

# Expectation-Maximization Bernoulli-Gaussian 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**—The approximate message passing (AMP) algorithm originally proposed by Donoho, Maleki, and Montanari yields a computationally attractive solution to the usual  $\ell_1$ -regularized least-squares problem faced in compressed sensing, whose solution is known to be robust to the signal distribution. When the signal is drawn i.i.d from a marginal distribution that is not least-favorable, better performance can be attained using a Bayesian variation of AMP. The latter, however, assumes that the distribution is perfectly known. In this paper, we navigate the space between these two extremes by modeling the signal as i.i.d Bernoulli-Gaussian (BG) with unknown prior sparsity, mean, and variance, and the noise as zero-mean Gaussian with unknown variance, and we simultaneously reconstruct the signal while learning the prior signal and noise parameters. To accomplish this task, we embed the BG-AMP algorithm within an expectation-maximization (EM) framework. Numerical experiments confirm the excellent performance of our proposed EM-BG-AMP on a range of signal types.<sup>12</sup>

## I. INTRODUCTION

We consider the problem of recovering a signal  $\mathbf{x} \in \mathbb{R}^N$  from noisy linear measurements  $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w} \in \mathbb{R}^M$  in the “undersampled” regime where  $M < N$ . Roughly speaking, when  $\mathbf{x}$  is sufficiently sparse (or compressible) and when the matrix  $\mathbf{A}$  is sufficiently well-conditioned, accurate signal recovery is possible with polynomial-complexity algorithms.

One of the best-known approaches to this problem is known as “Lasso” [1], which minimizes the convex criterion

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

with  $\lambda_{\text{lasso}}$  a tuning parameter. When  $\mathbf{A}$  is constructed from i.i.d Gaussian entries, the so-called phase transition curve (PTC) gives a sharp characterization of Lasso performance for  $K$ -sparse  $\mathbf{x}$  in the large system limit, i.e., as  $K, M, N \rightarrow \infty$  with fixed undersampling ratio  $M/N$  and sparsity ratio  $K/M$  [2]. (The Lasso PTC is illustrated in Figs. 1-3.) For noiseless observations, the PTC partitions the  $M/N$ -versus- $K/M$  plane into two regions: one where Lasso reconstructs the signal perfectly (with high probability), and one where it does not. For noisy observations, the same curve indicates whether the noise sensitivity (i.e., the ratio of estimation-error power to measurement-noise power under the worst-case signal distribution) of Lasso remains bounded [3].

One remarkable feature of the noiseless Lasso PTC is that it is invariant to signal distribution. In other words, if we adopt a

probabilistic viewpoint where the elements of  $\mathbf{x}$  are drawn i.i.d from the marginal pdf  $p_X(x) = \lambda f(x) + (1-\lambda)\delta(x)$ , for Dirac delta  $\delta(x)$ , active-coefficient pdf  $f(x)$ , and  $\lambda \triangleq K/N$ , then the Lasso PTC is not affected by  $f(\cdot)$ . This PTC invariance implies that Lasso is robust, but that it cannot benefit from the restriction of  $\mathbf{x}$  to an “easier” signal class. For example, if the coefficients in  $\mathbf{x}$  are known to be non-negative, then there exists a polynomial-complexity algorithm whose PTC is better than that of Lasso [2].

Although, in some applications, robustness to worst-case signals may be the dominant concern, in many other applications the goal is to maximize average-case performance over an anticipated signal class. When the signal  $\mathbf{x}$  is drawn i.i.d from an arbitrary *known* marginal distribution  $p_X(\cdot)$  and the noise  $\mathbf{w}$  is i.i.d Gaussian with *known* variance, there exist very-low-complexity iterative Bayesian algorithms to generate approximately maximum a posteriori (MAP) and minimum mean-squared error (MMSE) signal estimates, notably the Bayesian version of Donoho, Maleki, and Montanari’s *approximate message passing* (AMP) algorithm [4]. AMP is formulated from a loopy-belief-propagation perspective, leveraging central-limit-theorem approximations that hold in the large system limit for suitably dense  $\mathbf{A}$ , and admits a rigorous analysis in the large-system limit [5]. Meanwhile, AMP’s complexity is remarkably low, dominated by one application of  $\mathbf{A}$  and  $\mathbf{A}^\top$  per iteration (which is especially cheap if  $\mathbf{A}$  is an FFT or other fast operation), with typically  $< 50$  iterations to convergence. More recently, a *generalized AMP* (GAMP) algorithm [6] was proposed that relaxes the requirements on the noise distribution and on the sensing matrix  $\mathbf{A}$ . (See Table I for a summary.)

Given that it is rare to know the signal and noise distributions perfectly, we take the approach of assuming signal and noise distributions that are known up to some statistical parameters, and then learning those unknown parameters while simultaneously recovering the signal. Examples of this “empirical Bayesian” approach include several algorithms based on Tipping’s *relevance vector machine* [7]–[9]. Although the average-case performance of those algorithms is often quite good (depending on the signal class, of course), their complexities are generally much larger than that of (G)AMP.

In this paper, we propose a GAMP-based empirical-Bayesian algorithm. In particular, we treat the signal as Bernoulli-Gaussian (BG) signal with unknown sparsity, mean, and variance, and the noise as Gaussian with unknown variance, and then we then learn these statistical parameters using

<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.

an expectation-maximization (EM) approach [10] that calls BG-GAMP once per EM-iteration.

## II. BERNOULLI-GAUSSIAN GAMP

A core component of our proposed method is the Bernoulli-Gaussian (BG) GAMP algorithm, which we now review. For BG-GAMP, the signal  $\mathbf{x} = [x_1, \dots, x_N]^T$  is assumed to be i.i.d BG, i.e., to have marginal pdf

$$p_X(x; \lambda, \theta, \phi) = (1 - \lambda)\delta(x) + \lambda\mathcal{N}(x; \theta, \phi), \quad (2)$$

where  $\delta(\cdot)$  denotes the Dirac delta,  $\lambda$  the sparsity rate,  $\theta$  the active-coefficient mean, and  $\phi$  the active-coefficient variance. The noise  $\mathbf{w}$  is assumed to be independent of  $\mathbf{x}$  and i.i.d zero-mean Gaussian with variance  $\psi$ :

$$p_W(w; \psi) = \mathcal{N}(w; 0, \psi) \quad (3)$$

In our approach, the parameters  $\mathbf{q} \triangleq [\lambda, \theta, \phi, \psi]$  that define these prior distributions are treated as deterministic unknowns, and learned through the EM algorithm, as detailed in Section III. Although above and in the sequel we assume real-valued Gaussians, all expressions can be converted to the circular-complex-Gaussian case by replacing all  $\mathcal{N}$  with  $\mathcal{CN}$  and removing all  $\frac{1}{2}$ 's.

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  $m^{\text{th}}$  row of  $\mathbf{A}$ . Our additive Gaussian noise assumption implies  $p_{Y|Z}(y|z) = \mathcal{N}(y; z, \psi)$ . To complete our description, we need only to specify  $g_{\text{in}}(\cdot)$ ,  $g'_{\text{in}}(\cdot)$ ,  $g_{\text{out}}(\cdot)$ , and  $g'_{\text{out}}(\cdot)$  in Table I. Using straightforward manipulations, our  $p_{Y|Z}(\cdot|\cdot)$  yields [6]

$$g_{\text{out}}(y, \hat{z}, \mu^z; \mathbf{q}) = \frac{y - \hat{z}}{\mu^z + \psi} \quad (4)$$

$$-g'_{\text{out}}(y, \hat{z}, \mu^z; \mathbf{q}) = \frac{1}{\mu^z + \psi}, \quad (5)$$

and our BG signal prior (2) yields

$$g_{\text{in}}(\hat{r}, \mu^r; \mathbf{q}) = \pi(\hat{r}, \mu^r; \mathbf{q}) \gamma(\hat{r}, \mu^r; \mathbf{q}) \quad (6)$$

$$\begin{aligned} \mu^r g'_{\text{in}}(\hat{r}, \mu^r; \mathbf{q}) &= \pi(\hat{r}, \mu^r; \mathbf{q}) (\nu(\hat{r}, \mu^r; \mathbf{q}) + |\gamma(\hat{r}, \mu^r; \mathbf{q})|^2) \\ &\quad - (\pi(\hat{r}, \mu^r; \mathbf{q}))^2 |\gamma(\hat{r}, \mu^r; \mathbf{q})|^2, \end{aligned} \quad (7)$$

where

$$\pi(\hat{r}, \mu^r; \mathbf{q}) \triangleq \frac{1}{1 + \left( \frac{\lambda}{1-\lambda} \frac{\mathcal{N}(\hat{r}; \theta, \phi + \mu^r)}{\mathcal{N}(\hat{r}; 0, \mu^r)} \right)^{-1}} \quad (8)$$

$$\gamma(\hat{r}, \mu^r; \mathbf{q}) \triangleq \frac{\hat{r}/\mu^r + \theta/\phi}{1/\mu^r + 1/\phi} \quad (9)$$

$$\nu(\hat{r}, \mu^r; \mathbf{q}) \triangleq \frac{1}{1/\mu^r + 1/\phi}. \quad (10)$$

Table I implies that BG-GAMP's marginal posteriors are

$$p(x_n | \mathbf{y}; \mathbf{q}) = \frac{1}{C_n} p_X(x_n; \mathbf{q}) \mathcal{N}(x_n; \hat{r}_n, \mu_n^r) \quad (11)$$

$$\begin{aligned} &= \frac{1}{C_n} \left( (1 - \lambda)\delta(x_n) + \lambda\mathcal{N}(x_n; \theta, \phi) \right) \\ &\quad \times \mathcal{N}(x_n; \hat{r}_n, \mu_n^r) \end{aligned} \quad (12)$$

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)
$g_{\text{out}}(y, \hat{z}, \mu^z) = \frac{1}{\mu^z} (\mathbb{E}_{Z Y} \{z y; \hat{z}, \mu^z\} - \hat{z})$	(D2)
$g'_{\text{out}}(y, \hat{z}, \mu^z) = \frac{1}{\mu^z} \left( \frac{\text{var}_{Z Y} \{z y; \hat{z}, \mu^z\}}{\mu^z} - 1 \right)$	(D3)
$p_{X Y}(x \mathbf{y}; \hat{r}, \mu^r) = \frac{p_X(x) \mathcal{N}(x; \hat{r}, \mu^r)}{\int_{x'} p_X(x') \mathcal{N}(x'; \hat{r}, \mu^r)}$	(D4)
$g_{\text{in}}(\hat{r}, \mu^r) = \int_{x'} x p_{X Y}(x \mathbf{y}; \hat{r}, \mu^r)$	(D5)
$g'_{\text{in}}(\hat{r}, \mu^r) = \frac{1}{\mu^r} \int_x  x - g_{\text{in}}(\hat{r}, \mu^r) ^2 p_{X Y}(x \mathbf{y}; \hat{r}, \mu^r)$	(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)
$\forall m : \hat{u}_m(t) = g_{\text{out}}(y_m, \hat{p}_m(t), \mu_m^z(t))$	(R4)
$\forall m : \mu_m^u(t) = -g'_{\text{out}}(y_m, \hat{p}_m(t), \mu_m^z(t))$	(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)
$\forall n : \mu_n^x(t+1) = \mu_n^r(t) g'_{\text{in}}(\hat{r}_n(t), \mu_n^r(t))$	(R8)
$\forall n : \hat{x}_n(t+1) = g_{\text{in}}(\hat{r}_n(t), \mu_n^r(t))$	(R9)
end	

TABLE I  
THE GAMP ALGORITHM [6]

for scaling factor  $C_n \triangleq \int p_X(x_n; \mathbf{q}) \mathcal{N}(x_n; \hat{r}_n, \mu_n^r)$ . From (12), it is straightforward to show that BG-GAMP yields the following posterior support probabilities:

$$\Pr\{x_n \neq 0 | \mathbf{y}; \mathbf{q}\} = \pi(\hat{r}_n, \mu_n^r; \mathbf{q}). \quad (13)$$

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

We use the expectation-maximization (EM) algorithm [10] to learn the statistical parameters  $\mathbf{q} \triangleq [\lambda, \theta, \phi, \psi]$ . The EM algorithm is an iterative technique that increases the likelihood at each iteration, guaranteeing convergence to a local maximum of the likelihood  $p(\mathbf{y}; \mathbf{q})$ . In our case, we choose the ‘‘hidden data’’ to be  $\{\mathbf{x}, \mathbf{w}\}$ , which yields the EM update

$$\mathbf{q}^{i+1} = \arg \max_{\mathbf{q}} \mathbb{E} \{ \ln p(\mathbf{x}, \mathbf{w}; \mathbf{q}) | \mathbf{y}; \mathbf{q}^i \}, \quad (14)$$

where  $i$  denotes EM iteration and  $\mathbb{E}\{\cdot | \mathbf{y}; \mathbf{q}^i\}$  denotes expectation conditioned on the observations  $\mathbf{y}$  under the parameter hypothesis  $\mathbf{q}^i$ . Moreover, we use the well-established ‘‘incremental’’ updating schedule [11], where  $\mathbf{q}$  is updated one element at a time while keeping the other elements fixed.

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

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

$$\lambda^{i+1} = \arg \max_{\lambda \in (0,1)} \sum_{n=1}^N \mathbb{E} \{ \ln p_X(x_n; \lambda, \theta^i, \phi^i) | \mathbf{y}; \mathbf{q}^i \}. \quad (15)$$

The maximizing value of  $\lambda$  in (15) is necessarily a value of  $\lambda$  that zeroes the derivative, i.e., that satisfies

$$\sum_{n=1}^N \int_{x_n} p(x_n | \mathbf{y}; \mathbf{q}^i) \frac{d}{d\lambda} \ln p_X(x_n; \lambda, \theta^i, \phi^i) = 0. \quad (16)$$

For the  $p_X(x_n; \lambda, \theta, \phi)$  given in (2), it is readily seen that

$$\begin{aligned} & \frac{d}{d\lambda} \ln p_X(x_n; \lambda, \theta^i, \phi^i) \\ &= \frac{\mathcal{N}(x_n; \theta^i, \phi^i) - \delta(x_n)}{p_X(x_n; \lambda, \theta^i, \phi^i)} = \begin{cases} \frac{1}{\lambda} & x_n \neq 0 \\ \frac{-1}{1-\lambda} & x_n = 0 \end{cases}. \end{aligned} \quad (17)$$

Plugging (17) and (12) into (16), it becomes evident that the neighborhood around the point  $x_n = 0$  should be treated differently than the remainder of  $\mathbb{R}$ . 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 is equivalent to (16):

$$\frac{1}{\lambda} \sum_{n=1}^N \underbrace{\int_{x_n \in \overline{\mathcal{B}_\epsilon}} p(x_n | \mathbf{y}; \mathbf{q}^i)}_{\stackrel{\epsilon \rightarrow 0}{=} \pi(\hat{r}_n, \mu_n^r; \mathbf{q}^i)} = \frac{1}{1-\lambda} \sum_{n=1}^N \underbrace{\int_{x_n \in \mathcal{B}_\epsilon} p(x_n | \mathbf{y}; \mathbf{q}^i)}_{\stackrel{\epsilon \rightarrow 0}{=} 1 - \pi(\hat{r}_n, \mu_n^r; \mathbf{q}^i)}. \quad (18)$$

To verify that the left integral converges to the  $\pi(\hat{r}_n, \mu_n^r; \mathbf{q}^i)$  defined in (8), it suffices to plug (12) into (18) and apply the Gaussian-pdf multiplication rule;<sup>3</sup> meanwhile, for any  $\epsilon$ , the right integral must equal one minus the left. Finally, the EM update for  $\lambda$  is the unique value satisfying (18) as  $\epsilon \rightarrow 0$ , i.e.,

$$\lambda^{i+1} = \frac{1}{N} \sum_{n=1}^N \pi(\hat{r}_n, \mu_n^r; \mathbf{q}^i). \quad (19)$$

Conveniently,  $\{\pi(\hat{r}_n, \mu_n^r; \mathbf{q}^i)\}_{n=1}^N$  are GAMP outputs.

### B. EM update for $\theta$

Similar to (15), the EM update<sup>4</sup> for  $\theta$  can be written as

$$\theta^{i+1} = \arg \max_{\theta \in \mathbb{R}} \sum_{n=1}^N \mathbb{E} \{ \ln p_X(x_n; \lambda^i, \theta, \phi^i) | \mathbf{y}; \mathbf{q}^i \}, \quad (20)$$

The maximizing value of  $\theta$  in (20) is necessarily a value of  $\theta$  that zeroes the derivative, i.e., that satisfies

$$\sum_{n=1}^N \int_{x_n} p(x_n | \mathbf{y}; \mathbf{q}^i) \frac{d}{d\theta} \ln p_X(x_n; \lambda^i, \theta, \phi^i) = 0. \quad (21)$$

For the  $p_X(x_n; \lambda, \theta, \phi)$  given in (2), it is readily seen that

$$\begin{aligned} & \frac{d}{d\theta} \ln p_X(x_n; \lambda^i, \theta, \phi^i) \\ &= \frac{(x_n - \theta)}{\phi^i} \frac{\lambda^i \mathcal{N}(x_n; \theta, \phi^i)}{p_X(x_n; \lambda^i, \theta, \phi^i)} = \begin{cases} \frac{x_n - \theta}{\phi^i} & x_n \neq 0 \\ 0 & x_n = 0 \end{cases}. \end{aligned} \quad (22)$$

Splitting the domain of integration in (21) into  $\mathcal{B}_\epsilon$  and  $\overline{\mathcal{B}_\epsilon}$ , and then plugging in (22), we find that the following is equivalent to (21) in the limit of  $\epsilon \rightarrow 0$ :

$$\sum_{n=1}^N \int_{x_n \in \overline{\mathcal{B}_\epsilon}} (x_n - \theta) p(x_n | \mathbf{y}; \mathbf{q}^i) = 0. \quad (23)$$

<sup>3</sup> $\mathcal{N}(x; a, A) \mathcal{N}(x; b, B) = \mathcal{N}(x; \frac{a/A+b/B}{1/A+1/B}, \frac{1}{1/A+1/B}) \mathcal{N}(0; a-b, A+B)$ .

<sup>4</sup>If the user has good reason to believe that the true signal pdf is symmetric around zero, then they may consider fixing  $\theta=0$  and avoiding this EM update.

The unique value of  $\theta$  satisfying (23) as  $\epsilon \rightarrow 0$  is then

$$\theta^{i+1} = \frac{\sum_{n=1}^N \lim_{\epsilon \rightarrow 0} \int_{x_n \in \overline{\mathcal{B}_\epsilon}} x_n p(x_n | \mathbf{y}; \mathbf{q}^i)}{\sum_{n=1}^N \lim_{\epsilon \rightarrow 0} \int_{x_n \in \overline{\mathcal{B}_\epsilon}} p(x_n | \mathbf{y}; \mathbf{q}^i)} \quad (24)$$

$$= \frac{1}{\lambda^{i+1} N} \sum_{n=1}^N \pi(\hat{r}_n, \mu_n^r; \mathbf{q}^i) \gamma(\hat{r}_n, \mu_n^r; \mathbf{q}^i) \quad (25)$$

for the GAMP outputs  $\{\gamma(\hat{r}_n, \mu_n^r; \mathbf{q}^i)\}_{n=1}^N$  defined in (9). The equality in (25) can be verified by plugging the GAMP posterior expression from (12) into (24) and simplifying via the Gaussian-pdf multiplication rule.

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

Similar to (15), the EM update for  $\phi$  can be written as

$$\hat{\phi}^{i+1} = \arg \max_{\phi > 0} \sum_{n=1}^N \mathbb{E} \{ \ln p_X(x_n; \lambda^i, \theta^i, \phi) | \mathbf{y}; \mathbf{q}^i \}. \quad (26)$$

The maximizing value of  $\phi$  in (26) is necessarily a value of  $\phi$  that zeroes the derivative, i.e., that satisfies

$$\sum_{n=1}^N \int_{x_n} p(x_n | \mathbf{y}; \mathbf{q}^i) \frac{d}{d\phi} \ln p_X(x_n; \lambda^i, \theta^i, \phi) = 0. \quad (27)$$

For the  $p_X(x_n; \lambda, \theta, \phi)$  given in (2), it is readily seen that

$$\begin{aligned} & \frac{d}{d\phi} \ln p_X(x_n; \lambda^i, \theta^i, \phi) \\ &= \frac{1}{2} \left( \frac{|x_n - \theta^i|^2}{(\phi)^2} - \frac{1}{\phi} \right) \frac{\lambda^i \mathcal{N}(x_n; \theta^i, \phi)}{p_X(x_n; \lambda^i, \theta^i, \phi)} \\ &= \begin{cases} \frac{1}{2} \left( \frac{|x_n - \theta^i|^2}{(\phi)^2} - \frac{1}{\phi} \right) & x_n \neq 0 \\ 0 & x_n = 0 \end{cases}. \end{aligned} \quad (29)$$

Splitting the domain of integration in (27) into  $\mathcal{B}_\epsilon$  and  $\overline{\mathcal{B}_\epsilon}$ , and then plugging in (29), we find that the following is equivalent to (27) in the limit of  $\epsilon \rightarrow 0$ :

$$\sum_{n=1}^N \int_{x_n \in \overline{\mathcal{B}_\epsilon}} (|x_n - \theta^i|^2 - \phi) p(x_n | \mathbf{y}; \mathbf{q}^i) = 0. \quad (30)$$

The unique value of  $\phi$  satisfying (30) as  $\epsilon \rightarrow 0$  is then

$$\phi^{i+1} = \frac{\sum_{n=1}^N \lim_{\epsilon \rightarrow 0} \int_{x_n \in \overline{\mathcal{B}_\epsilon}} |x_n - \theta^i|^2 p(x_n | \mathbf{y}; \mathbf{q}^i)}{\sum_{n=1}^N \lim_{\epsilon \rightarrow 0} \int_{x_n \in \overline{\mathcal{B}_\epsilon}} p(x_n | \mathbf{y}; \mathbf{q}^i)} \quad (31)$$

$$\begin{aligned} &= \frac{1}{\lambda^{i+1} N} \sum_{n=1}^N \pi(\hat{r}_n, \mu_n^r; \mathbf{q}^i) \left( |\theta^i - \gamma(\hat{r}_n, \mu_n^r; \mathbf{q}^i)|^2 \right. \\ &\quad \left. + \nu(\hat{r}_n, \mu_n^r; \mathbf{q}^i) \right) \end{aligned} \quad (32)$$

for the GAMP outputs  $\{\nu(\hat{r}_n, \mu_n^r; \mathbf{q}^i)\}_{n=1}^N$  defined in (10). The equality in (32) can be verified by plugging the GAMP posterior expression from (12) into (31) and simplifying using the Gaussian-pdf multiplication rule.

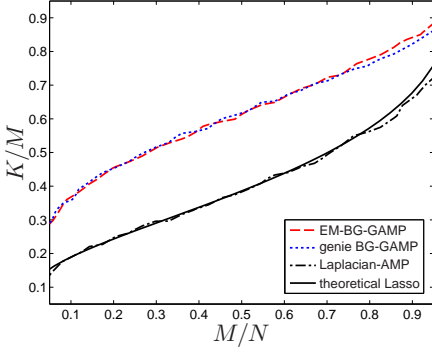


Fig. 1. Empirical noiseless PTCs for Bernoulli-Gaussian signals and theoretical PTC for Lasso.

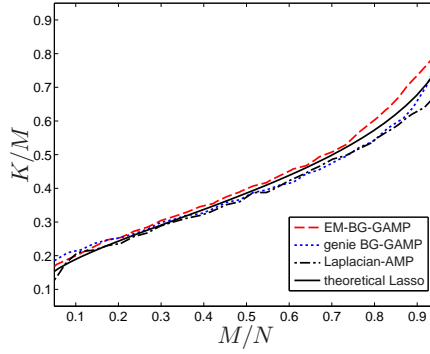


Fig. 2. Empirical noiseless PTCs for Bernoulli-Rademacher case and theoretical PTC for Lasso.

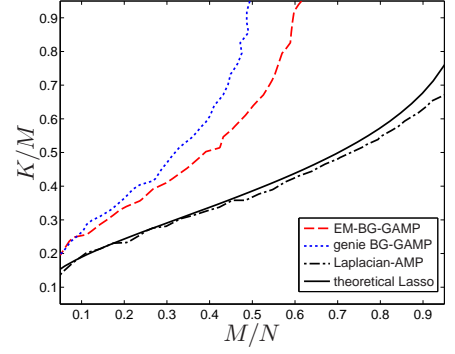


Fig. 3. Empirical noiseless PTCs for Bernoulli signals and theoretical PTC for Lasso.

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

Finally, we derive the EM update for  $\psi$  given previous parameters  $\mathbf{q}^i$ . Because  $\mathbf{w}$  is a priori independent of  $\mathbf{x}$  and i.i.d., the joint pdf  $p(\mathbf{x}, \mathbf{w}; \mathbf{q})$  decouples into  $C \prod_{m=1}^M p_W(w_m; \psi)$  for a  $\psi$ -invariant constant  $C$ , and so

$$\psi^{i+1} = \arg \max_{\psi > 0} \sum_{m=1}^M \mathbb{E} \{ \ln p_W(w_m; \psi) | \mathbf{y}; \mathbf{q}^i \}. \quad (33)$$

The maximizing value of  $\psi$  in (33) is necessarily a value of  $\psi$  that zeroes the derivative, i.e., that satisfies

$$\sum_{m=1}^M \int_{w_m} p(w_m | \mathbf{y}; \mathbf{q}^i) \frac{d}{d\psi} \ln p_W(w_m; \psi) = 0. \quad (34)$$

Because  $p_W(w_m; \psi) = \mathcal{N}(w_m; 0, \psi)$ , it is readily seen that

$$\frac{d}{d\psi} \ln p_W(w_m; \psi) = \frac{1}{2} \left( \frac{|w_m|^2}{(\psi)^2} - \frac{1}{\psi} \right), \quad (35)$$

which, when plugged into (34), yields the unique solution

$$\psi^{i+1} = \frac{1}{M} \sum_{m=1}^M \int_{w_m} |w_m|^2 p(w_m | \mathbf{y}; \mathbf{q}^i). \quad (36)$$

Since  $w_m = y_m - z_m$  for  $z_m \triangleq \mathbf{a}_m^\top \mathbf{x}$ , we can also write<sup>5</sup>

$$\psi^{i+1} = \frac{1}{M} \sum_{m=1}^M \int_{z_m} |y_m - z_m|^2 p(z_m | \mathbf{y}; \mathbf{q}^i) \quad (37)$$

$$= \frac{1}{M} \sum_{m=1}^M (|y_m - \hat{z}_m|^2 + \mu_m^z) \quad (38)$$

where  $\hat{z}_m$  and  $\mu_m^z$ , the posterior mean and variance of  $z_m$ , are available from GAMP (see steps (R1)-(R2) in Table I).

#### E. EM Initialization

Since the EM algorithm converges only to a local maximum of the likelihood function, proper initialization is essential. We initialize the sparsity as  $\lambda^0 = \frac{M}{N} \rho_{\text{SE}}(\frac{M}{N})$ , where  $\rho_{\text{SE}}(\frac{M}{N})$  is the sparsity ratio  $\frac{K}{M}$  achieved by the Lasso PTC [2]

$$\rho_{\text{SE}}(\frac{M}{N}) = \max_{a \geq 0} \frac{1 - \frac{2N}{M} [(1 + a^2)\Phi(a) - a\phi(a)]}{1 - a^2 - 2[(1 + a^2)\Phi(a) - a\phi(a)]}. \quad (39)$$

<sup>5</sup>Empirically, we have observed that the EM update for  $\psi$  works better with the  $\mu_m^z$  term in (38) weighted by  $\frac{M}{N}$  and suppressed until later EM iterations. We conjecture that this is due to bias in the GAMP variance estimates  $\mu_m^z$ .

We initialize the active mean as  $\theta^0 = 0$ , which effectively assumes that the active pdf  $f(\cdot)$  is symmetric. Finally, noting that  $\mathbb{E}\{\|\mathbf{y}\|_2^2\} = (\text{SNR} + 1)M\psi$  for  $\text{SNR} \triangleq \text{tr}(\mathbf{A}^\top \mathbf{A})\lambda\phi/(M\psi)$ , we see that the variances,  $\phi$  and  $\psi$ , can be initialized based on  $\|\mathbf{y}\|_2^2$  and a given hypothesis  $\text{SNR}^0 \geq 0$ . In particular,

$$\psi^0 = \frac{\|\mathbf{y}\|_2^2}{(\text{SNR}^0 + 1)M}, \quad \phi^0 = \frac{\|\mathbf{y}\|_2^2 - M\psi^0}{\text{tr}(\mathbf{A}^\top \mathbf{A})\lambda^0}, \quad (40)$$

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

## IV. NUMERICAL RESULTS

### A. Noiseless Phase Transitions

First, we describe the results of experiments that computed noiseless empirical phase transition curves (PTCs) under various sparse-signal distributions. To compute 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  $M \times N$  i.i.d.-Gaussian measurement matrix  $\mathbf{A}$ . 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 considered a success (i.e.,  $S_r = 1$ ) if  $\text{NMSE} \triangleq \|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 / \|\mathbf{x}\|_2^2 < 10^{-4}$ , where the average success rate is 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 three recovery algorithms: the proposed EM-BG-GAMP algorithm,<sup>6</sup> a “genie-aided” BG-GAMP that knew the true  $[\lambda, \theta, \phi, \psi]$ , and the Laplacian-AMP from [2]. For comparison, Figs. 1–3 also display the theoretical Lasso PTC (39). The signal was generated as BG with zero mean ( $\theta = 0$ ) and unit variance ( $\phi = 1$ ) in Fig. 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 or, equivalently, BG with  $\theta = 1$  and  $\phi = 0$ ).

Figures 1–3 demonstrate that, for all three signal types, the empirical PTC of EM-BG-GAMP improves on that for Laplacian-AMP as well as the theoretical Lasso PTC. (The latter two are known to converge in the large system limit

<sup>6</sup>Matlab code available at <http://www.ece.osu.edu/~schniter/EMturboGAMP>

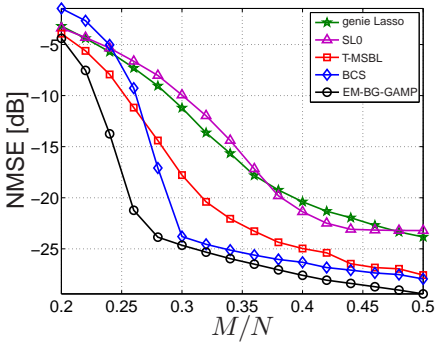


Fig. 4. NMSE for noisy recovery of a Bernoulli-Gaussian signal.

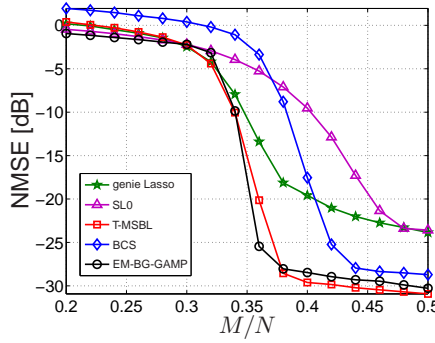


Fig. 5. NMSE for noisy recovery of a Bernoulli-Rademacher signal.

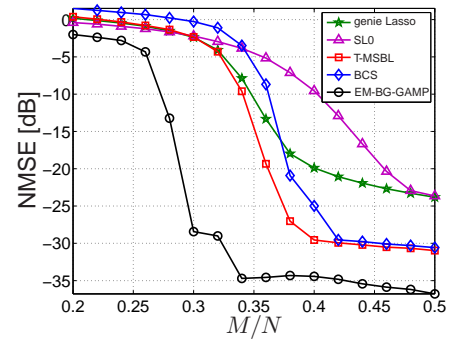


Fig. 6. NMSE for noisy recovery of a Bernoulli signal.

[2].) The smallest gains over Lasso appear when the signal is BR (i.e., the least-favorable distribution [3]), whereas the largest gains appear when the signal is Bernoulli. Amazingly, EM-BG-AMP perfectly recovered almost every Bernoulli realization when  $\frac{M}{N} \geq 0.65$ . The figures suggest that the EM algorithm does a decent job of learning the parameters  $\lambda, \theta, \phi$ . In fact, EM-BG-GAMP slightly outperforms genie-BG-GAMP in Figs. 1–2 due to realization-specific data fitting.

### B. Noisy Signal Recovery

Figures 4–6 show NMSE for noisy recovery of the same three sparse signal types considered in Figs. 1–3. To construct these new 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-BG-GAMP, Bayesian Compressive Sensing (BCS) [9], Sparse Bayesian Learning [8] (via T-MSBL), debiased genie-aided<sup>7</sup> Lasso (via SPGL1 [12]), and Smoothed- $\ell_0$  (SL0) [13]. All algorithms were run under the suggested defaults, with `'noise' = 'small'` in T-MSBL.

In Fig. 4 and Fig. 6, we see EM-BG-GAMP outperforming all other algorithms for all meaningful values of undersampling ratio  $\frac{M}{N}$ . In fact, for Bernoulli signals (Fig. 6), we see a significant improvement, especially when  $\frac{M}{N} \in [0.3, 0.38]$ . We have verified (in simulations not shown here) that similar behavior persists at lower SNRs. In Fig. 5, we see EM-BG-GAMP outperforming all algorithms except T-MSBL, which does  $\approx 1$  dB better for large values of  $\frac{M}{N}$ . Apparently, the prior assumed by T-MSBL is a better fit to BR than the BG prior.

Admittedly, the near-dominant EM-BG-GAMP performance observed for *perfectly* sparse signals in Figs. 1–6 does not hold for all signal classes. As an example, Fig. 7 shows noisy recovery NMSE for a Student’s-t signal with pdf

$$p_X(x; q) \triangleq \frac{\Gamma((q+1)/2)}{\sqrt{2\pi}\Gamma(q/2)} (1 + x^2)^{-(q+1)/2} \quad (41)$$

under the *non-compressible* parameter choice  $q = 1.67$  [14]. There, we see EM-BG-GAMP outperformed by SL0 and genie-aided Lasso, although not by T-MSBL and BCS. In fact, among the competing algorithms, those that performed best for exactly sparse signals seem to do worst for this non-compressible signal, and vice versa. We attribute these

<sup>7</sup>We ran SPGL1 in ‘BPDN’ mode:  $\min_{\mathbf{x}} \|\mathbf{x}\|_1$  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.

behaviors to a poor fit between the assumed and actual signal priors, motivating future work on an EM *Gaussian-Mixture* GAMP with automatic selection of the mixture order.

### REFERENCES

- [1] R. Tibshirani, “Regression shrinkage and selection via the lasso,” *J. Roy. Statist. Soc. B*, vol. 58, no. 1, pp. 267 – 288, 1996.
- [2] D. L. Donoho, A. Maleki, and A. Montanari, “Message passing algorithms for compressed sensing,” *Proc. National Academy of Sciences*, vol. 106, pp. 18914–18919, Nov. 2009.
- [3] D. L. Donoho, A. Maleki, and A. Montanari, “The noise-sensitivity phase transition in compressed sensing,” *arXiv:1004.1218*, Apr. 2010.
- [4] 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.
- [5] M. Bayati and A. Montanari, “The dynamics of message passing on dense graphs, with applications to compressed sensing,” *IEEE Trans. Inform. Theory*, vol. 57, pp. 764–785, Feb. 2011.
- [6] S. Rangan, “Generalized approximate message passing for estimation with random linear mixing,” *arXiv:1010.5141*, Oct. 2010.
- [7] M. E. Tipping, “Sparse Bayesian learning and the relevance vector machine,” *J. Machine Learning Research*, vol. 1, pp. 211–244, 2001.
- [8] D. Wipf and B. Rao, “Sparse Bayesian learning for basis selection,” *IEEE Trans. Signal Process.*, vol. 52, pp. 2153 – 2164, Aug. 2004.
- [9] S. Ji, Y. Xue, and L. Carin, “Bayesian compressive sensing,” *IEEE Trans. Signal Process.*, vol. 56, pp. 2346–2356, June 2008.
- [10] 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.
- [11] 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.
- [12] 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.
- [13] 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.
- [14] V. Cevher, “Learning with compressible priors,” in *Proc. Neural Inform. Process. Syst. Conf.*, (Vancouver, B.C.), Dec. 2009.

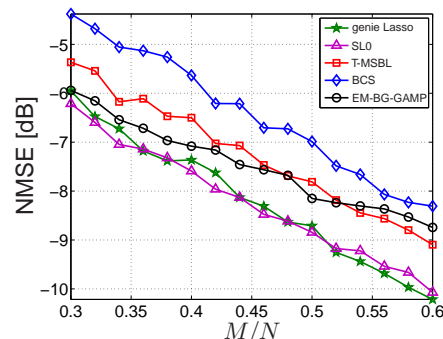


Fig. 7. NMSE for noisy recovery of a non-compressible Student’s-t signal.