1 and MAP2. The hybrid thresholding rule is observed in both of these sparse estimators. When used in the iterative thresholding framework, the resulting penalty on $\underline{\theta}$ is quadratic around the origin, and linear away from the origin, cf. (20). In order to apply lasso and the hybrid estimator to data, an empirical means of estimating the hyperparameters is required. This is achieved via Stein’s unbiased risk estimate.

A numerical study shows that MAP1 and MAP2 perform well at low SNR, but the performance deteriorates at higher SNR. While StOMP demonstrates competitive results in [11], such is not the case in the simulation study conducted in this paper. The SBL estimate is not sparse; despite this, the estimates look visually sparse due to many non-zero values being small. In the high SNR regime for the LAZE $\underline{\theta}$, SBL has good performance. When the hyperparameters are estimated via SURE, the hybrid estimator achieves a sparser estimate with lower l_p reconstruction error for $p = 0, 1, 2$ as compared to lasso. In addition, the hybrid estimator has lower detection error E_d . The numerical study suggests that sparse estimators based on sparse priors may achieve superior performance to the lasso.

The paper did not compare the MAP/ML and SURE estimates of the hyperparameters to other estimates, e.g., GCV, the method of [20] for lasso, etc. This is primarily due to a lack of space. In the case when $\underline{\theta}$ is a linear function of \underline{y} , SURE is equivalent to the C_p statistic, while GCV is the C_p statistic with σ^2 replaced by an estimated version [29]. Unfortunately,

the sparse estimators considered in the paper are all nonlinear in \underline{y} . Another issue that should be looked in future work is how to improve MAP1/2 to rectify the deteriorating performance at higher SNR. The estimates $\hat{\mathbf{a}}, \hat{\mathbf{w}}$ generally become more biased as the SNR increases [27]. This has been noted in [30]. With MAP2, the degree of bias is affected by the selection of g^* .

Implementation considerations were not discussed, although they are critical in the implementation of a deconvolution algorithm. The interested reader is referred to [27]. In terms of increasing complexity, the estimators can be approximately ordered as: StOMP, LS/oracular LS, MAP1 and MAP2, lasso-SURE, H-SURE, and SBL. Thanks to LARS, evaluating a goodness-of-fit criterion for lasso whether it be a SURE criterion, a GCV criterion, etc. has low computational complexity. Although LARS requires the selection of individual columns of \mathbf{H} , this is not an issue when \mathbf{H} represents a convolution operator. The selection can be efficiently implemented using the fast Fourier transform (FFT). Solving for the H-SURE hyperparameters has higher computational complexity since an efficient implementation of the H-SURE estimator is currently lacking. In this paper, the iterative thresholding framework is used for part of the solution; however, a LARS-like method would be a welcomed improvement.

REFERENCES

- [1] Y. Modis, S. Ogata, D. Clements, and S. C. Harrison, “Structure of the dengue virus envelope protein after membrane fusion,” *Nature*, vol. 427, pp. 313–319, 2004.
- [2] D. J. Müller, D. Fotiadis, S. Scheuring, S. A. Müller, and A. Engel, “Electrostatically balanced subnanometer imaging of biological specimens by atomic force microscope,” *Biophysical Journal*, vol. 76, pp. 1101–1111, 1999.
- [3] Y. G. Kuznetsov, A. J. Malkin, R. W. Lucas, M. Plomp, and A. McPherson, “Imaging of virus by atomic force microscopy,” *Journal of General Virology*, vol. 82, pp. 2025–2034, 2001.
- [4] D. Rugar, R. Budakian, H. J. Mamin, and B. W. Chui, “Single spin detection by magnetic resonance force microscopy,” *Nature*, vol. 430, no. 6997, pp. 329–332, 2004.
- [5] C. L. Degen, M. Poggio, H. J. Mamin, C. T. Rettner, and D. Rugar, “Magnetic resonance imaging of a biological sample with nanometer resolution,” *Science*, submitted.
- [6] P. Markiewicz and M. C. Goh, “Atomic force microscope tip deconvolution using calibration arrays,” *Rev. Sci. Instrum.*, vol. 66, pp. 3186–3190, 1995.
- [7] J. A. Tropp, “Greed is good: algorithmic results for sparse approximation,” *IEEE Trans. Inform. Theory*, vol. 50, no. 10, pp. 2231–2241, 2004.
- [8] J. J. Fuchs, “Recovery of exact sparse representations in the presence of bounded noise,” *IEEE Trans. Inform. Theory*, vol. 51, no. 10, pp. 3601–3608, 2005.
- [9] D. L. Donoho, M. Elad, and V. N. Temlyakov, “Stable recovery of sparse overcomplete representations in the presence of noise,” *IEEE Trans. Inform. Theory*, vol. 52, no. 1, pp. 6–18, 2006.
- [10] J. A. Tropp, “Just Relax: convex programming methods for identifying sparse signals in noise,” *IEEE Trans. Inform. Theory*, vol. 52, no. 3, pp. 1030–1051, 2006.
- [11] D. L. Donoho, Y. Tsaig, I. Drori, and J.-L. Starck, “Sparse solution of underdetermined linear equations by stagewise orthogonal matching pursuit,” Stanford University, Tech. Rep., 2006.
- [12] R. Tibshirani, “Regression shrinkage and selection via the lasso,” *Journal of the Royal Statistical Society, Series B*, vol. 58, no. 1, pp. 267–288, 1996.
- [13] S. Alliney and S. A. Ruzinsky, “An algorithm for the minimization of mixed l_1 and l_2 norms with application to Bayesian estimation,” *IEEE Trans. Signal Processing*, vol. 42, no. 3, pp. 618–627, 1994.
- [14] D. P. Wipf and B. D. Rao, “Sparse Bayesian learning for basis selection,” *IEEE Trans. Signal Processing*, vol. 52, no. 8, pp. 2153–2164, 2004.

- [15] I. M. Johnstone and B. W. Silverman, "Needles and straw in haystacks: empirical Bayes estimates of possibly sparse sequences," *The Annals of Statistics*, vol. 32, no. 4, pp. 1594–1649, 2004.
- [16] M. A. T. Figueiredo and R. D. Nowak, "An EM Algorithm for Wavelet-Based Image Restoration," *IEEE Trans. Image Processing*, vol. 12, no. 8, pp. 906–916, 2003.
- [17] 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*, vol. 57, no. 11, pp. 1413–1457, 2004.
- [18] D. L. Donoho and I. M. Johnstone, "Adapting to unknown smoothness via wavelet shrinkage," *Journal of the American Statistical Association*, vol. 90, no. 423, pp. 1200–1224, 1995.
- [19] L. Ng and V. Solo, "Optical flow estimation using adaptive wavelet zeroing," in *Proceedings of the IEEE Intl. Conf. on Image Processing*, vol. 3, 1999, pp. 722–726.
- [20] M. Yuan and Y. Lin, "Efficient empirical Bayes variable selection and estimation in linear models," *Journal of the American Statistical Association*, vol. 100, no. 472, pp. 1215–1225, 2005.
- [21] A. M. Thompson, J. C. Brown, J. W. Kay, and D. M. Titterington, "A study of methods of choosing the smoothing parameter in image restoration by regularization," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 13, no. 4, pp. 326–339, 1991.
- [22] J. A. Fessler, "Image Reconstruction: Algorithms and Analysis," draft of book.
- [23] C. M. Stein, "Estimation of the mean of a multivariate normal distribution," *The Annals of Statistics*, vol. 9, no. 6, pp. 1135–1151, 1981.
- [24] C. Byrne, "A unified treatment of some iterative algorithms in signal processing and image reconstruction," *Inverse Problems*, vol. 20, no. 1, pp. 103–120, 2004.
- [25] E. Candes, J. Romberg, and T. Tao, "Stable signal recovery from incomplete and inaccurate measurements," *Comm. Pure Appl. Math.*, vol. 59, pp. 1207–1223, 2005.
- [26] K. K. Herrity, A. C. Gilbert, and J. A. Tropp, "Sparse approximation via iterative thresholding," in *Proceedings of the IEEE Intl. Conf. on Acoustics, Speech, and Signal Processing*, 2006.
- [27] M. Y. J. Ting, "Signal processing for magnetic resonance force microscopy," Ph.D. dissertation, The University of Michigan, 2006.
- [28] S. H. Chou, L. Zhu, and B. R. Redi, "The Unusual Structure of the Human Centromere (GGA)2 Motif: Unpaired Guanosine Residues Stacked Between Sheared G-A Pairs," *J. Molecular Biology*, vol. 244, no. 3, pp. 259–268, 1994.
- [29] B. Efron, "Selection criteria for scatterplot smoothers," *The Annals of Statistics*, vol. 29, no. 2, pp. 470–504, 2001.
- [30] D. J. C. Mackay, "Comparison of Approximate Methods for Handling Hyperparameters," *Neural Computation*, vol. 11, pp. 1035–1068, 1999.
- [31] K. Lange, D. R. Hunter, and I. Yang, "Optimization transfer using surrogate objective functions," *Journal of Computational and Graphical Statistics*, vol. 9, no. 1, pp. 1–20, 2000.
- [32] V. Solo, "A sure-fired way to choose smoothing parameters in ill-conditioned inverse problems," in *Proceedings of the IEEE Intl. Conf. on Image Processing*, vol. 3, 1996, pp. 89–92.

APPENDIX I PROOFS OF SECTION IV

A more general result is derived here. Consider the iteration

$$\hat{\underline{\theta}}^{(n+1)} = S_{T, \underline{\zeta}} \left(\hat{\underline{\theta}}^{(n)} + (\alpha/\sigma)^2 \mathbf{H}'(\underline{y} - \mathbf{H}\hat{\underline{\theta}}^{(n)}) \right), \quad (31)$$

where $T(\underline{x}; \underline{\zeta}) = \sum_i T(x_i; \underline{\zeta}) \underline{e}_i$ is a thresholding rule [15, Sec. 2.3] with the following condition. Suppose that $T(\cdot; \underline{\zeta})$ has *threshold* $t > 0$; then, $T(\cdot; \underline{\zeta})$ is strictly increasing on $\mathbb{R} \setminus (-t, t)$. Note that $T^{-1}(x; \underline{\zeta})$ is only defined for $x \neq 0$. Extend the definition at $x = 0$ to get

$$T^\dagger(x; \underline{\zeta}) = \begin{cases} 0 & x = 0 \\ T^{-1}(x; \underline{\zeta}) & x \neq 0 \end{cases} \quad (32)$$

$T^\dagger(x; \underline{\zeta})$ is continuous on $x \in \mathbb{R} \setminus \{0\}$. For the remainder of this section, the dependency of T , T^{-1} , and T^\dagger on $\underline{\zeta}$ will be omitted for the sake of brevity.

Proposition 2 The function

$$J_1(x) \triangleq 2T^\dagger(x)x - x^2 - 2 \int_0^\xi T(q) dq \Big|_{\xi=T^\dagger(x)} \quad (33)$$

is continuous for $x \in \mathbb{R}$.

Proof. Since $T^\dagger(x)$ is continuous in $\mathbb{R} \setminus \{0\}$, the only place that should be checked is $x = 0$. The second term in (33) is continuous, so it remains to check the first and third terms. By definition of a threshold function, $T^\dagger(0^+) = t$ and $T^\dagger(0^-) = -t$.

Consider $\epsilon > 0$. Since $T^\dagger(\cdot)$ is right continuous at 0^+ , there exists $\delta_1 > 0$ s.t. $x \in (0, \delta_1)$ implies that $|T^\dagger(x) - t| < 0.1t$. Likewise, since $T^\dagger(\cdot)$ is left continuous at 0^- , there exists $\delta_2 > 0$ s.t. $x \in (-\delta_2, 0)$ implies that $|T^\dagger(x) + t| < 0.1t$. Set

$$\delta = \frac{1}{2} \min(\delta_1, \delta_2, \frac{\epsilon}{1.1t})$$

so that $|x| < \delta \implies |xT^\dagger(x)| < \epsilon$.

Consider the third term. Define $A(\xi) \triangleq \int_0^\xi T(q) dq$: since $T(\cdot)$ is continuous, so is $A(\cdot)$. Moreover, for $|x| \leq t$, $A(x) = 0$. For $\epsilon > 0$, there exists $\kappa > t$ s.t. $|\xi| < \kappa \implies |A(\xi)| < \epsilon$. Since $T^\dagger(\cdot)$ is right continuous at 0^+ , there exists $\delta_1 > 0$ s.t. $x \in (0, \delta_1) \implies |T^\dagger(x) - t| < \kappa - t$. In a similar fashion, since $T^\dagger(\cdot)$ is left continuous at 0^- , there exists $\delta_2 > 0$ s.t. $x \in (-\delta_2, 0) \implies |T^\dagger(x) + t| < \kappa - t$. Set $\delta = \min(\delta_1, \delta_2)$. From $|x| < \delta$, one gets $|T^\dagger(x)| < \kappa$ whence $|A(T^\dagger(x))| < \epsilon$. ■

Proposition 3 The minimizer of $\varphi(x) = x^2 - 2cx + J_1(x)$ is $\tilde{x} = T(c)$.

Proof. Let $\varphi_1(x) \triangleq x^2 - 2cx$: $\varphi'_1(x) = 2(x - c)$, and is lower bounded. Similarly, consider $J_1(x)$: for $x \neq 0$, $J'_1(x) = 2(T^\dagger(x) - x)$. Since $0 \leq T(x) \leq x$ for all $x \geq 0$ and $x \leq T(x) \leq 0$ for all $x < 0$,

$$J'_1(x) = \begin{cases} \geq 0 & x > 0 \\ \leq 0 & x < 0 \end{cases}$$

$J_1(x)$ is also lower bounded. Applying Prop. 2 results in $\varphi(x)$ being a continuous, lower bounded function. Consider now two cases.

Case 1: $|c| > t$, where recall that t is the threshold of $T(\cdot)$. For $x \neq 0$, $\varphi'(x) = 2x - 2c + J'_1(x) = 2[T^\dagger(x) - c]$. So $\varphi'(x) = 0$ iff $T^\dagger(x) = c$, which occurs uniquely at $\tilde{x} = T(c) > 0$. Consider

$$\varphi'(T(c) + \delta) = 2[T^\dagger(T(c) + \delta) - c] \quad (34)$$

Since we assume that $T(\cdot)$ is strictly increasing on $\mathbb{R} \setminus (-t, t)$, $T^\dagger(x)$ is also strictly increasing for $x \neq 0$. For sufficiently small $\delta > 0$, $\varphi'(T(c) + \delta) > 0$ and $\varphi'(T(c) - \delta) < 0$. So $\tilde{x} = T(c)$ is a local minimum. At this value of x , $\varphi(x) = -2A(c) < 0$. To verify that \tilde{x} is the global minimum, it is necessary to compute $\varphi(0) = 0$. So indeed, $x = T(c)$ minimizes $\varphi(x)$.

Case 2: $|c| \leq t$. Suppose that the minimizer $x \neq 0$. Then, the analysis in Case 1 applies, resulting in $x = T(c)$. But since $|c| \leq t$ by assumption, one gets $x = 0$. This is a contradiction: it must therefore be the case that $\tilde{x} = 0$. ■

Theorem 3 Suppose that $\|\mathbf{H}\|_2 < 1$ and $\alpha = \sigma$. Consider the iteration (31), where $T(\cdot)$ is a thresholding rule with threshold $t > 0$, and $T(\cdot)$ is strictly increasing in $\mathbb{R} \setminus (-t, t)$. Then, the iterations (31) converge to a stationary point of $\Psi(\underline{\theta})$, where

$$\Psi(\underline{\theta}) = \|\mathbf{H}\underline{\theta} - \underline{y}\|_2^2 + J(\underline{\theta})$$

$$\text{where: } J(\underline{\theta}) \triangleq \sum_{i=1}^M J_1(\theta_i) \quad (35)$$

Proof. Use the following definitions, which appear in [17]:

$$\Xi(\underline{\theta}; \underline{a}) \triangleq C\|\underline{\theta} - \underline{a}\|_2^2 - \|\mathbf{H}\underline{\theta} - \mathbf{H}\underline{a}\|_2^2 \quad (36)$$

$$\Phi^{\text{SUR}}(\underline{\theta}; \underline{a}) \triangleq \Phi(\underline{\theta}) + \Xi(\underline{\theta}; \underline{a}), \quad (37)$$

where C is chosen to ensure that $\Xi(\underline{\theta}; \underline{a})$ is strictly positive and convex in $\underline{\theta}$ for any choice of \underline{a} . By assumption, $\|\mathbf{H}\|_2 < 1$, and so select $C = 1$ [17]. The function $\Phi^{\text{SUR}}(\underline{\theta}; \underline{a})$ is the surrogate function that is minimized in place of $\Phi(\underline{\theta})$. Consider the minimization of $\Phi^{\text{SUR}}(\underline{\theta}; \underline{a})$, which can be simplified as

$$\Phi^{\text{SUR}}(\underline{\theta}; \underline{a}) = \|\underline{\theta}\|_2^2 - 2(\underline{a} + \mathbf{H}'(\underline{y} - \mathbf{H}\underline{a}))'\underline{\theta} + J(\underline{\theta}) + \|\underline{y}\|_2^2 + \|\underline{a}\|_2^2 - \|\mathbf{H}\underline{a}\|_2^2 \quad (38)$$

Since $J(\underline{\theta}) = \sum_i J_1(\theta_i)$, the minimization of $\Phi^{\text{SUR}}(\underline{\theta}; \underline{a})$ can be decomposed into M subproblems, where each θ_i is separately minimized. Indeed, each θ_i should minimize

$$\varphi(\theta_i) \triangleq \theta_i^2 - 2s_i\theta_i + J_1(\theta_i), \quad (39)$$

where $s_i \triangleq \underline{a} + \mathbf{H}'(\underline{y} - \mathbf{H}\underline{a})$. Apply Prop. 3 to get the minimizing θ_i , i.e., $\theta_i = T(s_i)$.

Let $\hat{\underline{\theta}}^{(n)}$ denote the sequence generated by

$$\hat{\underline{\theta}}^{(n+1)} = \underset{\underline{\theta}}{\operatorname{argmin}} \Phi^{\text{SUR}}(\underline{\theta}; \hat{\underline{\theta}}^{(n)}) \quad (40)$$

where $\hat{\underline{\theta}}^{(0)}$ is the initial estimate. Then, $\hat{\underline{\theta}}^{(n)}$ is generated by (31), where recall that $\alpha/\sigma = 1$. Any limit point of the iterations (31) is a stationary point of (35) [31]. ■

APPENDIX II PROOFS OF SECTION V

A. Proof of Thm. 1

Recall that $\mathbf{G}(\mathbf{H}) = \mathbf{H}'\mathbf{H}$ is the Gram matrix of \mathbf{H} . In order to simplify notation, for $\mathbf{A} \in \mathbb{R}^{M \times M}$, denote by $\mathbf{A}_{11} = \mathbf{A}[1 : r, 1 : r]$, $\mathbf{A}_{12} = \mathbf{A}[1 : r, r+1 : M]$, $\mathbf{A}_{21} = \mathbf{A}[r+1 : M, 1 : r]$, and $\mathbf{A}_{22} = \mathbf{A}[r+1 : M, r+1 : M]$. The following proposition is needed. Its proof is omitted due to a lack of space.

Proposition 4 If \mathbf{H} has linearly independent columns,

$$\det((\mathbf{P}\mathbf{G}(\mathbf{H})\mathbf{P}')_{22}) \neq 0. \quad (41)$$

where \mathbf{P} is a matrix that orders the zero and non-zero components of $\hat{\underline{\theta}}$.

For $\hat{\underline{\theta}}$, an unbiased estimate of the l_2 risk (22) is [23], [32]

$$\hat{R}(\underline{\zeta}) = \sigma^2 + \frac{\|\underline{\epsilon}\|_2^2}{N} - \frac{2\sigma^2}{N} \sum_{n=1}^N \frac{\partial e_n}{\partial y_n} \quad (42)$$

where $\underline{\epsilon} = \underline{y} - \mathbf{H}\hat{\underline{\theta}}$. If $\hat{\underline{\theta}}$ is obtained via a minimization $\hat{\underline{\theta}} = \operatorname{argmin}_{\underline{\theta}} \Psi_{\underline{\zeta}}(\underline{\theta})$, (42) can be evaluated as [32, (2)]

$$\hat{R}(\underline{\zeta}) = \sigma^2 + \frac{\|\underline{\epsilon}\|_2^2}{N} - \frac{2\sigma^2}{N} \operatorname{tr}(\mathbf{H}(\mathbf{D}_{\underline{\theta}\underline{\theta}}\Psi_{\underline{\zeta}})^{-1}\mathbf{D}_{\underline{\theta}\underline{y}}\Psi_{\underline{\zeta}}) \Big|_{\underline{\theta}=\hat{\underline{\theta}}}, \quad (43)$$

where $\mathbf{D}_{\underline{u}, \underline{v}}(\cdot) \triangleq \partial^2(\cdot)/\partial \underline{u} \partial \underline{v}'$.

Let $\Psi_{\underline{\zeta}, l}(\underline{\theta}) = \|\mathbf{H}\underline{\theta} - \underline{y}\|_2^2 + \zeta_1 \|\underline{\theta}\|_1$ denote the cost function of lasso. Since $\Psi_{\underline{\zeta}, l}(\underline{\theta})$ is not twice differentiable on \mathbb{R}^M , (43) cannot be directly applied. Consider

$$\Psi_{\underline{\zeta}, l}(\underline{\theta}; a) = \|\mathbf{H}\underline{\theta} - \underline{y}\|_2^2 + \frac{2\zeta_1}{\pi} \sum_{m=1}^M \left\{ \theta_m \arctan\left(\frac{\theta_m}{a}\right) - \frac{a}{2} \ln\left(1 + \frac{\theta_m^2}{a^2}\right) \right\} \quad (44)$$

which is twice differentiable on \mathbb{R}^M . It can be shown that $\lim_{a \rightarrow 0} \Psi_{\underline{\zeta}, l}(\underline{\theta}; a) = \Psi_{\underline{\zeta}, l}(\underline{\theta})$ pointwise. The minimizer of $\Psi_{\underline{\zeta}, l}(\underline{\theta}; a)$ therefore equals the minimizer of $\Psi_{\underline{\zeta}, l}(\underline{\theta})$ in the limit as $a \rightarrow 0$. Denote by $\hat{R}_l(\underline{\zeta}; a)$ the unbiased estimate of (22) when $\hat{\underline{\theta}}$ is obtained by minimizing $\Psi_{\underline{\zeta}, l}(\underline{\theta}; a)$. As the RHS of (42) is solely a function of $\hat{\underline{\theta}}$ (recall that \underline{y} , \mathbf{H} , and σ^2 are known), $\lim_{a \rightarrow 0} \hat{R}_l(\underline{\zeta}; a) = \hat{R}_l(\underline{\zeta})$ pointwise.

Applying (43),

$$\hat{R}_l(\underline{\zeta}; a) = \sigma^2 + \frac{\|\underline{\epsilon}\|_2^2}{N} + \frac{2\sigma^2}{N} \operatorname{tr}(\mathbf{G}(\mathbf{H})[\mathbf{G}(\mathbf{H}) + \frac{1}{2}\mathbf{Z}_a(\hat{\underline{\theta}})]^{-1}) \quad (45)$$

where

$$\mathbf{Z}_a(\underline{\theta}) \triangleq \frac{2\zeta_1}{\pi} \operatorname{diag}\left(\frac{a}{a^2 + \theta_1^2}, \dots, \frac{a}{a^2 + \theta_M^2}\right). \quad (46)$$

Consider the $\operatorname{tr}(\cdot)$ expression in (45). As \mathbf{P} is orthogonal and matrix multiplication is commutative under the trace operator,

$$\operatorname{tr}(\mathbf{G}(\mathbf{H}) \left[\mathbf{G}(\mathbf{H}) + \frac{\mathbf{Z}_a(\hat{\underline{\theta}})}{2} \right]^{-1}) = \operatorname{tr}(\mathbf{P}\mathbf{G}(\mathbf{H})\mathbf{P}' \left[\mathbf{P}\mathbf{G}(\mathbf{H})\mathbf{P}' + \frac{\mathbf{P}\mathbf{Z}_a(\hat{\underline{\theta}})\mathbf{P}'}{2} \right]^{-1})$$

Without loss of generality, suppose that $\mathbf{Z}_a(\hat{\underline{\theta}})$ is ordered so that $\hat{\theta}_1 = \dots = \hat{\theta}_r = 0$, where $r = M - \|\hat{\underline{\theta}}\|_0$ and $\hat{\theta}_m \neq 0$ for $m > r$. Let $\mathbf{K} \triangleq \mathbf{P}\mathbf{G}(\mathbf{H})\mathbf{P}'$. Then, $[\mathbf{K} + \mathbf{Z}_a(\hat{\underline{\theta}})/2]^{-1}$ equals

$$\begin{pmatrix} \mathbf{F}_{11}^{-1} & -\tilde{\mathbf{K}}_{11}^{-1}\mathbf{K}_{12}\mathbf{F}_{22}^{-1} \\ -\mathbf{F}_{22}^{-1}\mathbf{K}_{21}\tilde{\mathbf{K}}_{11}^{-1} & \mathbf{F}_{22}^{-1} \end{pmatrix}$$

where $\tilde{\mathbf{K}}_{11} = \mathbf{K}_{11} + \mathbf{Z}_a(\hat{\underline{\theta}})_{11}$, $\tilde{\mathbf{K}}_{22} = \mathbf{K}_{22} + \mathbf{Z}_a(\hat{\underline{\theta}})_{22}$, $\mathbf{F}_{11} = \tilde{\mathbf{K}}_{11} - \mathbf{K}_{12}\tilde{\mathbf{K}}_{22}^{-1}\mathbf{K}_{21}$, and $\mathbf{F}_{22} = \tilde{\mathbf{K}}_{22} - \mathbf{K}_{21}\tilde{\mathbf{K}}_{11}^{-1}\mathbf{K}_{12}$. $\tilde{\mathbf{K}}_{11}$ is invertible for sufficiently small a . Likewise, for sufficiently small a , $\tilde{\mathbf{K}}_{22}$ is invertible by Prop. 4.

As $a \rightarrow 0$, $\tilde{\mathbf{K}}_{11}^{-1} \rightarrow \mathbf{0}$ and $\tilde{\mathbf{K}}_{22} \rightarrow \mathbf{G}(\mathbf{H})_{22}$. In addition, $\mathbf{F}_{11}^{-1} \rightarrow \mathbf{0}$ and $\mathbf{F}_{22} \rightarrow \tilde{\mathbf{K}}_{22}$. So

$$\mathbf{K} \left[\mathbf{K} + \frac{\mathbf{Z}_a(\hat{\underline{\theta}})}{2} \right]^{-1} \rightarrow \begin{pmatrix} \mathbf{0} & \mathbf{K}_{12}\mathbf{K}_{22}^{-1} \\ \mathbf{0} & \mathbf{I}_{22} \end{pmatrix} \quad (47)$$

as $a \rightarrow 0$. Consequently,

$$\lim_{a \rightarrow 0} \hat{R}_l(\underline{\zeta}; a) = \sigma^2 + \frac{1}{N} \|\underline{\epsilon}\|_2^2 + \frac{2\sigma^2}{N} \|\hat{\underline{\theta}}\|_0 = \hat{R}_l(\underline{\zeta}). \quad (48)$$

B. Proof of Thm. 2

Earlier notation from this appendix will be retained. The proof of the following proposition is omitted due to a lack of space.

Proposition 5 *Suppose that \mathbf{H} has linearly independent columns. If $\det(\mathbf{P}[\mathbf{G}(\mathbf{H}) - \frac{1}{2}\mathbf{U}(\hat{\theta})]\mathbf{P}') = 0$, then $\mathbf{G}(\mathbf{H})$ has an eigenvalue of $1/2$.*

The proof of Thm. 2 parallels the proof of Thm. 1. As $\Psi_{\underline{\zeta},hy}(\underline{\theta})$ is not twice differentiable on \mathbf{R}^M , consider instead

$$\begin{aligned} \Psi_{\underline{\zeta},hy}(\underline{\theta}; a) &= \Psi_{\underline{\zeta},l}(\underline{\theta}; a) \\ &+ \sum_{m=1}^M [G_1(\theta_m - \Delta_\zeta; a) - G_1(\theta_m + \Delta_\zeta; a)] \end{aligned} \quad (49)$$

where

$$G_1(x; a) \triangleq \frac{(a^2 + x^2) \arctan(x/a) - ax}{2\pi} + \frac{1}{4}x^2. \quad (50)$$

$\Psi_{\underline{\zeta},hy}(\underline{\theta}; a)$ is twice differentiable in \mathbf{R}^M and $\lim_{a \rightarrow 0} \Psi_{\underline{\zeta},hy}(\underline{\theta}; a) = \Psi_{\underline{\zeta},hy}(\underline{\theta})$ pointwise. Result (43) can be applied to get

$$\begin{aligned} \hat{R}_{hy}(\underline{\zeta}; a) &= \sigma^2 + \frac{\|\underline{e}\|_2^2}{N} \\ &+ \frac{2\sigma^2}{N} \text{tr}(\mathbf{G}(\mathbf{H})[\mathbf{G}(\mathbf{H}) + \frac{1}{2}\mathbf{Z}_a(\hat{\theta}) - \frac{1}{2}\mathbf{U}_a(\hat{\theta})]^{-1}) \end{aligned} \quad (51)$$

with

$$\mathbf{U}_a(\underline{\theta}) \triangleq \text{diag}((\ddot{G}_1(\theta_m - \Delta_\zeta; a) - \ddot{G}_1(\theta_m + \Delta_\zeta; a))_{m=1}^M) \quad (52)$$

Notice that similarity between $\Psi_{\underline{\zeta},l}(\underline{\theta}; a)$ and $\Psi_{\underline{\zeta},hy}(\underline{\theta}; a)$; the same applies to $\hat{R}_l(\underline{\zeta}; a)$ and $\hat{R}_{hy}(\underline{\zeta}; a)$. The steps of Thm. 1 can be carried out to evaluate the $\text{tr}(\cdot)$ expression in (51) as $a \rightarrow 0$. One arrives at

$$\lim_{a \rightarrow 0} \text{tr} \left(\mathbf{K}_{22} [\mathbf{K}_{22} - \frac{1}{2}(\mathbf{P}\mathbf{U}_a(\hat{\theta}))_{22}]^{-1} \right). \quad (53)$$

Now $\lim_{a \rightarrow 0} \mathbf{U}_a(\hat{\theta}) = \mathbf{U}(\hat{\theta})$. By assumption, $\mathbf{G}(\mathbf{H})$ does not have an eigenvalue of $1/2$. Therefore, application of Prop. 5 implies that the inverse in (53) exists.