

# Expectation-Maximization Gaussian-Mixture Approximate Message Passing

Jeremy Vila and Philip Schniter

Dept. of ECE, The Ohio State University, Columbus, OH 43210. (Email: vilaj@ece.osu.edu, schniter@ece.osu.edu)

**Abstract**—When recovering a sparse signal from noisy compressive linear measurements, the distribution of the signal’s non-zero coefficients can have a profound affect on recovery mean-squared error (MSE). If this distribution was apriori known, one could use efficient approximate message passing (AMP) techniques for nearly minimum MSE (MMSE) recovery. In practice, though, the distribution is unknown, motivating the use of robust algorithms like Lasso—which is nearly minimax optimal—at the cost of significantly larger MSE for non-least-favorable distributions. As an alternative, we propose an empirical-Bayesian technique that simultaneously learns the signal distribution while MMSE-recovering the signal—according to the learned distribution—using AMP. In particular, we model the non-zero distribution as a Gaussian mixture, and learn its parameters through expectation maximization, using AMP to implement the expectation step. Numerical experiments confirm the state-of-the-art performance of our approach on a range of signal classes.<sup>1 2</sup>

## I. INTRODUCTION

We consider estimating a  $K$ -sparse (or compressible) signal  $\mathbf{x} \in \mathbb{R}^N$  from  $M < N$  linear measurements  $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w} \in \mathbb{R}^M$ , where  $\mathbf{A}$  is known and  $\mathbf{w}$  is additive white Gaussian noise (AWGN). For this problem, accurate (relative to the noise variance) signal recovery is possible with polynomial-complexity algorithms when  $\mathbf{x}$  is sufficiently sparse and when  $\mathbf{A}$  satisfies certain restricted isometry properties [1].

A well-known approach to the sparse-signal recovery problem is Lasso [2], which solves the convex problem

$$\hat{\mathbf{x}}_{\text{lasso}} = \arg \min_{\hat{\mathbf{x}}} \|\mathbf{y} - \mathbf{A}\hat{\mathbf{x}}\|_2^2 + \lambda_{\text{lasso}} \|\hat{\mathbf{x}}\|_1, \quad (1)$$

with  $\lambda_{\text{lasso}}$  a tuning parameter. When  $\mathbf{A}$  is constructed from i.i.d entries, the performance of Lasso can be sharply characterized in the large system limit (i.e., as  $K, M, N \rightarrow \infty$  with fixed undersampling ratio  $M/N$  and sparsity ratio  $K/M$ ) using the so-called phase transition curve (PTC) [3]. When the observations are noiseless, the PTC bisects the  $M/N$ -versus- $K/M$  plane into the region where Lasso reconstructs the signal perfectly (with high probability) and the region where it does not. (See Figs. 1-3.) When the observations are noisy, the same PTC bisects the plane into the regions where Lasso’s noise sensitivity (i.e., the ratio of estimation-error power to measurement-noise power under the worst-case signal distribution) is either finite or infinite [4]. An important

fact about Lasso’s noiseless PTC is that it is invariant to signal distribution. In other words, if we consider the elements of the vector  $\mathbf{x}$  to be drawn i.i.d from the marginal pdf

$$p_X(x) = \lambda f_X(x) + (1 - \lambda)\delta(x), \quad (2)$$

where  $\delta(\cdot)$  is the Dirac delta,  $f_X(\cdot)$  is the active-coefficient pdf (with zero probability mass at  $x = 0$ ), and  $\lambda \triangleq K/N$ , then the Lasso PTC is invariant to  $f_X(\cdot)$ . While this implies that Lasso is robust to “difficult” sparse-signal distributions, it also implies that Lasso cannot benefit from the sparse-signal distribution being an “easy” one.

At the other end of the spectrum is minimum mean-squared error (MMSE)-optimal signal recovery under *known* marginal pdfs of the form (2). The PTC of MMSE recovery has been recently characterized [5] and shown to be well above that of Lasso. In particular, for *any*  $f_X(\cdot)$ , the PTC on the  $M/N$ -versus- $K/M$  plane equals the line  $K/M = 1$  in both the noiseless and noisy cases. Moreover, efficient algorithms for approximate MMSE-recovery have been proposed, such as the Bayesian version of Donoho, Maleki, and Montanari’s *approximate message passing* (AMP) algorithm from [6], which performs loopy belief-propagation on the underlying factor graph using central-limit-theorem approximations that become exact in the large-system limit under i.i.d  $\mathbf{A}$ . Although AMP’s complexity is remarkably low (e.g., dominated by one application of  $\mathbf{A}$  and  $\mathbf{A}^T$  per iteration with typically  $< 50$  iterations to convergence), it offers rigorous performance guarantees in the large-system limit. To handle arbitrary noise distributions and a wider class of matrices  $\mathbf{A}$ , Rangan proposed a *generalized AMP* (GAMP) [7] that forms the starting point of this work. (See Table I.)

In practice, one desires a recovery algorithm that does not need to know  $p_X(\cdot)$  a priori, yet offers performance on par with MMSE recovery, which (by definition) knows  $p_X(\cdot)$  a priori. Towards this aim, we propose a recovery scheme that aims to *learn* the prior signal distribution  $p_X(\cdot)$  (as well as the variance of the AWGN) while simultaneously recovering the signal vector  $\mathbf{x}$  from the noisy compressed measurements  $\mathbf{y}$ . To do so, we model the active component  $f_X(\cdot)$  in (2) using a generic  $L$ -term Gaussian mixture (GM) and then learn the prior signal and noise parameters using the expectation-maximization (EM) algorithm [8]. As we will see, the EM expectation is naturally implemented using the GAMP algorithm, which also provides approximately MMSE estimates of  $\mathbf{x}$ .

<sup>1</sup>This work has been supported in part by NSF-IUCRC grant IIP-0968910, by NSF grant CCF-1018368, and by DARPA/ONR grant N66001-10-1-4090.

<sup>2</sup>Portions of this work were presented in a poster at the Duke Workshop on Sensing and Analysis of High-Dimensional Data, July 2011.

Since we treat the prior pdf parameters as deterministic unknowns, our proposed EM-GM-GAMP algorithm can be considered as an “empirical-Bayesian” approach. Compared with previous empirical-Bayesian approaches (e.g., [9]–[12]), ours has a more flexible signal model, and thus is able to better match a wide range of signal pdfs  $p_X(\cdot)$ , as we demonstrate numerically in the sequel. Moreover, due to the computationally efficient nature of GAMP, our algorithm is significantly faster than empirical-Bayesian algorithms based on Tipping’s relevance vector machine [9]–[11]. Finally, we note that our EM-GM-GAMP algorithm can be considered as a generalization of our previously proposed EM-BG-GAMP algorithm [12] from a Bernoulli-Gaussian signal model to a Bernoulli-GM signal model.

## II. GAUSSIAN-MIXTURE GAMP

We first introduce Gaussian-mixture (GM) GAMP, a key component of our overall algorithm. In GM-GAMP, the signal  $\mathbf{x} = [x_1, \dots, x_N]^T$  is assumed to be i.i.d with marginal pdf

$$p_X(x; \lambda, \boldsymbol{\omega}, \boldsymbol{\theta}, \boldsymbol{\phi}) = (1 - \lambda)\delta(x) + \lambda \sum_{\ell=1}^L \omega_\ell \mathcal{N}(x; \theta_\ell, \phi_\ell), \quad (3)$$

where  $\delta(\cdot)$  denotes the Dirac delta,  $\lambda$  the sparsity rate, and, for the  $k^{th}$  GM component,  $\omega_k$ ,  $\theta_k$ , and  $\phi_k$  are the weight, mean, and variance, respectively. The AWGN<sup>3</sup>  $\mathbf{w}$  is then assumed to be independent of  $\mathbf{x}$  and have variance  $\psi$ :

$$p_W(w; \psi) = \mathcal{N}(w; 0, \psi) \quad (4)$$

Although above and in the sequel we assume real-valued Gaussians, all expressions can be converted to the circular-complex case by replacing  $\mathcal{N}$  with  $\mathcal{CN}$  and removing all  $\frac{1}{2}$ ’s. We emphasize that, from GM-GAMP’s perspective, the prior parameters  $\mathbf{q} \triangleq [\lambda, \boldsymbol{\omega}, \boldsymbol{\theta}, \boldsymbol{\phi}, \psi]$  are all known.

GAMP can handle an arbitrary probabilistic relationship  $p_{Y|Z}(y_m|z_m)$  between the observed output  $y_m$  and the noiseless output  $z_m \triangleq \mathbf{a}_m^T \mathbf{x}$ , where  $\mathbf{a}_m^T$  is the  $m^{th}$  row of  $\mathbf{A}$ . Our additive Gaussian noise assumption implies  $p_{Y|Z}(y|z) = \mathcal{N}(y; z, \psi)$ . To complete our description of GM-GAMP, we only need to derive  $g_{in}(\cdot)$ ,  $g'_{in}(\cdot)$ ,  $g_{out}(\cdot)$ , and  $g'_{out}(\cdot)$  in Table I. Using straightforward calculations, our  $p_{Y|Z}(\cdot|\cdot)$  yields [7]

$$g_{out}(y, \hat{z}, \mu^z; \mathbf{q}) = \frac{y - \hat{z}}{\mu^z + \psi} \quad (5)$$

$$-g'_{out}(y, \hat{z}, \mu^z; \mathbf{q}) = \frac{1}{\mu^z + \psi}, \quad (6)$$

and our GM signal prior (3) yields

$$g_{in}(\hat{r}, \mu^r; \mathbf{q}) = \frac{\sum_{\ell=1}^L \beta_\ell(\hat{r}, \mu^r; \mathbf{q}) \gamma_\ell(\hat{r}, \mu^r; \mathbf{q})}{(1 - \lambda)\mathcal{N}(0; \hat{r}, \mu^r) + \sum_{\ell=1}^L \beta_\ell(\hat{r}, \mu^r; \mathbf{q})} \quad (7)$$

$$\begin{aligned} \mu^r g'_{in}(\hat{r}, \mu^r; \mathbf{q}) &= -|g_{in}(\hat{r}, \mu^r; \mathbf{q})|^2 \\ &+ \frac{\sum_{\ell=1}^L \beta_\ell(\hat{r}, \mu^r; \mathbf{q}) (|\gamma_\ell(\hat{r}, \mu^r; \mathbf{q})|^2 + \nu_\ell(\hat{r}, \mu^r; \mathbf{q}))}{(1 - \lambda)\mathcal{N}(0; \hat{r}, \mu^r) + \sum_{\ell=1}^L \beta_\ell(\hat{r}, \mu^r; \mathbf{q})} \end{aligned} \quad (8)$$

<sup>3</sup>To model heavy-tailed noise, a different choice of  $p_W(\cdot; \cdot)$  may be more appropriate. However, the way it is handled in GAMP and learned by EM would remain essentially the same.

definitions:	
$p_{Z Y}(z y; \hat{z}, \mu^z) = \frac{p_{Y Z}(y z) \mathcal{N}(z; \hat{z}, \mu^z)}{\int_{z'} p_{Y Z}(y z') \mathcal{N}(z'; \hat{z}, \mu^z)}$	(D1)
$g_{out}(y, \hat{z}, \mu^z) = \frac{1}{\mu^z} (\mathbb{E}_{Z Y} \{z y; \hat{z}, \mu^z\} - \hat{z})$	(D2)
$g'_{out}(y, \hat{z}, \mu^z) = \frac{1}{\mu^z} \left( \frac{\text{var}_{Z Y} \{z y; \hat{z}, \mu^z\}}{\mu^z} - 1 \right)$	(D3)
$p_{X Y}(x y; \hat{r}, \mu^r) = \frac{p_X(x) \mathcal{N}(x; \hat{r}, \mu^r)}{\int_{x'} p_X(x') \mathcal{N}(x'; \hat{r}, \mu^r)}$	(D4)
$g_{in}(\hat{r}, \mu^r) = \int_x x p_{X Y}(x y; \hat{r}, \mu^r)$	(D5)
$g'_{in}(\hat{r}, \mu^r) = \frac{1}{\mu^r} \int_x  x - g_{in}(\hat{r}, \mu^r) ^2 p_{X Y}(x y; \hat{r}, \mu^r)$	(D6)
initialize:	
$\forall n : \hat{x}_n(1) = \int_x x p_X(x)$	(I1)
$\forall n : \mu_n^x(1) = \int_x  x - \hat{x}_n(1) ^2 p_X(x)$	(I2)
$\forall m : \hat{u}_m(0) = 0$	(I3)
for $t = 1, 2, 3, \dots$	
$\forall m : \hat{z}_m(t) = \sum_{n=1}^N A_{mn} \hat{x}_n(t)$	(R1)
$\forall m : \mu_m^z(t) = \sum_{n=1}^N  A_{mn} ^2 \mu_n^x(t)$	(R2)
$\forall m : \hat{p}_m(t) = \hat{z}_m(t) - \mu_m^z(t) \hat{u}_m(t-1)$	(R3)
$\forall m : \hat{u}_m(t) = g_{out}(y_m, \hat{p}_m(t), \mu_m^z(t))$	(R4)
$\forall m : \mu_m^u(t) = -g'_{out}(y_m, \hat{p}_m(t), \mu_m^z(t))$	(R5)
$\forall n : \mu_n^r(t) = \left( \sum_{m=1}^N  A_{mn} ^2 \mu_m^u(t) \right)^{-1}$	(R6)
$\forall n : \hat{r}_n(t) = \hat{x}_n(t) + \mu_n^r(t) \sum_{m=1}^M A_{mn}^* \hat{u}_m(t)$	(R7)
$\forall n : \mu_n^x(t+1) = \mu_n^r(t) g'_{in}(\hat{r}_n(t), \mu_n^r(t))$	(R8)
$\forall n : \hat{x}_n(t+1) = g_{in}(\hat{r}_n(t), \mu_n^r(t))$	(R9)
end	

TABLE I  
THE GAMP ALGORITHM [7]

where

$$\beta_\ell(\hat{r}, \mu^r; \mathbf{q}) \triangleq \lambda \omega_\ell \mathcal{N}(\hat{r}; \theta_\ell, \phi_\ell + \mu^r) \quad (9)$$

$$\gamma_\ell(\hat{r}, \mu^r; \mathbf{q}) \triangleq \frac{\hat{r}/\mu^r + \theta_\ell/\phi_\ell}{1/\mu^r + 1/\phi_\ell} \quad (10)$$

$$\nu_\ell(\hat{r}, \mu^r; \mathbf{q}) \triangleq \frac{1}{1/\mu^r + 1/\phi_\ell}. \quad (11)$$

Table I implies that GM-GAMP’s marginal posteriors are

$$p(x_n | \mathbf{y}; \mathbf{q}) = p_X(x_n; \mathbf{q}) \mathcal{N}(x_n; \hat{r}_n, \mu_n^r) / \zeta(\hat{r}_n, \mu_n^r; \mathbf{q}) \quad (12)$$

$$\begin{aligned} &= \left( (1 - \lambda)\delta(x_n) + \lambda \sum_{\ell=1}^L \omega_\ell \mathcal{N}(x_n; \theta_\ell, \phi_\ell) \right) \\ &\quad \times \mathcal{N}(x_n; \hat{r}_n, \mu_n^r) / \zeta(\hat{r}_n, \mu_n^r; \mathbf{q}) \end{aligned} \quad (13)$$

$$\zeta(\hat{r}, \mu^r; \mathbf{q}) \triangleq \int_x p_X(x; \mathbf{q}) \mathcal{N}(x; \hat{r}, \mu^r). \quad (14)$$

From (13), it is straightforward to show that the posterior support probabilities returned by GM-GAMP are

$$\Pr\{x_n \neq 0 | \mathbf{y}; \mathbf{q}\} = \pi(\hat{r}_n, \mu_n^r; \mathbf{q}) \quad (15)$$

$$\pi(\hat{r}, \mu^r; \mathbf{q}) \triangleq \frac{1}{1 + \left( \frac{\sum_{\ell=1}^L \beta_\ell(\hat{r}, \mu^r; \mathbf{q})}{(1 - \lambda)\mathcal{N}(0; \hat{r}, \mu^r)} \right)^{-1}}. \quad (16)$$

## III. EM LEARNING OF THE PRIOR PARAMETERS $\mathbf{q}$

We use the expectation-maximization (EM) algorithm [8] to learn the prior parameters  $\mathbf{q} \triangleq [\lambda, \boldsymbol{\omega}, \boldsymbol{\theta}, \boldsymbol{\phi}, \psi]$ . The EM algorithm is an iterative technique that increases a lower bound on the likelihood  $p(\mathbf{y}; \mathbf{q})$  at each iteration, thus guaranteeing that the likelihood converges to a local maximum. In our case, the “hidden data” is chosen as  $\{\mathbf{x}, \mathbf{w}\}$ , implying the iteration- $i$  EM update

$$\mathbf{q}^{i+1} = \arg \max_{\mathbf{q}} \mathbb{E} \{ \ln p(\mathbf{x}, \mathbf{w}; \mathbf{q}) | \mathbf{y}; \mathbf{q}^i \}, \quad (17)$$

where  $E\{\cdot|\mathbf{y}; \mathbf{q}^i\}$  denotes expectation conditioned on the observations  $\mathbf{y}$  under the parameter hypothesis  $\mathbf{q}^i$ . Since it is impractical to update the entire vector  $\mathbf{q}$  at once, we update  $\mathbf{q}$  one element at a time (while holding the others fixed), which can be recognized as the “incremental” technique from [13]. In the sequel, we use “ $\mathbf{q}_{\setminus\lambda}^i$ ” to denote the vector  $\mathbf{q}^i$  with  $\lambda^i$  removed (and similar for the other parameters).

#### A. EM update for $\lambda$

We now derive the EM update for  $\lambda$  given previous parameters  $\mathbf{q}^i \triangleq [\lambda^i, \boldsymbol{\omega}^i, \boldsymbol{\theta}^i, \boldsymbol{\phi}^i, \boldsymbol{\psi}^i]$ . Since  $\mathbf{x}$  is apriori independent of  $\boldsymbol{\omega}$  and i.i.d, the joint pdf  $p(\mathbf{x}, \mathbf{w}; \mathbf{q})$  decouples into  $C \prod_{n=1}^N p_X(x_n; \mathbf{q})$  for a  $\lambda$ -invariant constant  $C$ , and so

$$\lambda^{i+1} = \arg \max_{\lambda \in (0,1)} \sum_{n=1}^N E \{ \ln p_X(x_n; \lambda, \mathbf{q}_{\setminus\lambda}^i) | \mathbf{y}; \mathbf{q}^i \}. \quad (18)$$

The maximizing value of  $\lambda$  in (18) is necessarily a value of  $\lambda$  that zeroes the derivative, i.e., that satisfies

$$\sum_{n=1}^N \int_{x_n} p(x_n | \mathbf{y}; \mathbf{q}^i) \frac{d}{d\lambda} \ln p_X(x_n; \lambda, \mathbf{q}_{\setminus\lambda}^i) = 0. \quad (19)$$

For the  $p_X(x_n; \mathbf{q})$  given in (3), it is readily seen that

$$\begin{aligned} \frac{d}{d\lambda} \ln p_X(x_n; \lambda, \mathbf{q}_{\setminus\lambda}^i) &= \frac{\sum_{\ell=1}^L \omega_\ell^i \mathcal{N}(x_n; \theta_\ell^i, \phi_\ell^i) - \delta(x_n)}{p_X(x_n; \lambda, \mathbf{q}_{\setminus\lambda}^i)} \\ &= \begin{cases} \frac{1}{\lambda} & x_n \neq 0 \\ \frac{-1}{1-\lambda} & x_n = 0 \end{cases}. \end{aligned} \quad (20)$$

Plugging (20) and (13) into (19), it becomes evident that the point  $x_n = 0$  must be treated differently than  $x_n \in \mathbb{R} \setminus 0$ . Thus, we define the closed ball  $\mathcal{B}_\epsilon = [-\epsilon, \epsilon]$  and  $\overline{\mathcal{B}}_\epsilon \triangleq \mathbb{R} \setminus \mathcal{B}_\epsilon$ , and note that, in the limit  $\epsilon \rightarrow 0$ , the following becomes equivalent to (19), given the definition of  $\pi(\hat{r}, \mu^r; \mathbf{q})$  in (16):

$$\frac{1}{\lambda} \sum_{n=1}^N \underbrace{\int_{x_n \in \overline{\mathcal{B}}_\epsilon} p(x_n | \mathbf{y}; \mathbf{q}^i)}_{\stackrel{\epsilon \rightarrow 0}{=} \pi(\hat{r}_n, \mu_n^r; \mathbf{q}^i)} = \frac{1}{1-\lambda} \sum_{n=1}^N \underbrace{\int_{x_n \in \mathcal{B}_\epsilon} p(x_n | \mathbf{y}; \mathbf{q}^i)}_{\stackrel{\epsilon \rightarrow 0}{=} 1 - \pi(\hat{r}_n, \mu_n^r; \mathbf{q}^i)}. \quad (21)$$

To verify that the left integral converges to the  $\pi(\hat{r}_n, \mu_n^r; \mathbf{q}^i)$  defined in (16), it suffices to plug (13) into (21) and apply the Gaussian-pdf multiplication rule.<sup>4</sup> Meanwhile, for any  $\epsilon$ , the right integral must equal one minus the left. Thus, the EM update for  $\lambda$  is the unique value satisfying (21) as  $\epsilon \rightarrow 0$ , i.e.,

$$\lambda^{i+1} = \frac{1}{N} \sum_{n=1}^N \pi(\hat{r}_n, \mu_n^r; \mathbf{q}^i). \quad (22)$$

Conveniently,  $\{\pi(\hat{r}_n, \mu_n^r; \mathbf{q}^i)\}_{n=1}^N$  are GM-GAMP outputs, as we recall from (16).

<sup>4</sup> $\mathcal{N}(x; a, A) \mathcal{N}(x; b, B) = \mathcal{N}(x; \frac{a/A+b/B}{1/A+1/B}, \frac{1}{1/A+1/B}) \mathcal{N}(0; a-b, A+B)$ .

#### B. EM updates for Gaussian Mixture Parameters

For each  $k = 1, \dots, L$ , we incrementally update each GM parameter  $\theta_k$ ,  $\phi_k$ , and  $\boldsymbol{\omega}$  while holding the others fixed. The EM updates become

$$\theta_k^{i+1} = \arg \max_{\theta_k \in \mathbb{R}} \sum_{n=1}^N E \{ \ln p_X(x_n; \theta_k, \mathbf{q}_{\setminus\theta_k}^i) | \mathbf{y}, \mathbf{q}^i \}, \quad (23)$$

$$\phi_k^{i+1} = \arg \max_{\phi_k > 0} \sum_{n=1}^N E \{ \ln p_X(x_n; \phi_k, \mathbf{q}_{\setminus\phi_k}^i) | \mathbf{y}, \mathbf{q}^i \} \quad (24)$$

$$\boldsymbol{\omega}^{i+1} = \arg \max_{\boldsymbol{\omega} > 0: \sum_k \omega_k = 1} \sum_{n=1}^N E \{ \ln p_X(x_n; \boldsymbol{\omega}, \mathbf{q}_{\setminus\boldsymbol{\omega}}^i) | \mathbf{y}, \mathbf{q}^i \} \quad (25)$$

Following (19), the maximizing value of  $\theta_k$  in (23) is necessarily a value of  $\theta_k$  that zeros

$$\sum_{n=1}^N \int_{x_n} p(x_n | \mathbf{y}, \mathbf{q}^i) \frac{d}{d\theta_k} \ln p_X(x_n; \theta_k, \mathbf{q}_{\setminus\theta_k}^i) = 0, \quad (26)$$

where  $p(x_n | \mathbf{y}, \mathbf{q}^i) = p_X(x_n; \mathbf{q}^i) \mathcal{N}(x_n; \hat{r}_n, \mu_n^r) / \zeta(\hat{r}_n, \mu_n^r; \mathbf{q}^i)$  from (D4), recalling  $p_X(x; \mathbf{q})$  from (3) and  $\zeta(\hat{r}, \mu^r; \mathbf{q})$  from (14). Taking the derivative, we find

$$\begin{aligned} \frac{d}{d\theta_k} \ln p_X(x_n; \theta_k, \mathbf{q}_{\setminus\theta_k}^i) &= \left( \frac{x_n - \theta_k}{\phi_k^i} \right) \\ &\times \frac{\lambda^i \omega_k^i \mathcal{N}(x_n; \theta_k, \phi_k^i)}{(1-\lambda^i) \delta(x_n) + \lambda^i (\omega_k^i \mathcal{N}(x_n; \theta_k, \phi_k^i) + \sum_{\ell \neq k} \omega_\ell^i \mathcal{N}(x_n; \theta_\ell^i, \phi_\ell^i))}. \end{aligned} \quad (27)$$

Integrating (26) separately over  $\mathcal{B}_\epsilon$  and  $\overline{\mathcal{B}}_\epsilon$ , as in (21), and taking  $\epsilon \rightarrow 0$ , we find that the  $\mathcal{B}_\epsilon$  portion vanishes, giving

$$\sum_{n=1}^N \int_{x_n} \frac{p(x_n | x_n \neq 0, \mathbf{y}, \mathbf{q}^i) \lambda^i \omega_k^i \mathcal{N}(x_n; \theta_k, \phi_k^i) (x_n - \theta_k)}{\zeta(\hat{r}_n, \mu_n^r; \mathbf{q}^i) (\omega_k^i \mathcal{N}(x_n; \theta_k, \phi_k^i) + \sum_{\ell \neq k} \omega_\ell^i \mathcal{N}(x_n; \theta_\ell^i, \phi_\ell^i))} = 0. \quad (28)$$

Since this integral is difficult to evaluate, we apply the approximation  $\mathcal{N}(x_n; \theta_k, \phi_k^i) \approx \mathcal{N}(x_n; \theta_k^i, \phi_k^i)$  and exploit the fact that  $p(x_n | x_n \neq 0, \mathbf{y}, \mathbf{q}^i) = \mathcal{N}(x_n; \hat{r}_n, \mu_n^r) \sum_{\ell} \omega_\ell^i \mathcal{N}(x_n; \theta_\ell^i, \phi_\ell^i)$  to cancel terms, giving

$$\sum_{n=1}^N \int_{x_n} \frac{\lambda^i \omega_k^i \mathcal{N}(x_n; \hat{r}_n, \mu_n^r) \mathcal{N}(x_n; \theta_k^i, \phi_k^i)}{\zeta_n(\hat{r}_n, \mu_n^r; \mathbf{q}^i)} (x_n - \theta_k) = 0. \quad (29)$$

We then simplify (29) using the Gaussian-pdf multiplication rule,<sup>4</sup> and set  $\theta_k^{i+1}$  equal to the value of  $\theta_k$  satisfying (29):

$$\theta_k^{i+1} = \frac{\sum_{n=1}^N \Pr\{x_n \neq 0, k_n = k | \mathbf{y}, \mathbf{q}^i\} \gamma_k(\hat{r}, \mu^r; \mathbf{q}^i)}{\sum_{n=1}^N \Pr\{x_n \neq 0, k_n = k | \mathbf{y}, \mathbf{q}^i\}}, \quad (30)$$

where the joint activity/mixture probabilities are

$$\begin{aligned} \Pr\{x_n \neq 0, k_n = k | \mathbf{y}, \mathbf{q}^i\} \\ \triangleq \frac{\beta_k(\hat{r}_n, \mu_n^r; \mathbf{q}^i)}{(1-\lambda) \mathcal{N}(0; \hat{r}_n, \mu_n^r) + \sum_{\ell=1}^L \beta_\ell(\hat{r}_n, \mu_n^r; \mathbf{q}^i)} \end{aligned} \quad (31)$$

with  $\beta_k(\hat{r}, \mu^r; \mathbf{q}^i)$  and  $\gamma_k(\hat{r}, \mu^r; \mathbf{q}^i)$  defined in (9)-(10). Above, “ $k_n = k$ ” represents the event that  $x_n$  was generated from mixture component  $k$ .

Following (26), the maximizing value of  $\phi_k$  in (24) is necessarily a value of  $\phi_k$  that zeroes the derivative

$$\sum_{n=1}^N \int_{x_n} p(x_n | \mathbf{y}, \mathbf{q}^i) \frac{d}{d\phi_k} \ln p_X(x_n; \phi_k, \mathbf{q}_{\setminus\phi_k}^i) = 0. \quad (32)$$

Taking the derivative, we find

$$\frac{d}{d\phi_k} \ln p_X(x_n; \phi_k, \mathbf{q}_{\setminus \phi_k}^i) = \frac{1}{2} \left( \frac{|x_n - \theta_k^i|^2}{\phi_k^2} - \frac{1}{\phi_k} \right) \times \frac{\lambda^i \omega_k^i \mathcal{N}(x_n; \theta_k^i, \phi_k)}{(1 - \lambda^i) \delta(x_n) + \lambda^i (\omega_k^i \mathcal{N}(x_n; \theta_k^i, \phi_k) + \sum_{\ell \neq k} \omega_\ell^i \mathcal{N}(x_n; \theta_\ell^i, \phi_\ell^i))}. \quad (33)$$

Integrating (32) separately over  $\mathcal{B}_\epsilon$  and  $\overline{\mathcal{B}_\epsilon}$ , as in (21), and taking  $\epsilon \rightarrow 0$ , we find that the  $\mathcal{B}_\epsilon$  portion vanishes, giving

$$\sum_{n=1}^N \int_{x_n} \frac{p(x_n | x_n \neq 0, \mathbf{y}, \mathbf{q}^i) \lambda^i \omega_k^i \mathcal{N}(x_n; \theta_k^i, \phi_k)}{\zeta_n(\hat{r}_n, \mu_n^r; \mathbf{q}^i) (\omega_k^i \mathcal{N}(x_n; \theta_k^i, \phi_k) + \sum_{\ell \neq k} \omega_\ell^i \mathcal{N}(x_n; \theta_\ell^i, \phi_\ell^i))} \times \left( \frac{|x_n - \theta_k^i|^2}{\phi_k} - 1 \right) = 0. \quad (34)$$

Similar to (28), this integral is difficult to evaluate, and so we apply the approximation  $\mathcal{N}(x_n; \theta_k^i, \phi_k) \approx \mathcal{N}(x_n; \theta_k^i, \phi_k^i)$ , after which certain terms cancel, yielding

$$\sum_{n=1}^N \int_{x_n} \frac{\mathcal{N}(x_n; \hat{r}_n, \mu_n^r) \lambda^i \omega_k^i \mathcal{N}(x_n; \theta_k^i, \phi_k^i)}{\zeta_n(\hat{r}_n, \mu_n^r; \mathbf{q}^i)} \left( \frac{|x_n - \theta_k^i|^2}{\phi_k} - 1 \right) = 0. \quad (35)$$

To find the value of  $\phi_k$  satisfying (35), we expand  $|x_n - \theta_k^i|^2 = |x_n|^2 - 2 \operatorname{Re}(x_n^* \theta_k^i) + |\theta_k^i|^2$  and apply the Gaussian-pdf multiplication rule,<sup>4</sup> which gives

$$\phi_k^{i+1} = \frac{\sum_{n=1}^N \Pr\{x_n \neq 0, k_n = k | \mathbf{y}, \mathbf{q}^i\} (|\theta_k^i - \gamma_k(\hat{r}, \mu^r; \mathbf{q}^i)|^2 + \nu_k(\hat{r}, \mu^r; \mathbf{q}))}{\sum_{n=1}^N \Pr\{x_n \neq 0, k_n = k | \mathbf{y}, \mathbf{q}^i\}} \quad (36)$$

where  $\Pr\{x_n \neq 0, k | \mathbf{y}, \mathbf{q}^i\}$  was given in (31).

Finally, the value of the pmf-constrained  $\omega$  maximizing (25) can be found by solving the unconstrained optimization problem  $\max_{\omega, \xi} J(\omega, \xi)$ , where  $\xi$  is a Lagrange multiplier and

$$J(\omega, \xi) \triangleq \sum_{n=1}^N \mathbb{E} \{ \ln p_X(x_n; \omega, \mathbf{q}_{\setminus \omega}^i) | \mathbf{y}, \mathbf{q}^i \} - \xi \left( \sum_{\ell=1}^L \omega_\ell - 1 \right). \quad (37)$$

We start by setting  $\frac{d}{d\omega_k} J(\omega, \xi) = 0$ , which yields

$$\sum_{n=1}^N \int_{x_n} \frac{p_X(x_n; \mathbf{q}^i) \mathcal{N}(x_n; \hat{r}_n, \mu_n^r)}{\zeta(\hat{r}_n, \mu_n^r; \mathbf{q}^i)} \frac{\lambda^i \mathcal{N}(x_n; \theta_k^i, \phi_k^i)}{p_X(x_n; \omega, \mathbf{q}_{\setminus \omega}^i)} = \xi. \quad (38)$$

Like (28) and (34), the above is difficult to evaluate, and so we approximate  $\omega \approx \omega^i$ , which leads to

$$\xi = \sum_{n=1}^N \int_{x_n} \frac{\lambda^i \mathcal{N}(x_n; \theta_k^i, \phi_k^i) \mathcal{N}(x_n; \hat{r}_n, \mu_n^r)}{\zeta(\hat{r}_n, \mu_n^r; \mathbf{q}^i)}. \quad (39)$$

Multiplying both sides by  $\omega_k^i$ , summing over  $k = 1, \dots, L$ , employing the fact  $1 = \sum_k \omega_k^i$ , and simplifying, we obtain

$$\xi = \sum_{n=1}^N \int_{x_n} \frac{\lambda^i \sum_{k=1}^L \omega_k^i \mathcal{N}(x_n; \theta_k^i, \phi_k^i) \mathcal{N}(x_n; \hat{r}_n, \mu_n^r)}{\zeta(\hat{r}_n, \mu_n^r; \mathbf{q}^i)} \quad (40)$$

$$= \sum_{n=1}^N \Pr\{x_n \neq 0 | \mathbf{y}, \mathbf{q}^i\}. \quad (41)$$

Plugging (41) into (39) and multiplying both sides by  $\omega_k$ , the derivative-zeroing value of  $\omega_k$  is seen to be

$$\omega_k = \frac{\sum_{n=1}^N \int_{x_n} \lambda^i \omega_k \mathcal{N}(x_n; \theta_k^i, \phi_k^i) \mathcal{N}(x_n; \hat{r}_n, \mu_n^r) / \zeta(\hat{r}_n, \mu_n^r; \mathbf{q}^i)}{\sum_{n=1}^N \Pr\{x_n \neq 0 | \mathbf{y}, \mathbf{q}^i\}}, \quad (42)$$

where, if we use  $\omega_k \approx \omega_k^i$  on the right of (42), then we obtain

$$\omega_k^{i+1} = \frac{\sum_{n=1}^N \Pr\{x_n \neq 0, k_n = k | \mathbf{y}, \mathbf{q}^i\}}{\sum_{n=1}^N \Pr\{x_n \neq 0 | \mathbf{y}, \mathbf{q}^i\}}. \quad (43)$$

Note that the numerator and denominator of (43) can be computed from GM-GAMP outputs via (16) and (31).

### C. EM update for $\psi$

The final parameter to estimate is the noise energy  $\psi$ . In this paper, the AWGN noise model given in (4) is identical to EM-BG-GAMP's [12]. There, the (exact) EM update for  $\psi$  is<sup>5</sup>

$$\psi^{i+1} = \frac{1}{M} \sum_{m=1}^M (|y_m - \hat{z}_m|^2 + \mu_m^z). \quad (44)$$

## IV. LEARNING THE MODEL ORDER

So far, the GM model order  $L$  has been treated as fixed and known. In practice, one could indeed choose a fixed value  $L$  that is thought to be large enough to capture the essential structure of  $p_X(\cdot)$  and, given an appropriate initialization of the  $3L+2$  parameters  $\mathbf{q}^0$ , apply the previously described EM-GM-GAMP algorithm to jointly estimate  $\mathbf{x}$  and  $\mathbf{q}$ .

As an alternative, one could instead start with the model order  $L = 1$  (i.e., a Bernoulli-Gaussian model for  $p_X(\cdot)$ ) and increment  $L$  one-by-one, stopping as soon as negligible benefits are observed (e.g.,  $\|\hat{\mathbf{x}}_L - \hat{\mathbf{x}}_{L-1}\|_2^2 / \|\hat{\mathbf{x}}_{L-1}\|_2^2 < \text{tol}$ ) or a predefined  $L_{\max}$  has been reached. Here, EM-GM-GAMP would be re-run as described in Sections II-III at each new value of  $L$ . This latter approach would relieve the user from the potentially difficult task of choosing both  $L$  and a many-parameter  $\mathbf{q}^0$  apriori. In the remainder of this section, we propose a particular implementation of this approach.

When growing the model-order from  $L$  to  $L+1$ , we propose to split the mixture component  $k_* \in \{1, \dots, L\}$  with the “worst fit” into two new components. To select the worst-fitting mixture component, one could use the approach of split-and-merge-EM [14], i.e., maximization of local Kullback-Leibler divergence, or similar. Given the mixture-component-to-split  $k_*$ , the subset of coefficient indices  $n$  that are most probably associated with  $k_*$  is identified, i.e.,

$$\mathfrak{N}_{k_*} \triangleq \{n \in \mathfrak{N} : \arg \max_k \Pr\{x_n \neq 0, k_n = k | \mathbf{y}, \mathbf{q}\} = k_*\}, \quad (45)$$

where  $\mathfrak{N} \triangleq \{n : \Pr\{x_n \neq 0 | \mathbf{y}, \mathbf{q}\} > 0.5\}$  is the subset of coefficient indices that are most probably non-zero. To simplify the notation, we henceforth assume, without loss of generality, that  $k_* = L$ .

To split the  $L^{\text{th}}$  mixture component, we replace the mean  $\theta_L$ , variance  $\phi_L$ , and weight  $\omega_L$  with two new values for each (e.g.,  $\theta_L^{\text{new}}$  and  $\theta_{L+1}^{\text{new}}$  replace  $\theta_L$ ), resulting in an  $L+1$ -term parameter vector  $\mathbf{q}^{\text{new}}$ . Rather than considering only a single possibility for  $\mathbf{q}^{\text{new}}$ , we consider  $S$  possibilities  $\{\mathbf{q}_s^{\text{new}}\}_{s=1}^S$  obtained as variations of the following two strategies:

<sup>5</sup>Sometimes we observe that the EM update for  $\psi$  works better with the  $\mu_m^z$  term in (44) weighted by  $\frac{M}{N}$  and suppressed until later EM iterations. We conjecture that this is due to bias in the GAMP variance estimates  $\mu_m^z$ .



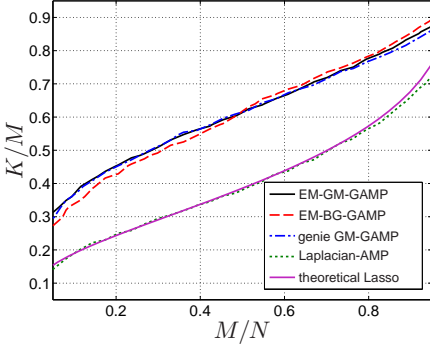


Fig. 1. Noiseless empirical PTCs and Lasso theoretical PTC for Bernoulli-Gaussian signals.

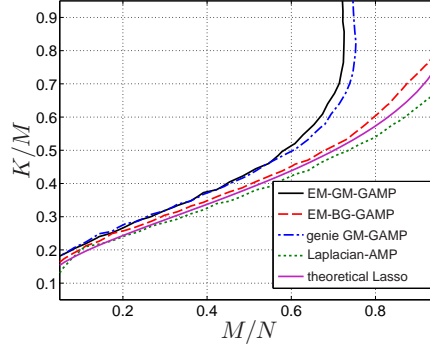


Fig. 2. Noiseless empirical PTCs and Lasso theoretical PTC for Bernoulli-Rademacher signals.

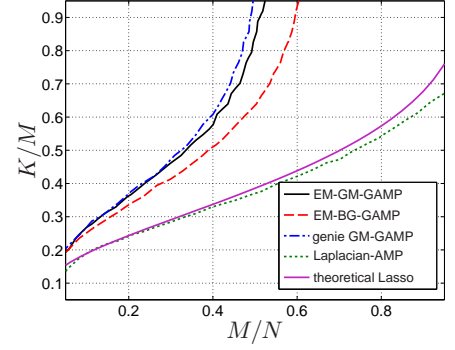


Fig. 3. Noiseless empirical PTCs and Lasso theoretical PTC for Bernoulli signals.

- 1) **Split-mean(a)**:  $\theta_L^{\text{new}} = \theta_L - a$ ,  $\theta_{L+1}^{\text{new}} = \theta_L + a$ ,  $\phi_L^{\text{new}} = \phi_{L+1}^{\text{new}} = \phi_L$ , and  $\omega_L^{\text{new}} = \omega_{L+1}^{\text{new}} = \omega_L/2$ ; and
- 2) **Split-variance(b)**:  $\phi_L^{\text{new}} = b\phi_L$  and  $\phi_{L+1}^{\text{new}} = b^{-1}\phi_L$ ,  $\theta_L^{\text{new}} = \theta_{L+1}^{\text{new}} = \theta_L$ , and  $\omega_L^{\text{new}} = \omega_{L+1}^{\text{new}} = \omega_L/2$ ,

where  $a, b > 0$  are design parameters. Note that, by considering several distinct values of  $a$  and/or  $b$ , we have  $S > 2$ . Finally, to judge which of the  $S$  possible splits is best, one could, e.g., evaluate the corresponding likelihoods, i.e., solve

$$\arg \max_s \sum_n \mathbb{E} \{ \ln p_X(\mathbf{x}; \mathbf{q}_s^{\text{new}}) | \mathbf{y}; \mathbf{q} \}. \quad (46)$$

Empirically, we have found that the incremental method of learning  $L$  described above<sup>6</sup> works very well for “sparse” signals like Bernoulli-Gaussian, Bernoulli-Rademacher, Bernoulli-Uniform, and Bernoulli. (See below for details.) For “heavy-tailed” signals like Student-t, however, it seems better to fix  $L$  at a reasonable value (e.g.,  $L = 4$ ), keep the means at zero (i.e.,  $\theta_k^i = 0 \forall k, i$ ), and jointly learn the  $L$  weights  $\omega$ , the  $L$  variances  $\phi$ , and the sparsity rate  $\lambda$ .

## V. EM INITIALIZATION

Since the EM algorithm is guaranteed to converge only to a local maximum of the likelihood function, proper initialization of  $\mathbf{q}$  is essential. Here, we describe initialization strategies for the “sparse” and “heavy-tailed” modes described above.

For the “sparse” mode, where initially  $L = 1$  (and thus  $\omega_1^0 = 0$ ), we use the same initializations that we proposed for EM-BG-GAMP [12], i.e.,  $\lambda^0 = \frac{M}{N} \rho_{\text{SE}}(\frac{M}{N})$ , where  $\rho_{\text{SE}}(\frac{M}{N})$  is the sparsity ratio  $\frac{K}{M}$  achieved by the noiseless Lasso PTC [3]

$$\rho_{\text{SE}}(\frac{M}{N}) = \max_{c \geq 0} \frac{1 - \frac{2N}{M} [(1 + c^2)\Phi(c) - c\phi(c)]}{1 - c^2 - 2[(1 + c^2)\Phi(c) - c\phi(c)]}, \quad (47)$$

with  $\Phi(\cdot)$  and  $\phi(\cdot)$  the cdf and pdf of the normal distribution, respectively, and

$$\psi^0 = \frac{\|\mathbf{y}\|_2^2}{(\text{SNR}^0 + 1)M}, \quad \phi_1^0 = \frac{\|\mathbf{y}\|_2^2 - M\psi^0}{\text{tr}(\mathbf{A}^T \mathbf{A})\lambda^0}, \quad \theta_1^0 = 0, \quad (48)$$

where, without other knowledge, we suggest  $\text{SNR}^0 = 100$ .

For the “heavy-tailed” mode, we suggest initializing  $\lambda^0$  and  $\psi^0$  as above and, with  $L = 4$ , choosing

$$\omega_k^0 = \frac{1}{L}, \quad \phi_k^0 = \frac{k}{\sqrt{L}} \frac{(\|\mathbf{y}\|_2^2 - M\psi^0)}{\text{tr}(\mathbf{A}^T \mathbf{A})\lambda^0}, \quad \theta_k^0 = 0, \quad k = 1 \dots L. \quad (49)$$

<sup>6</sup>For the simulations, we used  $S = 3$  splitting methods in the “sparse mode”:  $a = \sqrt{\phi_L}$ ,  $b_1 = 3$ ,  $b_2 = 6$ .

## VI. NUMERICAL RESULTS

### A. Noiseless Phase Transitions

First, we describe the results of experiments that computed noiseless empirical phase transition curves (PTCs) under three sparse-signal distributions. To evaluate each empirical PTC, we constructed a  $30 \times 30$  grid of oversampling ratio  $\frac{M}{N} \in [0.05, 0.95]$  and sparsity ratio  $\frac{K}{M} \in [0.05, 0.95]$  for fixed signal length  $N = 1000$ . At each grid point, we generated  $R = 100$  independent realizations of  $K$ -sparse signal  $\mathbf{x}$  and  $M \times N$  measurement matrix with i.i.d  $\mathcal{N}(0, M^{-1})$  entries. From the measurements  $\mathbf{y} = \mathbf{A}\mathbf{x}$ , we attempted to reconstruct the signal  $\mathbf{x}$  using various algorithms. A recovery  $\hat{\mathbf{x}}$  from realization  $r \in \{1, \dots, R\}$  was defined a success (i.e.,  $S_r = 1$ ) if the NMSE  $\triangleq \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 / \|\mathbf{x}\|_2^2 < 10^{-4}$ , and the average success rate was defined as  $\bar{S} \triangleq \frac{1}{R} \sum_{r=1}^R S_r$ . The empirical PTC was then plotted, using Matlab’s `contour` command, as the  $\bar{S} = 0.5$  contour over the sparsity-undersampling grid.

Figures 1-3 show the empirical PTCs for four recovery algorithms: the proposed EM-GM-GAMP algorithm<sup>7</sup> (in “sparse” mode), the EM-BG-GAMP algorithm from [12], a “genie-aided” GM-GAMP with the true  $[\lambda, \omega, \theta, \phi, \psi]$ , and the Laplacian-AMP from [3]. For comparison, Figs. 1-3 also display the theoretical Lasso PTC (47). The signals were generated as Bernoulli-Gaussian (BG) in Fig. 1 ( $\theta = 0$ ,  $\phi = 1$ ), as Bernoulli-Rademacher (BR) in Fig. 2 (i.e., non-zero coefficients chosen uniformly from  $\{-1, 1\}$ ), and as Bernoulli in Fig. 3 (i.e., all non-zero coefficients set equal to 1).

For all three signal types, Figs. 1-3 show that the empirical PTC of EM-GM-GAMP significantly improves on those of Laplacian-AMP and theoretical Lasso. (The latter two converge in the large system limit [3].) For Bernoulli-Gaussian signals, EM-GM-GAMP performed very similarly to genie-GM-GAMP and EM-BG-GAMP. Such behavior is expected, because all three can accurately model the signal distribution. For BR signals, however, EM-GM-GAMP performed significantly better than EM-BG-GAMP, since it can better model the BR distribution (using an  $L = 2$  GM). It even performed slightly better than genie-GM-GAMP here, since it is able to perform realization-specific parameter fitting. For Bernoulli signals, EM-GM-GAMP performed moderately better than EM-BG-GAMP, and nearly the same as genie-GM-GAMP.

<sup>7</sup>Matlab code available at <http://www.ece.osu.edu/~schniter/EMturboGAMP>

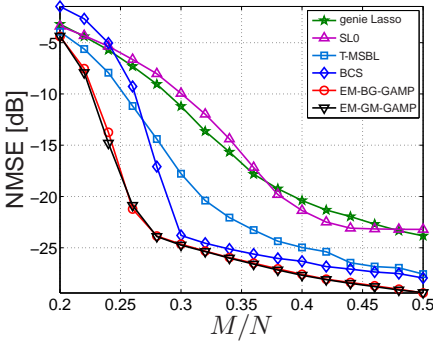


Fig. 4. NMSE for noisy recovery of Bernoulli-Gaussian signals.

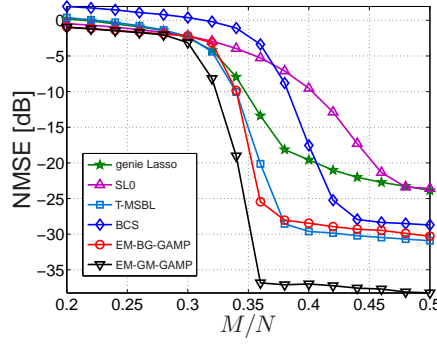


Fig. 5. NMSE for noisy recovery of Bernoulli-Rademacher signals.

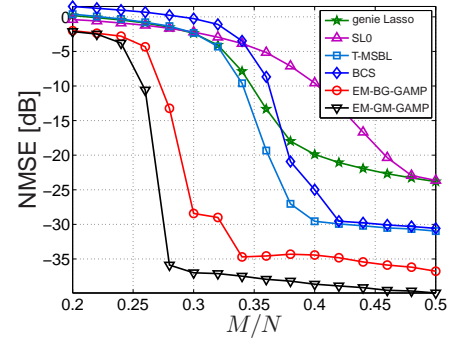


Fig. 6. NMSE for noisy recovery of Bernoulli signals.

### B. Noisy Signal Recovery

Figures 4-6 show NMSE for noisy recovery of BG, BR, and Bernoulli signals. To construct these plots, we fixed  $N = 1000$ ,  $K = 100$ ,  $\text{SNR} = 25\text{dB}$ , and varied  $M$ . Each data point represents NMSE averaged over  $R = 500$  realizations. For comparison, we show the performance of the proposed EM-GM-GAMP (in “sparse” mode), EM-BG-GAMP [12], Bayesian Compressive Sensing (BCS) [11], Sparse Bayesian Learning (via T-MSBL [15]), debiased genie-aided<sup>8</sup> Lasso (via SPGL1 [16]), and Smoothed- $\ell_0$  (SLO) [17]. All algorithms were run under the suggested defaults, with ‘noise=small’ in T-MSBL.

For BG signals, Fig. 4 shows that EM-GM-GAMP exhibits the best performance (together with EM-BG-GAMP). For BR and Bernoulli signals, however, Figs. 5-6 show that EM-GM-GAMP significantly outperforms the other algorithms. Relative to EM-BG-GAMP, EM-GM-GAMP’s greatest improvement comes with BR signals, which are not well-modeled using a BG prior. We have verified, using all three signal types, that EM-GM-GAMP’s excellent behavior persists at lower SNRs, as well as on sparse  $L$ -ary discrete signals with  $L > 2$ .

Perhaps most impressive is EM-GM-GAMP’s performance in recovering heavy-tailed signals. As an example, Fig. 7 shows noisy recovery NMSE for a Student’s-t signal with pdf

$$p_X(x; q) \triangleq \frac{\Gamma((q+1)/2)}{\sqrt{2\pi}\Gamma(q/2)} (1 + x^2)^{-(q+1)/2} \quad (50)$$

under the *non-compressible* parameter choice  $q = 1.67$  [18]. Here, EM-GM-GAMP was run in “heavy-tailed” mode and outperformed all other algorithms under test—even genie-aided Lasso. Although the algorithms that perform best on the sparse signals in Figs. 4-6 usually perform worst on heavy-tailed signals like that in Fig. 7, and vice versa, Figs. 4-7 show EM-GM-GAMP excelling on *all* signal types.

In conclusion, we attribute EM-GM-GAMP’s excellent performance, on a wide range of signal types, to its ability to near-optimally learn and exploit a wide range of signal priors.

### REFERENCES

- [1] Y. C. Eldar and G. Kutyniok, *Compressed Sensing: Theory and Applications*. New York: Cambridge Univ. Press, 2012.
- [2] R. Tibshirani, “Regression shrinkage and selection via the lasso,” *J. Roy. Statist. Soc. B*, vol. 58, no. 1, pp. 267 – 288, 1996.

<sup>8</sup>We ran SPGL1 in ‘BPDN’ mode:  $\min_{\hat{\mathbf{x}}} \|\mathbf{x}\|_1$  s.t.  $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2 \leq \sigma$ , for tolerances  $\sigma^2 \in \{0.1, 0.2, \dots, 1.5\} \times M\psi$ , and reported the lowest NMSE.

- [3] D. L. Donoho, A. Maleki, and A. Montanari, “Message passing algorithms for compressed sensing,” *Proc. Nat. Acad. Sci.*, vol. 106, pp. 18914–18919, Nov. 2009.
- [4] D. L. Donoho, A. Maleki, and A. Montanari, “The noise-sensitivity phase transition in compressed sensing,” *arXiv:1004.1218*, Apr. 2010.
- [5] Y. Wu and S. Verdú, “Optimal phase transitions in compressed sensing,” *arXiv:1111.6822*, Nov. 2011.
- [6] D. L. Donoho, A. Maleki, and A. Montanari, “Message passing algorithms for compressed sensing: I. Motivation and construction,” in *Proc. Inform. Theory Workshop*, (Cairo, Egypt), Jan. 2010.
- [7] S. Rangan, “Generalized approximate message passing for estimation with random linear mixing,” *arXiv:1010.5141*, Oct. 2010.
- [8] A. Dempster, N. M. Laird, and D. B. Rubin, “Maximum-likelihood from incomplete data via the EM algorithm,” *J. Roy. Statist. Soc.*, vol. 39, pp. 1–17, 1977.
- [9] M. E. Tipping, “Sparse Bayesian learning and the relevance vector machine,” *J. Mach. Learn. Res.*, vol. 1, pp. 211–244, 2001.
- [10] D. P. Wipf and B. D. Rao, “Sparse Bayesian learning for basis selection,” *IEEE Trans. Signal Process.*, vol. 52, pp. 2153 – 2164, Aug. 2004.
- [11] S. Ji, Y. Xue, and L. Carin, “Bayesian compressive sensing,” *IEEE Trans. Signal Process.*, vol. 56, pp. 2346–2356, June 2008.
- [12] J. P. Vila and P. Schniter, “Expectation-maximization Bernoulli-Gaussian approximate message passing,” in *Proc. Asilomar Conf. Signals Syst. Comput.*, (Pacific Grove, CA), Nov. 2011.
- [13] R. Neal and G. Hinton, “A view of the EM algorithm that justifies incremental, sparse, and other variants,” in *Learning in Graphical Models* (M. I. Jordan, ed.), pp. 355–368, MIT Press, 1999.
- [14] N. Ueda, R. Nakano, Z. Ghahramani, and G. E. Hinton, “SMEM algorithm for mixture models,” *Neural Comput.*, vol. 12, pp. 2109–2128, Sept. 2000.
- [15] Z. Zhang and B. D. Rao, “Sparse signal recovery with temporally correlated source vectors using sparse Bayesian learning,” *IEEE J. Sel. Topics Signal Process.*, vol. 5, pp. 912–926, Sept. 2011.
- [16] E. van den Berg and M. P. Friedlander, “Probing the Pareto frontier for basis pursuit solutions,” *SIAM J. Scientific Comput.*, vol. 31, no. 2, pp. 890–912, 2008.
- [17] H. Mohimani, M. Babaie-Zadeh, and C. Jutten, “A fast approach for overcomplete sparse decomposition based on smoothed norm,” *IEEE Trans. Signal Process.*, vol. 57, pp. 289–301, Jan. 2009.
- [18] V. Cevher, “Learning with compressible priors,” in *Proc. Neural Inform. Process. Syst. Conf.*, (Vancouver, B.C.), Dec. 2009.

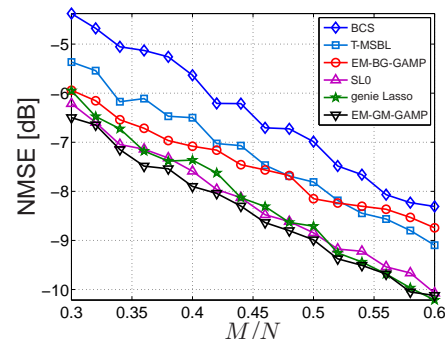


Fig. 7. NMSE for noisy recovery of non-compressible Student’s-t signals.