

# Binary Linear Classification and Feature Selection via Generalized Approximate Message Passing

Justin Ziniel, *Student Member, IEEE*, Philip Schniter, *Fellow, IEEE*, and Per Sederberg

**Abstract**—For the problem of binary linear classification and feature selection, we propose algorithmic approaches to classifier design based on the generalized approximate message passing (GAMP) algorithm, recently proposed in the context of compressive sensing. We are particularly motivated by problems where the number of features greatly exceeds the number of training examples, but where only a few features suffice for accurate classification. We show that sum-product GAMP can be used to (approximately) minimize the classification error rate and max-sum GAMP can be used to minimize a wide variety of regularized loss functions. Furthermore, we describe an expectation-maximization (EM)-based scheme to learn the associated model parameters online, as an alternative to cross-validation, and we show that GAMP’s state-evolution framework can be used to accurately predict the misclassification rate. Finally, we present a detailed numerical study to confirm the accuracy, speed, and flexibility afforded by our GAMP-based approaches to binary linear classification and feature selection.

**Index Terms**—Belief propagation, classification, feature selection, message passing, one-bit compressed sensing.

## I. INTRODUCTION

IN this work we consider binary linear classification and feature selection [1]. The objective of *binary linear classification* is to learn the weight vector  $\mathbf{w} \in \mathbb{R}^N$  that best predicts an unknown binary class label  $y \in \{-1, 1\}$  associated with a given vector of quantifiable features  $\mathbf{x} \in \mathbb{R}^N$  from the sign of a linear “score”  $z \triangleq \langle \mathbf{x}, \mathbf{w} \rangle$ .<sup>1</sup> The goal of *linear feature selection* is to identify which subset of the  $N$  weights in  $\mathbf{w}$  are necessary for accurate prediction of the unknown class label  $y$ , since in some

applications (e.g., multi-voxel pattern analysis) this subset itself is of primary concern.

In formulating this linear feature selection problem, we assume that there exists a  $K$ -sparse weight vector  $\mathbf{w}$  (i.e.,  $\|\mathbf{w}\|_0 = K \ll N$ ) such that  $y = \text{sgn}(\langle \mathbf{x}, \mathbf{w} \rangle - \epsilon)$ , where  $\text{sgn}(\cdot)$  is the signum function and  $\epsilon \sim p_\epsilon$  is a random perturbation accounting for model inaccuracies. For the purpose of learning  $\mathbf{w}$ , we assume the availability of  $M$  labeled training examples generated independently according to this model:

$$y_m = \text{sgn}(\langle \mathbf{x}_m, \mathbf{w} \rangle - \epsilon_m), \quad \forall m = 1, \dots, M, \quad (1)$$

with  $\epsilon_m \sim \text{i.i.d. } p_\epsilon$ . It is common to express the relationship between the label  $y_m$  and the score  $z_m \triangleq \langle \mathbf{x}_m, \mathbf{w} \rangle$  in (1) via the conditional pdf  $p_{y_m|z_m}(y_m|z_m)$ , known as the “activation function,” which can be related to the perturbation pdf  $p_\epsilon$  via

$$p_{y_m|z_m}(1|z_m) = \int_{-\infty}^{z_m} p_\epsilon(\epsilon) d\epsilon = 1 - p_{y_m|z_m}(-1|z_m). \quad (2)$$

We are particularly interested in classification problems in which the number of potentially discriminatory features  $N$  drastically exceeds the number of available training examples  $M$ . Such computationally challenging problems are of great interest in a number of modern applications, including text classification [2], multi-voxel pattern analysis (MVPA) [3]–[5], conjoint analysis [6], and micro-array gene expression [7]. In MVPA, for instance, neuro-scientists attempt to infer which regions in the human brain are responsible for distinguishing between two cognitive states by measuring neural activity via fMRI at  $N \sim 10^4$  voxels. Due to the expensive and time-consuming nature of working with human subjects, classifiers are routinely trained using only  $M \sim 10^2$  training examples, and thus  $N \gg M$ .

In the  $N \gg M$  regime, the model of (1) coincides with that of *noisy one-bit compressed sensing* (CS) [8], [9]. In that setting, it is typical to write (1) in matrix-vector form using  $\mathbf{y} \triangleq [y_1, \dots, y_M]^T$ ,  $\mathbf{e} \triangleq [\epsilon_1, \dots, \epsilon_M]^T$ ,  $\mathbf{X} \triangleq [\mathbf{x}_1, \dots, \mathbf{x}_M]^T$ , and element-wise  $\text{sgn}(\cdot)$ , yielding

$$\mathbf{y} = \text{sgn}(\mathbf{X}\mathbf{w} - \mathbf{e}), \quad (3)$$

where  $\mathbf{w}$  embodies the signal-of-interest’s sparse representation,  $\mathbf{X} = \Phi\Psi$  is a concatenation of a linear measurement operator  $\Phi$  and a sparsifying signal dictionary  $\Psi$ , and  $\mathbf{e}$  is additive noise.<sup>2</sup> Importantly, in the  $N \gg M$  setting, [9] established performance guarantees on the estimation of  $K$ -sparse  $\mathbf{w}$  from  $O(K \log N/K)$  binary measurements of the form (3), under i.i.d. Gaussian  $\{\mathbf{x}_m\}$  and mild conditions on the perturbation process  $\{\epsilon_m\}$ , even when the entries within  $\mathbf{x}_m$  are correlated.

<sup>2</sup>For example, the common case of additive white Gaussian noise (AWGN)  $\{\epsilon_m\} \sim \text{i.i.d. } \mathcal{N}(0, v)$  corresponds to the “probit” activation function, i.e.,  $p_{y_m|z_m}(1|z_m) = \Phi(z_m/v)$ , where  $\Phi(\cdot)$  is the standard-normal cdf.

Manuscript received August 22, 2014; revised December 15, 2014; accepted January 29, 2015. Date of publication February 26, 2015; date of current version March 13, 2015. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Fauzia Ahmad. This work was supported by NSF grant CCF-1218754, NSF grant CCF-1018368, DARPA/ONR grant N66001-10-1-4090, and an allocation of computing time from the Ohio Supercomputer Center.

Portions of this work were presented at the Workshop on Information Theory and its Applications, San Diego, CA, USA, 2013 and the 2014 Conference on Information Sciences and Systems, Princeton, NJ, USA, 2014. (*Corresponding Author: P. Schniter.*)

J. Ziniel and P. Schniter are with the Department of Electrical and Computer Engineering, The Ohio State University, Columbus, OH 43210, USA (e-mail: schniter@ece.osu.edu).

P. Sederberg is with the Department of Psychology, The Ohio State University, Columbus, OH 43210 USA.

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TSP.2015.2407311

<sup>1</sup>We note that one could also compute the score from a fixed non-linear transformation  $\psi(\cdot)$  of the original feature  $\mathbf{x}$  via  $z \triangleq \langle \psi(\mathbf{x}), \mathbf{w} \rangle$  as in kernel-based classification. Although the methods we describe here are directly compatible with this approach, we write  $z = \langle \mathbf{x}, \mathbf{w} \rangle$  for simplicity.

This result implies that, in large binary linear classification problems, accurate feature selection is indeed possible from  $M \ll N$  training examples, as long as the underlying weight vector  $\mathbf{w}$  is sufficiently sparse. Not surprisingly, many techniques have been proposed to find such weight vectors [10]–[17].

In addition to theoretical analyses, the CS literature also offers a number of high-performance algorithms for the inference of  $\mathbf{w}$  in (3), e.g., [8], [9], [18]–[21]. Thus, the question arises as to whether these algorithms also show advantages in the domain of binary linear classification and feature selection. In this paper, we answer this question in the affirmative by focusing on the *generalized approximate message passing* (GAMP) algorithm [22], which extends the AMP algorithm [23], [24] from the case of linear, AWGN-corrupted observations (i.e.,  $\mathbf{y} = \mathbf{X}\mathbf{w} - \mathbf{e}$  for  $\mathbf{e} \sim \mathcal{N}(\mathbf{0}, v\mathbf{I})$ ) to the case of generalized-linear observations, such as (3). AMP and GAMP are attractive for several reasons: (i) For i.i.d sub-Gaussian  $\mathbf{X}$  in the large-system limit (i.e.,  $M, N \rightarrow \infty$  with fixed ratio  $\delta = \frac{M}{N}$ ), they are rigorously characterized by a state-evolution whose fixed points, when unique, are optimal [25]; (ii) Their state-evolutions predict fast convergence rates; (iii) They are very flexible with regard to data-modeling assumptions (see, e.g., [26]); (iv) Their model parameters can be learned online using an expectation-maximization (EM) approach that has been shown to yield state-of-the-art mean-squared reconstruction error in CS problems [27].

In this work, we develop a GAMP-based approach to binary linear classification and feature selection that makes the following contributions: 1) in Section II, we show that GAMP implements a particular approximation to the error-rate minimizing linear classifier under the assumed model (1); 2) in Section III, we show that GAMP’s state evolution framework can be used to characterize the misclassification rate in the large-system limit; 3) in Section IV, we develop methods to implement logistic, probit, and hinge-loss-based regression using both max-sum and sum-product versions of GAMP, and we further develop a method to make these classifiers robust in the face of corrupted training labels; and 4) in Section V, we present an EM-based scheme to learn the model parameters online, as an alternative to cross-validation. The numerical study presented in Section VI then confirms the efficacy, flexibility, and speed afforded by our GAMP-based approaches to binary classification and feature selection.

*Notation.* Random quantities are typeset in sans-serif (e.g.,  $e$ ) while deterministic quantities are typeset in serif (e.g.,  $e$ ). The pdf of random variable  $e$  under deterministic parameters  $\boldsymbol{\theta}$  is written as  $p_e(e; \boldsymbol{\theta})$ , where the subscript and parameterization are sometimes omitted for brevity. Column vectors are typeset in boldface lower-case (e.g.,  $\mathbf{y}$  or  $\mathbf{y}$ ), matrices in boldface upper-case (e.g.,  $\mathbf{X}$  or  $\mathbf{X}$ ), and their transpose is denoted by  $(\cdot)^T$ . For vector  $\mathbf{y} = [y_1, \dots, y_N]^T$ ,  $\mathbf{y}_{m:n}$  refers to the subvector  $[y_m, \dots, y_n]^T$ . Finally,  $\mathcal{N}(\mathbf{a}; \mathbf{b}, \mathbf{C})$  is the multivariate normal distribution as a function of  $\mathbf{a}$ , with mean  $\mathbf{b}$ , and with covariance matrix  $\mathbf{C}$ , while  $\phi(\cdot)$  and  $\Phi(\cdot)$  denote the standard normal pdf and cdf, respectively.

## II. GAMP FOR CLASSIFICATION

In this section, we introduce generalized approximate message passing (GAMP) from the perspective of binary linear classification. In particular, we show that the *sum-product* variant of

GAMP is a loopy belief propagation (LBP) approximation of the classification-error-rate minimizing linear classifier and that the *max-sum* variant of GAMP is a LBP implementation of the standard regularized-loss-minimization approach to linear classifier design.

### A. Sum-Product GAMP

Suppose that we are given  $M$  labeled training examples  $\{y_m, \mathbf{x}_m\}_{m=1}^M$ , and  $T$  test feature vectors  $\{\mathbf{x}_t\}_{t=M+1}^{M+T}$  associated with unknown test labels  $\{y_t\}_{t=M+1}^{M+T}$ , all obeying the noisy linear model (1) under some known error pdf  $p_e$ , and thus known  $p_{y_m|z_m}$ . We then consider the problem of computing the classification-error-rate minimizing hypotheses  $\{\hat{y}_t\}_{t=M+1}^{M+T}$ ,

$$\hat{y}_t = \arg \max_{y_t \in \{-1, 1\}} p_{y_t|\mathbf{y}_{1:M}}(y_t | \mathbf{y}_{1:M}; \mathbf{X}), \quad (4)$$

with  $\mathbf{y}_{1:M} \triangleq [y_1, \dots, y_M]^T$  and  $\mathbf{X} \triangleq [\mathbf{x}_1, \dots, \mathbf{x}_{M+T}]^T$ . Note that we treat the labels  $\{y_m\}_{m=1}^{M+T}$  as random but the features  $\{\mathbf{x}_m\}_{m=1}^{M+T}$  as deterministic parameters. The probabilities in (4) can be computed via the marginalization

$$\begin{aligned} p_{y_t|\mathbf{y}_{1:M}}(y_t | \mathbf{y}_{1:M}; \mathbf{X}) &= p_{y_t, \mathbf{y}_{1:M}}(y_t, \mathbf{y}_{1:M}; \mathbf{X}) C_{\mathbf{y}}^{-1} \\ &= C_{\mathbf{y}}^{-1} \sum_{\mathbf{y} \in \mathcal{Y}_t(y_t)} \int p_{\mathbf{y}, \mathbf{w}}(\mathbf{y}, \mathbf{w}; \mathbf{X}) d\mathbf{w} \end{aligned} \quad (5)$$

with scaling constant  $C_{\mathbf{y}} \triangleq p_{\mathbf{y}_{1:M}}(\mathbf{y}_{1:M}; \mathbf{X})$ , label vector  $\mathbf{y} = [y_1, \dots, y_{M+T}]^T$ , and constraint set

$$\begin{aligned} \mathcal{Y}_t(y) &\triangleq \{\tilde{\mathbf{y}} \in \{-1, 1\}^{M+T} \text{ s.t. } [\tilde{\mathbf{y}}]_t = y \\ &\quad \text{and } [\tilde{\mathbf{y}}]_m = y_m \forall m = 1, \dots, M\} \end{aligned}$$

which fixes the  $t$ th element of  $\mathbf{y}$  at the value  $y$  and the first  $M$  elements of  $\mathbf{y}$  at the values of the corresponding training labels. The joint pdf in (6) factors as

$$p_{\mathbf{y}, \mathbf{w}}(\mathbf{y}, \mathbf{w}; \mathbf{X}) = \prod_{m=1}^{M+T} p_{y_m|z_m}(y_m | \mathbf{x}_m^T \mathbf{w}) \prod_{n=1}^N p_{w_n}(w_n) \quad (7)$$

due to the model (1) and assuming a separable prior, i.e.,

$$p_{\mathbf{w}}(\mathbf{w}) = \prod_{n=1}^N p_{w_n}(w_n). \quad (8)$$

Although the separability assumption can also be relaxed (see, e.g., [26], [28]), we do not consider such extensions in this work.

The factorization (7) is illustrated using the *factor graph* in Fig. 1(a), which connects the various random variables to the pdf factors in which they appear. Although exact computation of the marginal posterior test-label probabilities via (6) is computationally intractable due to the high-dimensional summation and integration, the factor graph in Fig. 1(a) suggests the use of loopy belief propagation (LBP) [29], and in particular the *sum-product algorithm* (SPA) [30], as a tractable way to approximate these marginal probabilities. Although the SPA guarantees exact marginal posteriors only under non-loopy (i.e., tree-structured graphs), it has proven successful in many applications with

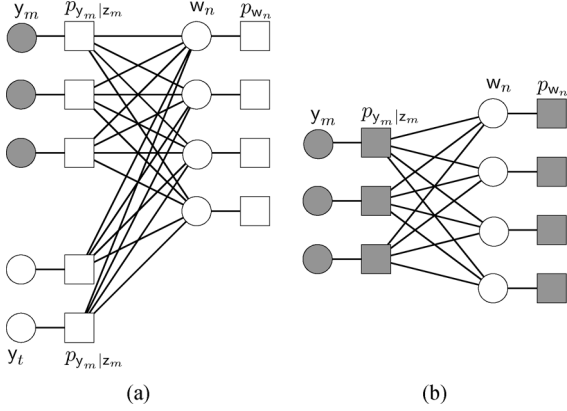


Fig. 1. Factor graph representations of the integrand of (7), with white/grey circles denoting unobserved/observed random variables, and rectangles denoting pdf “factors”. (a) Full. (b) Reduced.

loopy graphs, such as turbo decoding [31], computer vision [32], and compressive sensing [22]–[24].

Because a direct application of the SPA to the factor graph in Fig. 1(a) is itself computationally infeasible in the high-dimensional case of interest, we turn to a recently developed approximation: the sum-product variant of GAMP [22], as specified in Algorithm 1. The GAMP algorithm is specified in Algorithm 1 for a given instantiation of  $\mathbf{X}$ ,  $p_{y|z}$ , and  $\{p_{w_n}\}$ . There, the expectation and variance in lines 5–6 and 16–17 are taken element-wise w.r.t the GAMP-approximated marginal posterior pdfs (with superscript  $k$  denoting the iteration)

$$q(z_m | \hat{p}_m^k, \tau_{p_m}^k) = p_{y_m|z_m}(y_m|z_m) \mathcal{N}(z_m; \hat{p}_m^k, \tau_{p_m}^k) C_z^{-1} \quad (9)$$

$$q(w_n | \hat{r}_n^k, \tau_{r_n}^k) = p_{w_n}(w_n) \mathcal{N}(w_n; \hat{r}_n^k, \tau_{r_n}^k) C_w^{-1} \quad (10)$$

with appropriate normalizations  $C_z$  and  $C_w$ , and the vector-vector multiplications and divisions in lines 3, 9, 11, 12, 14, 13, 20 are performed element-wise. Due to space limitations, we refer the interested reader to [22] for an overview and derivation of GAMP, to [25] for rigorous analysis under large i.i.d sub-Gaussian  $\mathbf{X}$ , and to [33], [34] for fixed-point and local-convergence analysis under arbitrary  $\mathbf{X}$ .

Applying GAMP to the classification factor graph in Fig. 1(a) and examining the resulting form of lines 5–6 in Algorithm 1, it becomes evident that the test-label nodes  $\{y_t\}_{t=M+1}^{M+T}$  do not affect the GAMP weight estimates ( $\hat{\mathbf{w}}^k, \tau_w^k$ ) and thus the factor graph can effectively be simplified to the form shown in Fig. 1(b), after which the (approximated) posterior test-label pdfs are computed via

$$p_{y_t|y_{1:M}}(y_t|\mathbf{y}_{1:M}; \mathbf{X}) \approx \int p_{y_t|z_t}(y_t|z_t) \mathcal{N}(z_t; \hat{z}_t^\infty, \tau_{z_t}^\infty) dz_t \quad (11)$$

where  $\hat{z}_t^\infty$  and  $\tau_{z_t}^\infty$  denote the  $t^{\text{th}}$  element of the GAMP vectors  $\hat{\mathbf{z}}^k$  and  $\tau_z^k$ , respectively, at the final iteration “ $k = \infty$ .”

### B. Max-Sum GAMP

An alternate approach to linear classifier design is through the minimization of a regularized loss function, e.g.,

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^N} \sum_{m=1}^M f_{z_m}(\mathbf{x}_m^\top \mathbf{w}) + \sum_{n=1}^N f_{w_n}(w_n), \quad (12)$$

### Algorithm 1: Generalized Approximate Message Passing

**Input:** Matrix  $\mathbf{X}$ , priors  $p_{w_n}(\cdot)$ , activation functions  $p_{y_m|z_m}(y_m|\cdot)$ , and mode  $\in \{\text{SumProduct}, \text{MaxSum}\}$

**Initialize:**  $k \leftarrow 0$ ;  $\hat{\mathbf{s}}^{-1} \leftarrow \mathbf{0}$ ;  $\mathbf{S} \leftarrow |\mathbf{X}|^2$ ;  $\hat{\mathbf{w}}^0 \leftarrow \mathbf{0}$ ;  $\tau_w^0 \leftarrow 1$

- 1: **repeat**
- 2:    $\tau_p^k \leftarrow \mathbf{S} \tau_w^k$
- 3:    $\hat{\mathbf{p}}^k \leftarrow \mathbf{X} \hat{\mathbf{w}}^k - \hat{\mathbf{s}}^{k-1} \tau_p^k$
- 4:   **if** SumProduct **then**
- 5:      $\hat{\mathbf{z}}^k \leftarrow \mathbb{E}\{\mathbf{z} | \hat{\mathbf{p}}^k, \tau_p^k\}$
- 6:      $\tau_z^k \leftarrow \text{var}\{\mathbf{z} | \hat{\mathbf{p}}^k, \tau_p^k\}$
- 7:   **else if** MaxSum **then**
- 8:      $\hat{\mathbf{z}}^k \leftarrow \text{prox}_{\tau_p^k f_{z_m}}(\hat{\mathbf{p}}^k)$
- 9:      $\tau_z^k \leftarrow \tau_p^k \text{prox}'_{\tau_p^k f_{z_m}}(\hat{\mathbf{p}}^k)$
- 10:   **end if**
- 11:    $\tau_s^k \leftarrow 1/\tau_p^k - \tau_z^k/(\tau_p^k)^2$
- 12:    $\hat{\mathbf{s}}^k \leftarrow (\hat{\mathbf{z}}^k - \hat{\mathbf{p}}^k)/\tau_p^k$
- 13:    $\tau_r^k \leftarrow 1/(\mathbf{S}^\top \tau_s^k)$
- 14:    $\hat{\mathbf{r}}^k \leftarrow \hat{\mathbf{w}}^k + \tau_r^k \mathbf{X}^\top \hat{\mathbf{s}}^k$
- 15:   **if** SumProduct **then**
- 16:      $\hat{\mathbf{w}}^{k+1} \leftarrow \mathbb{E}\{\mathbf{w} | \hat{\mathbf{r}}^k, \tau_r^k\}$
- 17:      $\tau_w^{k+1} \leftarrow \text{var}\{\mathbf{w} | \hat{\mathbf{r}}^k, \tau_r^k\}$
- 18:   **else if** MaxSum **then**
- 19:      $\hat{\mathbf{w}}^{k+1} \leftarrow \text{prox}_{\tau_r^k f_{w_n}}(\hat{\mathbf{r}}^k)$
- 20:      $\tau_w^{k+1} \leftarrow \tau_r^k \text{prox}'_{\tau_r^k f_{w_n}}(\hat{\mathbf{r}}^k)$
- 21:   **end if**
- 22:    $k \leftarrow k + 1$
- 23: **until** Terminated

where  $f_{z_m}(\cdot)$  are  $y_m$ -dependent convex loss functions (e.g., logistic, probit, or hinge based) and where  $f_{w_n}(\cdot)$  are convex regularization terms (e.g.,  $f_{w_n}(w) = \lambda w^2$  for  $\ell_2$  regularization and  $f_{w_n}(w) = \lambda|w|$  for  $\ell_1$  regularization).

The solution to (12) can be recognized as the *maximum a posteriori* (MAP) estimate of random vector  $\mathbf{w}$  given a separable prior  $p_{\mathbf{w}}(\cdot)$  and likelihood corresponding to (1), i.e.,

$$p_{\mathbf{y}|\mathbf{w}}(\mathbf{y}|\mathbf{w}; \mathbf{X}) = \prod_{m=1}^M p_{y_m|z_m}(y_m|\mathbf{x}_m^\top \mathbf{w}), \quad (13)$$

when  $f_{z_m}(z) = -\log p_{y_m|z_m}(y_m|z)$  and  $f_{w_n}(w) = -\log p_{w_n}(w)$ . Importantly, this statistical model is exactly the one yielding the reduced factor graph in Fig. 1(b).

Similar to how sum-product LBP can be used to compute (approximate) marginal posteriors in loopy graphs, *max-sum* LBP can be used to compute the MAP estimate [35]. Since max-sum LBP is itself intractable for the high-dimensional problems of interest, we turn to the max-sum variant of GAMP [22], which is also specified in Algorithm 1. There, lines 8–9 are to be interpreted as

$$\hat{z}_m^k = \text{prox}_{\tau_p^k f_{z_m}}(\hat{p}_m^k), \quad m = 1, \dots, M, \quad (14)$$

$$\tau_{z_m}^k = \tau_{p_m}^k \text{prox}'_{\tau_{p_m}^k f_{z_m}}(\hat{p}_m^k), \quad m = 1, \dots, M, \quad (15)$$

with  $(\cdot)'$  and  $(\cdot)''$  denoting first and second derivatives and

$$\text{prox}_{\tau f}(v) \triangleq \arg \min_{u \in \mathbb{R}} \left[ f(u) + \frac{1}{2\tau}(u - v)^2 \right] \quad (16)$$

$$\text{prox}'_{\tau f}(v) = \left( 1 + \tau f''(\text{prox}_{\tau f}(v)) \right)^{-1}, \quad (17)$$

and lines 19–20 are to be interpreted similarly. It is known [33] that, for *arbitrary*  $\mathbf{X}$ , the fixed points of GAMP correspond to the critical points of the optimization objective (12).

### C. GAMP Summary

In summary, the sum-product and max-sum variants of the GAMP algorithm provide tractable methods of approximating the posterior test-label probabilities  $\{p_{y_t | \mathbf{y}_{1:M}}(y_t | \mathbf{y}_{1:M})\}_{t=T+1}^{M+T}$  and finding the MAP weight vector  $\hat{\mathbf{w}} = \arg \max_{\mathbf{w}} p_{\mathbf{w} | \mathbf{y}_{1:M}}(\mathbf{w} | \mathbf{y}_{1:M})$ , respectively, under the label-generation model (13) [equivalently, (1)] and the separable weight-vector prior (8), assuming that the distributions  $p_{y|z}$  and  $\{p_{w_n}\}$  are known and facilitate tractable scalar-nonlinear update steps 5–6, 8–9, 16–17, and 19–20. In Section IV, we discuss the implementation of these update steps for several popular activation functions, and in Section V, we discuss how the parameters of  $p_{y_m | z_m}$  and  $p_{w_n}$  can be learned online.

## III. MISCLASSIFICATION RATE VIA STATE EVOLUTION

As mentioned earlier, the behavior of GAMP in the large-system limit (i.e.,  $M, N \rightarrow \infty$  with fixed ratio  $\delta = \frac{M}{N}$ ) under i.i.d sub-Gaussian  $\mathbf{X}$  is characterized by a scalar state evolution [22], [25]. We now describe how this state evolution can be used to characterize the test-error rate of the linear-classification GAMP algorithms described in Section II.

The GAMP state evolution characterizes average GAMP performance over an ensemble of (infinitely sized) problems, each associated with one realization  $(\mathbf{y}, \mathbf{X}, \mathbf{w})$  of the random triple  $(\mathbf{y}, \mathbf{X}, \mathbf{w})$ . Recall that, for a given problem realization  $(\mathbf{y}, \mathbf{X}, \mathbf{w})$ , the GAMP iterations in Algorithm 1 yields the sequence of estimates  $\{\hat{\mathbf{w}}^k\}_{k=1}^{\infty}$  of the true weight vector  $\mathbf{w}$ . Then, according to the state evolution,  $p_{\mathbf{w}, \hat{\mathbf{w}}^k}(\mathbf{w}, \hat{\mathbf{w}}^k) \sim \prod_n p_{w_n, \hat{w}_n^k}(w_n, \hat{w}_n^k)$  and the first two moments of the joint pdf  $p_{w_n, \hat{w}_n^k}$  can be computed using [22, Algorithm 3].

Suppose that the  $(\mathbf{y}, \mathbf{X})$  above represent training examples associated with a true weight vector  $\mathbf{w}$ , and that  $(\mathbf{y}, \mathbf{x})$  represents a test pair also associated with the same  $\mathbf{w}$  and with  $\mathbf{x}$  having i.i.d elements distributed identically to those of  $\mathbf{X}$  (with, say, variance  $\frac{1}{M}$ ). The true and iteration- $k$ -estimated test scores are then  $z \triangleq \mathbf{x}^T \mathbf{w}$  and  $\hat{z}^k \triangleq \mathbf{x}^T \hat{\mathbf{w}}^k$ , respectively. The corresponding test-error rate<sup>3</sup>  $\mathcal{E}^k \triangleq \Pr\{y \neq \text{sgn}(\hat{z}^k)\}$  can be computed as follows. Letting  $I_{\{\cdot\}}$  denote an indicator function that assumes the value 1 when its Boolean argument is true and the value 0 otherwise, we have

$$\mathcal{E}^k = \mathbb{E}\{I_{\{y \neq \text{sgn}(\hat{z}^k)\}}\} \quad (18)$$

$$= \sum_{y \in \{-1, 1\}} \int I_{\{y \neq \text{sgn}(\hat{z}^k)\}} \int p_{y, \hat{z}^k, z}(y, \hat{z}^k, z) dz d\hat{z}^k \quad (19)$$

<sup>3</sup>For simplicity we assume a decision rule of the form  $\hat{y}^k = \text{sgn}(\hat{z}^k)$ , although other decision rules can be accommodated in our analysis.

$$= \sum_{y \in \{-1, 1\}} \iint I_{\{y \neq \text{sgn}(\hat{z}^k)\}} p_{y|z}(y|z) p_{z, \hat{z}^k}(z, \hat{z}^k) dz d\hat{z}^k. \quad (20)$$

Furthermore, from the definitions of  $(z, \hat{z}^k)$  and the bivariate central limit theorem, we have that

$$\begin{bmatrix} z \\ \hat{z}^k \end{bmatrix} \xrightarrow{d} \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_z^k) = \mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \Sigma_{11}^k & \Sigma_{12}^k \\ \Sigma_{21}^k & \Sigma_{22}^k \end{bmatrix}\right), \quad (21)$$

where  $\xrightarrow{d}$  indicates convergence in distribution. In [36], it is shown that the above matrix components are

$$\Sigma_{11}^k = \delta^{-1}(\text{var}\{w_n\} + \mathbb{E}[w_n]^2), \quad (22)$$

$$\Sigma_{12}^k = \Sigma_{21}^k = \delta^{-1}(\text{cov}\{w_n, \hat{w}_n^k\} + \mathbb{E}[w_n]\mathbb{E}[\hat{w}_n^k]), \quad (23)$$

$$\Sigma_{22}^k = \delta^{-1}(\text{var}\{\hat{w}_n^k\} + \mathbb{E}[\hat{w}_n^k]^2) \quad (24)$$

for label-to-feature ratio  $\delta$ . As described earlier, the above moments can be computed using [22, Algorithm 3]. The integral in (20) can then be computed (numerically if needed) for a given activation function  $p_{y|z}$ , yielding an estimate of GAMP's test-error rate at the  $k^{\text{th}}$  iteration.

To validate the accuracy of the above asymptotic analysis, we conducted a Monte-Carlo experiment with data synthetically generated in accordance with the assumed model. In particular, for each of 1000 problem realizations, a true weight vector  $\mathbf{w} \in \mathbb{R}^N$  was drawn i.i.d zero-mean Bernoulli-Gaussian and a feature matrix  $\mathbf{X}$  was drawn i.i.d Gaussian, yielding true scores  $z = \mathbf{X}\mathbf{w}$ , from which the true labels  $\mathbf{y}$  were randomly drawn using a probit activation function  $p_{y|z}$ . A GAMP weight-vector estimate  $\hat{\mathbf{w}}^\infty$  was then computed using the training data  $(\mathbf{y}_{1:M}, \mathbf{X}_{1:M})$ , from which the test-label estimates  $\{\hat{y}_t^\infty\}_{t=M+1}^{M+T}$  with  $\hat{y}_t^\infty = \text{sgn}(\mathbf{x}_t^T \hat{\mathbf{w}}^\infty)$  were computed and compared to the true test-labels in order to calculate the test-error rate for that realization. Fig. 2(a) plots the Monte-Carlo averaged empirical test-error rates (dashed) and state-evolution predicted rates (solid) as level curves over different combinations of training ratio  $\frac{M}{N}$  and discriminative-feature ratio  $\frac{K}{N}$ , where  $K = \|\mathbf{w}\|_0$  and  $N = 1024$ . Similarly, Fig. 2(b) plots average empirical mean-squared error (MSE) versus state-evolution predicted MSE, where  $\text{MSE} = \frac{1}{N} \mathbb{E}\{\|\hat{\mathbf{w}}^\infty - \mathbf{w}\|_2^2\}$ .

In both Fig. 2(a) and (b), the training-to-feature ratio  $\frac{M}{N}$  increases from left to right, and the discriminative-feature ratio  $\frac{K}{N}$  increases from bottom to top. The region to the upper-left of the dash-dotted black line contains ill-posed problems (where the number of discriminative features  $K$  exceeds the number of training samples  $M$ ) for which data was not collected. The remainders of Fig. 2(a) and (b) show very close agreement between empirical averages and state-evolution predictions.

## IV. GAMP NONLINEAR STEPS

Section II gave a high-level description of how the GAMP iterations in Algorithm 1 can be applied to binary linear classification and feature selection. In this section, we detail the nonlinear steps used to compute  $(\hat{z}, \tau_z)$  and  $(\hat{x}, \tau_x)$  in lines 5–6, 8–9, 16–17, and 19–20 of Algorithm 1. For sum-product GAMP, we recall that the mean and variance computations in lines 5–6

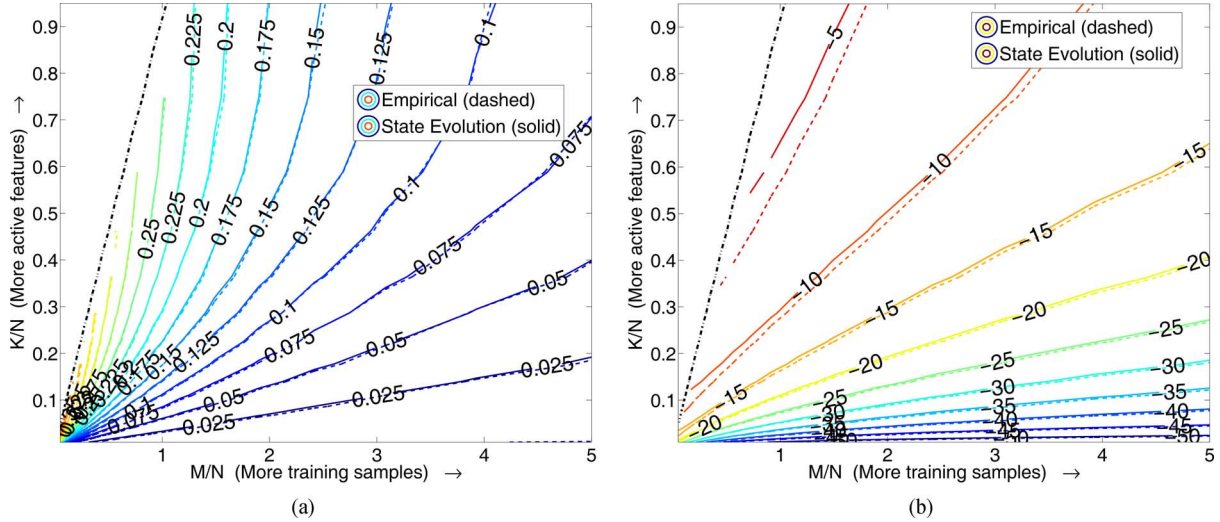


Fig. 2. Test-error rate (a) and weight-vector MSE (b), versus training-to-feature ratio  $M/N$  and discriminative-feature ratio  $K/N$ , calculated using empirical averaging (dashed) and state-evolution prediction (solid), assuming i.i.d Bernoulli-Gaussian weight vectors and a probit activation function. (a) Test-Error Rate. (b) Weight-Vector MSE (dB).

and 16–17 are computed based on the pdfs in (9) and (10), respectively, and for max-sum GAMP the prox steps in 8–9 are computed using (14)–(15) and those in 19–20 are computed similarly.

#### A. Logistic Activation Function

Arguably the most popular activation function for binary linear classification is the logistic sigmoid [1, §4.3.2],[37]:

$$p_{y|z}(y|z; \alpha) = \frac{1}{1 + \exp(-y\alpha z)}, \quad y \in \{-1, 1\} \quad (25)$$

where  $\alpha > 0$  controls the steepness of the transition.

For logistic sum-product GAMP, we propose to compute the mean and variance ( $\hat{z}, \tau_z$ ) of the marginal posterior approximation (9) using the variational approach in Algorithm 2, whose derivation is relegated to [36] for reasons of space. We note that Algorithm 2 is reminiscent of the one presented in [1, §10.6], but is more general in that it handles  $\alpha \neq 1$ .

For logistic max-sum GAMP,  $\hat{z}$  from (14) solves the scalar minimization problem (16) with  $f(u) = -\log p_{y|z}(y|u; \alpha)$  from (25), which is convex. To find this  $\hat{z}$ , we use bisection search to locate the root of  $\frac{d}{du}[f(u) + \frac{1}{2\tau}(u - v)^2]$ . The max-sum  $\tau_z$  from (15) can then be computed in closed form using  $\hat{z}$  and  $f''(\cdot)$  via (17). Note that, unlike the classical ML-based approach to logistic regression (e.g., [1, §4.3.3]), GAMP performs only scalar minimizations and thus does not need to construct or invert a Hessian matrix.

#### B. Probit Activation Function

Another popular activation function is the probit [1, §4.3.5]:

$$p_{y|z}(1|z; v) = \int_{-\infty}^z \mathcal{N}(\tau; 0, v) d\tau = \Phi\left(\frac{z}{\sqrt{v}}\right) \quad (26)$$

where  $p_{y|z}(-1|z) = 1 - p_{y|z}(1|z) = \Phi\left(-\frac{z}{\sqrt{v}}\right)$  and where  $v > 0$  controls the steepness of the sigmoid.

Unlike the logistic case, the probit case leads to closed-form sum-product GAMP computations. In particular, the density (9)

TABLE I  
SUM-PRODUCT GAMP COMPUTATIONS FOR PROBIT  
ACTIVATION FUNCTION

Quantity	Value
$c$	$\frac{\hat{p}}{\sqrt{v + \tau_p}}$
$\hat{z}$	$\hat{p} + \frac{y\tau_p\phi(c)}{\Phi(yc)\sqrt{v + \tau_p}}$
$\tau_z$	$\tau_p - \frac{\tau_p^2\phi(c)}{\Phi(yc)(v + \tau_p)} \left( yc + \frac{\phi(c)}{\Phi(c)} \right)$

#### Algorithm 2: A Variational Approach to Logistic Activation Functions for Sum-Product GAMP

**Input:** Class label  $y \in \{-1, 1\}$ , logistic scale  $\alpha$ , and GAMP-computed parameters  $\hat{p}$  and  $\tau_p$  (see (9))

**Initialize:**  $\xi \leftarrow \sqrt{\tau_p + |\hat{p}|^2}$

1: **repeat**

2:  $\sigma \leftarrow (1 + \exp(-\alpha\xi))^{-1}$

3:  $\lambda \leftarrow \frac{\alpha}{2\xi} \left( \sigma - \frac{1}{2} \right)$

4:  $\tau_z \leftarrow \tau_p(1 + 2\tau_p\lambda)^{-1}$

5:  $\hat{z} \leftarrow \tau_z(\hat{p}/\tau_p + \alpha y/2)$

6:  $\xi \leftarrow \sqrt{\tau_z + |\hat{z}|^2}$

7: **until** Terminated

8: **return**  $\hat{z}, \tau_z$

corresponds to the posterior pdf of a random variable  $z$  with prior  $\mathcal{N}(\hat{p}, \tau_p)$  from an observation  $y = y$  measured under the likelihood model (26). A derivation in [38, §3.9] provides the necessary expressions for these moments when  $y = 1$ , and a similar exercise tackles the  $y = -1$  case. For completeness, the sum-product computations are summarized in Table I. Max-sum GAMP computation of  $(\hat{z}, \tau_z)$  can be performed using a bisection search akin to that described in Section IV-A.

TABLE II  
SUM-PRODUCT GAMP COMPUTATIONS FOR THE HINGE-LOSS ACTIVATION  
FUNCTION. SEE APPENDIX A FOR DEFINITIONS OF  $\gamma_y$ ,  $\underline{\mu}_y$ ,  $\bar{\mu}_y$ ,  $\underline{v}_y$ ,  $\bar{v}_y$

Quantity	Value
$\hat{z}$	$(1 + \gamma_y)^{-1} \underline{\mu}_y + (1 + \gamma_y^{-1})^{-1} \bar{\mu}_y$
$\tau_z$	$(1 + \gamma_y)^{-1} (\underline{v}_y + \underline{\mu}_y^2) + (1 + \gamma_y^{-1})^{-1} (\bar{v}_y + \bar{\mu}_y^2) - \hat{z}^2$

TABLE III  
SUM-PRODUCT GAMP COMPUTATIONS FOR A ROBUSTIFIED ACTIVATION  
FUNCTION. SEE TEXT FOR DEFINITIONS OF  $C_y^*$ ,  $\hat{z}^*$ , AND  $\tau_z^*$

Quantity	Value
$C_y$	$\frac{\gamma}{\gamma + (1 - 2\gamma)C_y^*}$
$\hat{z}$	$C_y \hat{p} + (1 - C_y) \hat{z}^*$
$\tau_z$	$C_y (\tau_p + \hat{p}^2) + (1 - C_y) (\tau_z^* + (\hat{z}^*)^2) - \hat{z}^2$

### C. Hinge-Loss Activation Function

The hinge loss  $f_{z_m}(z) \triangleq \max(0, 1 - y_m z)$  is commonly used in the support vector machine (SVM) approach to maximum-margin classification [1, §7.1], i.e.,

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \sum_{m=1}^M f_{z_m}(\mathbf{x}_m^T \mathbf{w}) + \lambda \|\mathbf{w}\|_2^2 \quad (27)$$

or variations where  $\|\mathbf{w}\|_2^2$  is replaced with a sparsity-inducing alternative like  $\|\mathbf{w}\|_1$  [39]. Recalling Section II-B, this loss leads to the activation function

$$p_{y_m|z_m}(y_m|z) \propto \exp(-\max(0, 1 - y_m z)). \quad (28)$$

For hinge-loss sum-product GAMP, the mean and variance  $(\hat{z}, \tau_z)$  of (9) can be computed in closed form using the procedure described in Appendix A, and summarized in Table II. Meanwhile, for max-sum GAMP, the proximal steps (14)-(15) can be efficiently computed using bisection search, as in the logistic and probit cases.

### D. A Method to Robustify Activation Functions

In some applications, a fraction  $\gamma \in (0, 1)$  of the training labels are known<sup>4</sup> to be corrupted, or at least highly atypical under a given activation model  $p_{y|z}^*(y|z)$ . As a robust alternative to  $p_{y|z}^*(y|z)$ , Oppor and Winther [40] proposed to use

$$p_{y|z}(y|z; \gamma) = (1 - \gamma) p_{y|z}^*(y|z) + \gamma p_{y|z}^*(-y|z) \quad (29)$$

$$= \gamma + (1 - 2\gamma) p_{y|z}^*(y|z). \quad (30)$$

We now describe how the GAMP nonlinear steps for an arbitrary  $p_{y|z}^*$  can be used to compute the GAMP nonlinear steps for a robust  $p_{y|z}$  of the form in (30).

In the sum-product case, knowledge of the non-robust quantities  $\hat{z}^* \triangleq \frac{1}{C_y^*} \int_z z p_{y|z}^*(y|z) \mathcal{N}(z; \hat{p}, \tau_p)$ ,  $\tau_z^* \triangleq \frac{1}{C_y^*} \int_z (z - \hat{z}^*)^2 p_{y|z}^*(y|z) \mathcal{N}(z; \hat{p}, \tau_p)$ , and  $C_y^* \triangleq \int_z p_{y|z}^*(y|z) \mathcal{N}(z; \hat{p}, \tau_p)$  is sufficient for computing the robust sum-product quantities  $(\hat{z}, \tau_z)$ , as summarized in Table III. (See [36] for details.)

In the max-sum case, computing  $\hat{z}$  in (14) involves solving the scalar minimization problem in (16) with  $f(u) = -\log p_{y|z}(y|u; \gamma) = -\log[\gamma + (1 - 2\gamma) p_{y|z}^*(y|u)]$ .

<sup>4</sup>A method to learn an unknown  $\gamma$  will be proposed in Section V.

TABLE IV  
SUM-PRODUCT GAMP (SPG) AND MAX-SUM GAMP (MSG) COMPUTATIONS  
FOR THE ELASTIC-NET REGULARIZER  $f_{w_n}(w) = \lambda_1 |w| + \lambda_2 w^2$ ,  
WHICH INCLUDES  $\ell_1$  OR LAPLACIAN-PRIOR (VIA  $\lambda_2 = 0$ ) AND  $\ell_2$  OR  
GAUSSIAN-PRIOR (VIA  $\lambda_1 = 0$ ) AS SPECIAL CASES. SEE  
TABLE V FOR DEFINITIONS OF  $\underline{C}$ ,  $\bar{C}$ ,  $\underline{\mu}$ ,  $\bar{\mu}$ , ETC

	Quantity	Value
SPG	$\hat{w}$	$(\underline{C} \underline{\mu} + \bar{C} \bar{\mu}) / (\underline{C} + \bar{C})$
	$\tau_w$	$(\underline{C}(\underline{v} + \underline{\mu}^2) + \bar{C}(\bar{v} + \bar{\mu}^2)) / (\underline{C} + \bar{C}) - \hat{w}^2$
MSG	$\hat{w}$	$\text{sgn}(\sigma \bar{r}) \max( \sigma \bar{r}  - \lambda_1 \sigma^2, 0)$
	$\tau_w$	$\sigma^2 \cdot \mathbf{I}_{\{\hat{w} \neq 0\}}$

As before, we use a bisection search to find  $\hat{z}$  and then we use  $f''(\hat{z})$  to compute  $\tau_z$  via (17).

### E. Weight Vector Priors

We now discuss the nonlinear steps used to compute  $(\hat{w}, \tau_w)$ , i.e., lines 16–17 and 19–20 of Algorithm 1. These steps are, in fact, identical to those used to compute  $(\hat{z}, \tau_z)$  except that the prior  $p_{w_n}(\cdot)$  is used in place of the activation function  $p_{y_m|z_m}(y_m|\cdot)$ . For linear classification and feature selection in the  $N \gg M$  regime, it is customary to choose a prior  $p_{w_n}(\cdot)$  that leads to sparse (or approximately sparse) weight vectors  $\mathbf{w}$ , as discussed below.

For sum-product GAMP, this can be accomplished by choosing a Bernoulli- $\tilde{p}$  prior, i.e.,

$$p_{w_n}(w) = (1 - \pi_n) \delta(w) + \pi_n \tilde{p}_{w_n}(w), \quad (31)$$

where  $\delta(\cdot)$  is the Dirac delta function,  $\pi_n \in [0, 1]$  is the prior<sup>5</sup> probability that  $w_n = 0$ , and  $\tilde{p}_{w_n}(\cdot)$  is the pdf of a non-zero  $w_n$ . While Bernoulli-Gaussian [28] and Bernoulli-Gaussian-mixture [27] are common choices, Section VI suggests that Bernoulli-Laplacian also performs well.

In the max-sum case, the GAMP nonlinear outputs  $(\hat{w}, \tau_w)$  are computed via

$$\hat{w} = \text{prox}_{\tau_r f_{w_n}}(\hat{r}) \quad (32)$$

$$\tau_w = \tau_r \text{prox}'_{\tau_r f_{w_n}}(\hat{r}) \quad (33)$$

for a suitably chosen regularizer  $f_{w_n}(w)$ . Common examples include  $f_{w_n}(w) = \lambda_1 |w|$  for  $\ell_1$  regularization [23],  $f_{w_n}(w) = \lambda_2 w^2$  for  $\ell_2$  regularization [22], and  $f_{w_n}(w) = \lambda_1 |w| + \lambda_2 w^2$  for the “elastic net” [41]. As described in Section II-B, any regularizer  $f_{w_n}$  can be interpreted as a (possibly improper) prior pdf  $p_{w_n}(w) \propto \exp(-f_{w_n}(w))$ . Thus,  $\ell_1$  regularization corresponds to a Laplacian prior,  $\ell_2$  to a Gaussian prior, and the elastic net to a product of Laplacian and Gaussian pdfs.

In Table VII, we give the sum-product and max-sum computations for the prior corresponding to the elastic net, which includes both Laplacian (i.e.,  $\ell_1$ ) and Gaussian (i.e.,  $\ell_2$ ) as special cases; a full derivation can be found in [36]. For the Bernoulli-Laplacian case, these results can be combined with the Bernoulli- $\tilde{p}$  extension in Table VII.

### F. The GAMPmatlab Software Suite

The GAMP iterations from Algorithm 1, including the nonlinear steps discussed in this section, have been implemented

<sup>5</sup>In Section V we describe how a common  $\pi = \pi_n \forall n$  can be learned.

TABLE V  
DEFINITIONS OF ELASTIC-NET QUANTITIES USED IN TABLE IV

$\sigma \triangleq \sqrt{\tau_r / (2\lambda_2\tau_r + 1)}$	$\tilde{r} \triangleq \hat{r} / (\sigma(2\lambda_2\tau_r + 1))$
$r \triangleq \tilde{r} + \lambda_1\sigma$	$\bar{r} \triangleq \tilde{r} - \lambda_1\sigma$
$C \triangleq \frac{\lambda_1}{2} \exp\left(\frac{r^2 - \tilde{r}^2}{2}\right) \Phi(-r)$	$\bar{C} \triangleq \frac{\lambda_1}{2} \exp\left(\frac{\bar{r}^2 - \tilde{r}^2}{2}\right) \Phi(\bar{r})$
$\mu \triangleq \sigma r - \sigma\phi(-r) / \Phi(-r)$	$\bar{\mu} \triangleq \sigma\bar{r} + \sigma\phi(\bar{r}) / \Phi(\bar{r})$
$v \triangleq \sigma^2 \left[1 - \frac{\phi(r)}{\Phi(r)} \left(\frac{\phi(r)}{\Phi(r)} - r\right)\right]$	$\bar{v} \triangleq \sigma^2 \left[1 - \frac{\phi(\bar{r})}{\Phi(\bar{r})} \left(\frac{\phi(\bar{r})}{\Phi(\bar{r})} + \bar{r}\right)\right]$

TABLE VI  
ACTIVITY-FUNCTIONS AND THEIR GAMP/MATLAB SUM-PRODUCT AND MAX-SUM IMPLEMENTATION METHOD: CF = closed form, VI = variational inference, RF = root-finding

Name	$p_{y z}(y z)$ Description	Sum-Product	Max-Sum
Logistic	$\propto (1 + \exp(-\alpha yz))^{-1}$	VI	RF
Probit	$\Phi\left(\frac{yz}{v}\right)$	CF	RF
Hinge Loss	$\propto \exp(-\max(0, 1 - yz))$	CF	RF
Robust- $p^*$	$\gamma + (1 - 2\gamma)p_{y z}^*(y z)$	CF	RF

TABLE VII  
WEIGHT-COEFFICIENT PRIORS AND THEIR GAMP/MATLAB SUM-PRODUCT AND MAX-SUM IMPLEMENTATION METHOD: CF = closed form, NI = not implemented, NA = not applicable

Name	$p_{w_n}(w)$ Description	Sum-Product	Max-Sum
Gaussian	$\mathcal{N}(w; \mu, \sigma^2)$	CF	CF
GM	$\sum_l \omega_l \mathcal{N}(w; \mu_l, \sigma_l^2)$	CF	NI
Laplacian	$\propto \exp(-\lambda w )$	CF	CF
Elastic Net	$\propto \exp(-\lambda_1 w  - \lambda_2 w^2)$	CF	CF
Bernoulli- $\tilde{p}$	$(1 - \pi_n)\delta(w) + \pi_n\tilde{p}_{w_n}(w)$	CF	NA

in the open-source ‘‘GAMPmatlab’’ software suite.<sup>6</sup> For convenience, the existing activation-function implementations are summarized in Table VI and relevant weight-prior implementations appear in Table VII.

## V. ONLINE PARAMETER TUNING

The activation functions and weight-vector priors described in Section IV depend on modeling parameters that, in practice, must be tuned. For example, the logistic sigmoid (25) depends on  $\alpha$ ; the probit depends on  $v$ ;  $\ell_1$  regularization depends on  $\lambda$ ; and the Bernoulli-Gaussian-mixture prior depends on  $\pi$  and  $\{\omega_l, \mu_l, \sigma_l^2\}_{l=1}^L$ , where  $\omega_l$  parameterizes the weight,  $\mu_l$  the mean, and  $\sigma_l^2$  the variance of the  $l$ th mixture component. Although cross-validation (CV) is the customary approach to tuning parameters such as these, it suffers from two major drawbacks: First, it can be very computationally costly, since each parameter must be tested over a grid of hypothesized values and over multiple data folds. For example,  $K$ -fold cross-validation tuning of  $P$  parameters using  $G$  hypothesized values of each requires the training and evaluation of  $KG^P$  classifiers. Second, leaving out a portion of the training data for CV can degrade classification performance, especially in the example-starved regime where  $M \ll N$  (see, e.g., [42]).

As an alternative to CV, we consider *online learning* of the unknown model parameters  $\theta$  using the methodology from [27],

<sup>6</sup>The latest source code can be obtained through the GAMPmatlab SourceForge Subversion repository at <http://sourceforge.net/projects/gampmatlab/>.

[43]. Here, the goal is to compute the maximum-likelihood estimate  $\hat{\theta}_{\text{ML}} = \arg \max_{\theta} p_{\mathbf{y}}(\mathbf{y}; \theta)$ , where our data model implies a likelihood function of the form

$$p_{\mathbf{y}}(\mathbf{y}; \theta) = \int_{\mathbf{w}} \prod_m p_{y_m|z_m}(y_m | \mathbf{x}^T \mathbf{w}; \theta) \prod_n p_{w_n}(w_n; \theta). \quad (34)$$

Because it is computationally infeasible to evaluate and/or maximize (34) directly, we apply the expectation-maximization (EM) algorithm [44]. For EM, we treat  $\mathbf{w}$  as the ‘‘hidden’’ data, giving the iteration- $j$  EM update

$$\theta^j = \arg \max_{\theta} E_{\mathbf{w}|\mathbf{y}} \{ \log p_{\mathbf{y}, \mathbf{w}}(\mathbf{y}, \mathbf{w}; \theta) \mid \mathbf{y}; \theta^{j-1} \} \quad (35)$$

$$= \arg \max_{\theta} \sum_m E_{z_m|\mathbf{y}} \{ \log p_{y_m|z_m}(y_m | z_m; \theta) \mid \mathbf{y}; \theta^{j-1} \} + \sum_n E_{w_n|\mathbf{y}} \{ \log p_{w_n}(w_n; \theta) \mid \mathbf{y}; \theta^{j-1} \}. \quad (36)$$

Furthermore, to evaluate the conditional expectations in (36), GAMP’s posterior approximations from (9)-(10) are used. It was shown in [45] that, in the large-system limit, the estimates generated by this procedure are asymptotically consistent (as  $j \rightarrow \infty$  and under certain identifiability conditions). Moreover, it was shown in [27], [43] that, for various priors and likelihoods of interest in compressive sensing (e.g., AWGN likelihood, Bernoulli-Gaussian-Mixture priors,  $\ell_1$  regularization), the quantities needed from the expectation in (36) are implicitly computed by GAMP, making this approach computationally attractive. However, because this EM procedure runs GAMP several times, once for each EM iteration (although not necessarily to convergence), the total runtime may be increased relative to that of GAMP without EM.

In this work, we propose EM-based learning of the activation-function parameters, i.e.,  $\alpha$  in the logistic model (25),  $v$  in the probit model (26), and  $\gamma$  in the robust model (30). Starting with  $\alpha$ , we find that a closed-form expression for the value maximizing (36) remains out of reach, due to the form of the logistic model (25). So, we apply the same variational lower bound used for Algorithm 2, and find that the lower-bound maximizing value of  $\alpha$  obeys (see [36])

$$0 = \sum_m \frac{1}{2} (\hat{z}_m y_m - \xi_m) + \frac{\xi_m}{1 + \exp(\alpha \xi_m)}, \quad (37)$$

where  $\xi_m$  is the variational parameter being used to optimize the lower-bound and  $\hat{z}_m \approx E\{z_m | \mathbf{y} = \mathbf{y}\}$  is output by Algorithm 2. We then solve for  $\alpha$  using Newton’s method.

To tune the probit parameter,  $v$ , we zero the derivative of (36) w.r.t  $v$  to obtain

$$0 = \sum_m E_{z_m|\mathbf{y}} \left\{ \frac{\partial}{\partial v} \log p_{y_m|z_m}(y_m | z_m; v^j) \mid \mathbf{y}; v^{j-1} \right\} \quad (38)$$

$$= \sum_m E_{z_m|\mathbf{y}} \left\{ \frac{-\ddot{c}_m(v^j)}{v^j} \phi(\ddot{c}_m(v^j)) \Phi(\ddot{c}_m(v^j))^{-1} \mid \mathbf{y}; v^{j-1} \right\}, \quad (39)$$

where  $\ddot{c}_m(v) \triangleq (y_m z_m) / v$ . We then numerically evaluate the expectation and apply an iterative root-finding procedure to find the EM update  $v^j$  that solves (39).

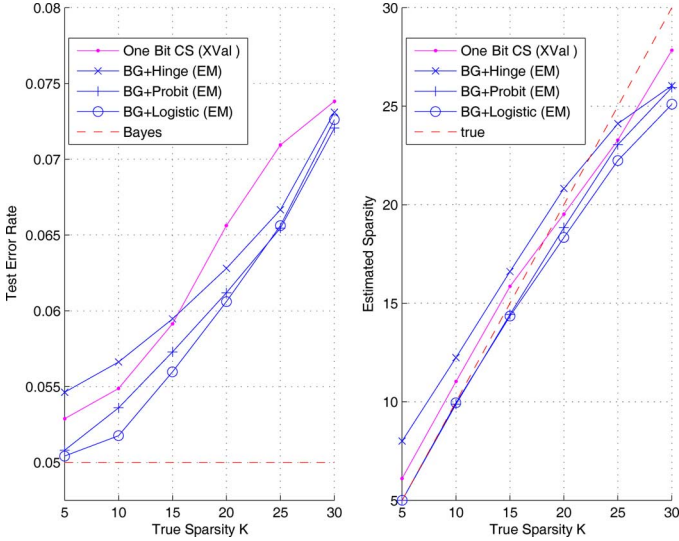


Fig. 3. Test error rate and estimated sparsity  $\hat{K}$  for cross-validation-tuned OneBitCS, and for EM-tuned sum-product GAMP classifiers based on the Bernoulli-Gaussian (BG) prior and the hinge ( $\times$ ), probit ( $+$ ), and logistic ( $\circ$ ) activation functions, as a function of the true sparsity  $K$ . Here,  $N = 30\,000$ ,  $M = 300$ , and Bayes error rate was  $\varepsilon_B = 0.05$ .

To learn  $\gamma$ , we include the corruption indicators  $\beta \in \{0, 1\}^M$  in the EM-algorithm’s hidden data (i.e.,  $\beta_m = 0$  indicates that  $y_m$  was corrupt and  $\beta_m = 1$  that it was not), where an i.i.d assumption on the corruption mechanism implies the prior  $p(\beta; \gamma) = \prod_{m=1}^M \gamma^{1-\beta_m} (1-\gamma)^{\beta_m}$ . In this case, it can be shown [36] that the update of the  $\gamma$  parameter reduces to

$$\gamma^j = \arg \max_{\gamma \in [0,1]} \sum_{m=1}^M \mathbb{E}_{\beta_m | \mathbf{y}} [\log p(\beta_m; \gamma) | \mathbf{y}; \boldsymbol{\theta}^{j-1}] \quad (40)$$

$$= \frac{1}{M} \sum_{m=1}^M p(\beta_m = 0 | \mathbf{y}; \boldsymbol{\theta}^{j-1}), \quad (41)$$

where (41) leveraged  $\mathbb{E}[\beta_m | \mathbf{y}; \boldsymbol{\theta}^{j-1}] = 1 - p(\beta_m = 0 | \mathbf{y}; \boldsymbol{\theta}^{j-1})$ . Moreover,  $p(\beta_m = 0 | \mathbf{y}; \boldsymbol{\theta}^{j-1})$  is easily computed using quantities returned by sum-product GAMP.

## VI. NUMERICAL STUDY

In this section we describe several synthetic and real-world classification problems to which GAMP was applied. Experiments were conducted on a workstation running Red Hat Enterprise Linux (r2.4), with an Intel Core i7–2600 CPU (3.4 GHz, 8 MB cache) and 8 GB DDR3 RAM.

### A. Synthetic Classification in the $N \gg M$ Regime

We first examine a synthetic problem where the number of features,  $N$ , greatly exceeds the number of training examples,  $M$ . As discussed in the Introduction, it is possible to perform accurate classification when  $N \gg M$  if the number of discriminatory features  $K$  is sufficiently small. In this experiment, we consider  $N = 30\,000$ ,  $M = 300$ , and  $K \in \{5, \dots, 30\}$ , where the range on  $K$  is chosen based on the following information-theoretic argument:  $M$  training labels bring  $\log_2 M$  bits of information, whereas at least  $K \log_2(N/K) \leq \log_2 \binom{N}{K}$  bits of information are needed to determine the  $N$ -length  $K$ -sparse Bayes weight vector, assuming that we have no prior knowledge of its

support, which takes on  $\binom{N}{K}$  possibilities. With  $N = 30\,000$  and  $M = 300$ , it turns out that  $K = 31$  is the largest value of  $K \leq N$  such that  $M \geq K \log_2(N/K)$ .

Our experiment was of a Monte-Carlo form. In each trial, we constructed a random  $K$ -sparse Bayes weight vector  $\mathbf{w}$  with a support drawn uniformly at random and with non-zero-coefficient amplitudes drawn uniformly in  $\{-1, 1\}$ . We used  $\pm 1$  amplitudes to eliminate the potential ambiguity about whether a given non-zero coefficient was effectively non-zero, since, e.g., Gaussian-distributed amplitudes can be arbitrarily close to zero. We then constructed a balanced set of training labels  $y_m \in \{-1, 1\}$  (i.e., exactly  $M/2$  labels were positive) and we drew  $M$  i.i.d random feature vectors  $\mathbf{x}_m$  from the class-conditional generative distribution  $\mathbf{x}_m | y_m \sim \mathcal{N}(y_m \mathbf{w}, v \mathbf{I})$ .

Fig. 3 shows both the average test error rate and the average estimated sparsity  $\hat{K}$  for cross-validation tuned “OneBitCS” from [9],<sup>7</sup> and for EM-tuned sum-product GAMP classifiers based on the Bernoulli-Gaussian (BG) prior and activation functions including hinge loss (HL), probit (PR), and logistic (LR). The average was computed over 1000 Monte-Carlo trials, where in each trial the expected error probability of the designed classifier  $\hat{\mathbf{w}}$  was computed in closed form as  $\Phi(-\mathbf{w}^T \hat{\mathbf{w}} / \sqrt{v \|\hat{\mathbf{w}}\|^2})$ . The figure shows all algorithms under test performing relatively close to the Bayes error rate, and for small  $K$  it shows BG-LR and BG-PR GAMP performing extremely close to the Bayes error rate. Comparing the classifiers, we see that GAMP’s BG-LR performs the best, which is not surprising since the logistic activation function is statistically matched to data model in this experiment [37]. Meanwhile, GAMP’s BG-PR classifier performed the second best, and the two remaining classifiers (GAMP’s BG-HL and OneBitCS) performed only slightly worse.

Fig. 3 also shows the sparsities estimated by cross-validation in the case of OneBitCS and by the EM-tuning in the case of GAMP. Since the weights returned by sum-product BG-GAMP are non-zero with probability one, the estimated sparsity is defined as the number of coefficients with posterior support probability  $p(w_n \neq 0 | \mathbf{y})$  exceeding 1/2. The figure shows that all algorithms under test returned accurate estimates of the true sparsity  $K$ . For small values of  $K$ , the estimates returned by BG-LR and BG-PR GAMP were extremely accurate while those for OneBitCS and BG-HL GAMP slightly overestimated the sparsity. Meanwhile, for large values of  $K$ , all algorithms underestimated the sparsity by about 15%.

### B. Text Classification and Adaptive Learning

We next consider a binary text classification problem based on the Reuter’s Corpus Volume I (RCV1) dataset [46]. As in [17], [47], newswire article topic codes CCAT and ECAT were combined to form the positive class while GCAT and MCAT were combined to form constitute the negative class.<sup>8</sup> Although the original dataset consisted of 20 242 balanced training examples of  $N = 47\,236$  features, with 677 399 examples reserved for testing, we followed the approach in [17], [47] and swapped training and testing sets in order to test computational efficiency on a large training dataset (and thus  $M = 677\,399$ ).

<sup>7</sup>For cross-validation of OneBitCS, we used 2 folds and searched over all sparsities in a radius of 10 from the true sparsity  $K$ .

<sup>8</sup>Data was taken from <http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/binary.html>.



TABLE VIII

A COMPARISON OF DIFFERENT CLASSIFIERS ON THE “SWAPPED” RCV1 BINARY DATASET (WHERE  $M \gg N$ ), SHOWING THE TEST-SET CLASSIFICATION ACCURACY, THE TOTAL AND POST-TUNING RUNTIMES, AND THE DENSITY OF THE WEIGHT VECTOR. BELOW, sp = sum-product; ms = max-sum; BG = Bernoulli-Gaussian; PR = Probit; HL = Hinge loss; L1 =  $\ell_1$  REGULARIZATION; LR = Logistic

Classifier	Tuning	Accuracy	Runtime (s)	Density
spGAMP: BG-PR	EM	97.4%	<b>105</b> / 105	8.6%
spGAMP: BG-HL	EM	97.3%	134 / 134	8.9%
msGAMP: L1-LR	EM	97.6%	684 / 123	9.8%
msGAMP: L1-LR	xval	97.6%	3068 / 278	19.6%
CDN [17]	xval	<b>97.7%</b>	1298 / 112	10.9%
TRON [49]	xval	<b>97.7%</b>	1682 / 133	10.8%
TFOCS [48]	xval	97.6%	1086 / 94	19.2%
OneBitCS [9]	xval	90.1%	193 / <b>1</b>	<b>1.3%</b>

As in [17], we constructed feature vectors as cosine-normalized logarithmic transformations of the TF-IDF (term frequency – inverse document frequency) data vectors. We note that the resulting features are very sparse; only 0.16% of the entries in  $\mathbf{X}$  are non-zero. Finally, we trained linear classifiers (i.e., weight vectors) using four GAMP-based methods and four existing state-of-the-art methods: TFOCS [48] in L1-LR mode, CDN [17], TRON [49], and OneBitCS [9]. In doing so, for EM learning we used 5 EM iterations, and for cross-validation we used 2 folds and a logarithmically spaced grid of size  $10^9$ .

Table VIII summarizes the performance achieved by the resulting classifiers, including the test-set classification accuracy, weight-vector density (i.e., the fraction of non-zero weights), and two runtimes: the *total* runtime needed to train the classifier, which includes EM- or cross-validation-based parameter tuning, and the *post-tuning* runtime. Although it is customary to report only the latter, we feel that the former better captures the true computational cost of classifier design. We note that, in the case of spGAMP, the total and post-tuning runtime are identical because EM tuning was performed once per GAMP iteration. In contrast, for msGAMP, we ran many GAMP iterations per EM iteration, and hence the total runtime (which avoids EM iterations) is much longer. We also note that the post-tuning runtime of OneBitCS is extremely fast because of a computational trick that we learned via personal communication with an author, Yaniv Plan: Given signed labels  $y_m \in \{-1, 1\}$  and a sparsity estimate  $\hat{K}$ , the OneBitCS weight vector  $\hat{\mathbf{w}}$  can be computed from the training pair  $(\mathbf{X}, \mathbf{y})$  via  $\hat{\mathbf{w}} = \text{thresh}_{\hat{K}}(\mathbf{X}^T \mathbf{y})$ , where  $\text{thresh}_{\hat{K}}(\cdot)$  is the mapping from  $\mathbb{R}^N \rightarrow \mathbb{R}^N$  that preserves the input components with the largest  $\hat{K}$  magnitudes and zeros the remainder.

Table VIII shows all 8 classifiers achieving nearly identical test-set classification accuracy, with the exception of cross-validated OneBitCS, which gives noticeably poorer accuracy. Interestingly, OneBitCS also gives by far the sparsest weight vectors, apparently at the cost of test-error rate. A better tradeoff between test accuracy and weight vector density is given by the EM-tuned GAMP algorithms, which return weight vectors that are about half as dense as those returned by CDN, TRON, TFOCS, and cross-validated GAMP, but that sacrifice only a fraction-of-a-percent in test accuracy.

<sup>9</sup>For OneBitCS, the cross-validation grid included sparsity rates between 0.1% and 15%.

TABLE IX

A COMPARISON OF DIFFERENT CLASSIFIERS ON THE “NON-SWAPPED” RCV1 BINARY DATASET (WITH  $N > M$ ), SHOWING THE TEST-SET CLASSIFICATION ACCURACY, THE TOTAL AND POST-TUNING RUNTIMES, AND THE DENSITY OF THE WEIGHT VECTOR. BELOW, sp = sum-product; ms = max-sum; BG = Bernoulli-Gaussian; PR = Probit; HL = Hinge loss; L1 =  $\ell_1$  REGULARIZATION; LR = Logistic

Classifier	Tuning	Accuracy	Runtime (s)	Density
spGAMP: BG-PR	EM	95.5%	<b>4</b> / 4	7.4%
spGAMP: BG-HL	EM	95.1%	6 / 6	3.1%
msGAMP: L1-LR	EM	95.6%	16 / 3	1.8%
msGAMP: L1-LR	xval	95.5%	134 / 16	4.6%
CDN [17]	xval	95.5%	11 / 2	5.0%
TRON [49]	xval	<b>96.0%</b>	19 / 3	12.4%
TFOCS [48]	xval	95.7%	17 / 2	4.3%
OneBitCS [9]	xval	89.7%	8 / <b>0.1</b>	<b>0.8%</b>

Table VIII also shows a wide range of runtimes. OneBitCS gives by far the fastest post-tuning runtime, for the reasons described earlier. Among the total runtimes, however, the two fastest are EM-GAMP based, with the best (at 105 seconds) beating the fastest high-accuracy non-GAMP algorithm (i.e., TFOCS at 1086 seconds) by more than a factor of 10. That said, some caution must be used when comparing runtimes. For example, while all algorithms were given a “stopping tolerance” of  $10^{-3}$ , the algorithms apply this tolerance in different ways. Also, CDN and TRON are implemented in C++, while GAMP is implemented in object-oriented MATLAB (and therefore is far from optimized).

To understand how performance is impacted in a data-starved regime (i.e.,  $N > M$ ), we tested each algorithm on the same RCV1 dataset, but *without* swapping the train/test datasets as was done in [17], [47] and our Table VIII. The results are shown in Table IX. Similar to our other RCV experiment, we see all classifiers yielding very similar test error rates, with the exception of OneBitCS, which does significantly worse. Again, however, OneBitCS generates an extremely sparse weight vector at the expense of test error rate, whereas some the EM-tuned BG-HL and L1-LR GAMP algorithms offer (milder) density reduction without a significant cost in test accuracy. Finally, the two fastest total runtimes are earned by the spGAMP algorithms, and the fastest (BG-PR at 4 seconds) is about 3 times as quick as the fastest high-accuracy non-GAMP algorithm (i.e., CDN at 11 seconds).

Finally, we note that, although GAMP was derived under the assumption that the elements of  $\mathbf{X}$  are realizations of an i.i.d sub-Gaussian distribution, it worked well even with the  $\mathbf{X}$  of this experiment, which was far from i.i.d sub-Gaussian. We attribute the robust performance of GAMP to the “damping” mechanism included in the GAMPmatlab implementation, which was first described in [50] and rigorously analyzed in [34]. Essentially, damping slows down the updates with the goal of preventing divergence.

### C. Robust Classification

In Section IV-D, we proposed an approach by which GAMP can be made robust to labels that are corrupted or otherwise highly atypical under a given activation model  $p_{y|z}^*$ . We now evaluate the performance of this robustification method. To

do so, we first generated examples<sup>10</sup>  $(y_m, \mathbf{x}_m)$  with balanced classes such that the Bayes-optimal classification boundary is a hyper-plane with a desired Bayes error rate of  $\varepsilon_B$ . Then, we flipped a fraction  $\gamma$  of the training labels (but not the test labels), trained several different varieties of GAMP classifiers, and measured their classification accuracy on the test data.

The first classifier we considered paired a genie-aided “standard logistic” activation function, (25), with an i.i.d. zero-mean, unit-variance Gaussian weight vector prior. Note that under a class-conditional Gaussian generative distribution with balanced classes, the corresponding activation function is logistic with scale parameter  $\alpha = 2 M\mu$  [37]. Therefore, the genie-aided logistic classifier was provided the true value of  $\mu$ , which was used to specify the logistic scale  $\alpha$ . The second classifier we considered paired a genie-aided robust logistic activation function, which possessed perfect knowledge of both  $\mu$  and the mislabeling probability  $\gamma$ , with the aforementioned Gaussian weight vector prior. To understand how performance is impacted by the parameter tuning scheme of Section V, we also trained EM variants of the preceding classifiers. The EM-enabled standard logistic classifier was provided a fixed logistic scale of  $\alpha = 100$ , and was allowed to tune the variance of the weight vector prior. The EM-enabled robust logistic classifier was similarly configured, and in addition was given an initial mislabeling probability of  $\gamma^0 = 0.01$ , which was updated according to (41).

In Fig. 4, we plot the test error rate for each of the four GAMP classifiers as a function of the mislabeling probability  $\gamma$ . For this experiment,  $\mu$  was set so as to yield a Bayes error rate of  $\varepsilon_B = 0.05$ .  $M = 8192$  training examples of  $N = 512$  training features were generated independently, with the test set error rate evaluated based on 1024 unseen (and uncorrupted) examples. Examining the figure, we can see that EM parameter tuning is beneficial for both the standard and robust logistic classifiers, although the benefit is more pronounced for the standard classifier. Remarkably, both the genie-aided and EM-tuned robust logistic classifiers are able to cope with an extreme amount of mislabeling while still achieving the Bayes error rate, thanks in part to the abundance of training data.

#### D. Multi-Voxel Pattern Analysis

Multi-voxel pattern analysis (MVPA) has become an important tool for analyzing functional MRI (fMRI) data [3]–[5]. Cognitive neuro-scientists, who study how the human brain functions at a physical level, employ MVPA not only to infer a subject’s cognitive state but to gather information about how the brain itself distinguishes between cognitive states. In particular, by identifying *which* brain regions are most important in discriminating between cognitive states, they hope to learn the underlying processes by which the brain operates. In this sense, the goal of MVPA is often feature selection, not classification.

To investigate the performance of GAMP for MVPA, we conducted an experiment using the well-known Haxby dataset [3]. The Haxby dataset consists of fMRI data collected from 6 subjects with 12 “runs” per subject. In each run, the subject pas-

<sup>10</sup>Data was generated according to a class-conditional Gaussian distribution with  $N$  discriminatory features. Specifically, given the label  $y \in \{-1, 1\}$  a feature vector  $\mathbf{x}$  was generated as follows: entries of  $\mathbf{x}$  were drawn i.i.d.  $\mathcal{N}(y\mu, M^{-1})$  for some  $\mu > 0$ . Under this model, with balanced classes, the Bayes error rate can be shown to be  $\varepsilon_B = \Phi(-\sqrt{NM}\mu)$ . The parameter  $\mu$  can then be chosen to achieve a desired  $\varepsilon_B$ .

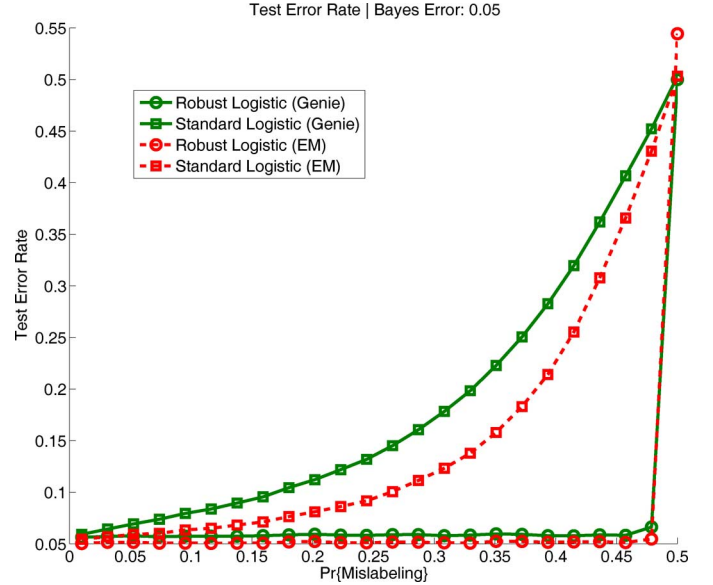


Fig. 4. Test error rate of genie-aided (solid curves) and EM-tuned (dashed curves) instances of standard logistic (□) and robust logistic (○) classifiers, as a function of mislabeling probability  $\gamma$ , with  $M = 8192$ ,  $N = 512$ , and Bayes error rate  $\varepsilon_B = 0.05$ .

sively viewed blocks of 9 greyscale images from each of 8 object categories (i.e., faces, houses, cats, bottles, scissors, shoes, chairs, and nonsense patterns), during which full-brain fMRI data was recorded over  $N = 31\,398$  voxels.

In our experiment, we designed classifiers that predict binary object category (e.g., cat vs. scissors) from  $M$  examples of  $N$ -voxel fMRI data collected from a single subject. For comparison, we tried four algorithms: i)  $\ell_1$ -penalized logistic regression (L1-LR) as implemented using cross-validation-tuned TFOCS [48], ii) L1-LR as implemented using EM-tuned max-sum GAMP, iii) sum-product GAMP under a Bernoulli-Laplace prior and logistic activation function (BL-LR), and iv) a cross-validation-tuned OneBitCS [9] classifier.

Algorithm performance (i.e., error-rate, sparsity, and consistency) was assessed using 12-fold leave-one-out cross-validation. In other words, for each algorithm, 12 separate classifiers were trained, each for a different combination of 1 testing fold (used to evaluate error-rate) and 11 training folds. The reported performance then represents an average over the 12 classifiers. Each fold comprised one of the runs described above, and thus contained 18 examples (i.e., 9 images from each of the 2 object categories constituting the pair), yielding a total of  $M = 11 \times 18 = 198$  training examples. Since  $N = 31\,398$ , the underlying problem is firmly in the  $N \gg M$  regime.

To tune each TFOCS classifier (i.e., select its  $\ell_1$  regularization weight  $\lambda$ ), we used a second level of leave-one-out cross-validation. For this, we first chose a fixed  $G = 10$ -element grid of logarithmically spaced  $\lambda$  hypotheses. Then, for each hypothesis, we designed 11 TFOCS classifiers, each of which used 10 of the 11 available folds for training and the remaining fold for error-rate evaluation. Finally, we chose the  $\lambda$  hypothesis that minimized the error-rate averaged over these 11 TFOCS classifiers. A similar two-level cross-validation strategy was applied for selection of the sparsity rate in OneBitCS, using a logarithmically spaced 50-point grid over sparsity rates between 0.1%

TABLE X  
PERFORMANCE OF CROSS-VALIDATION TUNED L1-LR TFOCS (“TFOCS”), EM-TUNED L1-LR MAX-SUM GAMP (“L1-LR”), EM-TUNED BG-LR SUM-PRODUCT GAMP (“BG-LR”), AND CROSS-VALIDATION TUNED ONE BITCS (“1-Bit”) CLASSIFIERS ON VARIOUS HAXBY PAIRWISE COMPARISONS

Comparison	Error Rate (%)				Sparsity (%)				Consistency (%)				Runtime (s)			
	TFOCS	L1-LR	BG-LR	1-Bit	TFOCS	L1-LR	BG-LR	1-Bit	TFOCS	L1-LR	BG-LR	1-Bit	TFOCS	L1-LR	BG-LR	1-Bit
Cat vs. Scissors	9.7	11.1	9.3	<b>5.1</b>	0.1	0.07	<b>0.01</b>	0.12	38	43	<b>60</b>	57	1318	137	<b>21</b>	202
Cat vs. Shoe	<b>6.1</b>	<b>6.1</b>	11.6	6.5	0.14	0.07	<b>0.01</b>	0.12	34	47	<b>60</b>	59	1347	191	<b>24</b>	205
Cat vs. House	0.4	<b>0.0</b>	1.4	3.7	0.04	0.02	<b>0.01</b>	0.12	53	<b>87</b>	84	75	1364	144	<b>18</b>	202
Bottle vs. Shoe	29.6	30.5	23.6	<b>20.4</b>	0.2	0.1	<b>0.01</b>	0.12	23	31	36	<b>53</b>	1417	166	<b>22</b>	205
Bottle vs. Chair	<b>13.9</b>	<b>13.9</b>	15.7	26.9	0.1	0.07	<b>0.01</b>	0.12	30	45	<b>61</b>	37	1355	150	<b>21</b>	203
Face vs. Chair	<b>0.9</b>	<b>0.9</b>	6.9	2.8	0.09	0.05	<b>0.01</b>	0.12	43	67	68	<b>76</b>	1362	125	<b>24</b>	205
Average	<b>10.1</b>	10.4	11.4	10.9	0.11	0.06	<b>0.01</b>	0.12	37	53	<b>62</b>	60	1358	152	<b>22</b>	204

and 15%. For EM-tuned GAMP, there was no need to perform the second level of cross-validation: we simply applied the EM tuning strategy described in Section V to the 11-fold training data.

Table X reports the results of the above-described experiment for six pairwise comparisons. For all but BG-LR GAMP, sparsity refers to the average percentage of non-zero elements in the learned weight vectors. But, since BG-LR GAMP’s weights are non-zero with probability one, we instead define BG-LR’s sparsity as the number of weights with posterior probability  $p(w_n \neq 0 | \mathbf{y}) > 1/2$ , as we did with the other sum-product-GAMP classifiers in earlier experiments. Consistency refers to the average Jaccard index between weight-vector supports, i.e.,

$$\text{consistency} \triangleq \frac{1}{12} \sum_{i=1}^{12} \frac{1}{11} \sum_{j \neq i} \frac{|\mathcal{S}_i \cap \mathcal{S}_j|}{|\mathcal{S}_i \cup \mathcal{S}_j|} \quad (42)$$

where  $\mathcal{S}_i$  denotes the support of the weight vector learned when holding out the  $i^{\text{th}}$  fold. Runtime refers to the total time used to complete the 12-fold cross-validation procedure.

Ideally, we would like an algorithm that quickly computes weight vectors with low estimated error rate, high consistency, and relatively low density. It should be emphasized that minimizing estimated error rate alone is not of sole importance, especially for this dataset, where the total number of samples is so few that the error rate estimates are understood to be very noisy. Moreover, since the goal of MVPA is to identify which voxels of the brain are most important in discriminating between cognitive states, consistency among folds is very important.

Unfortunately, Table X reveals no clear winner among the algorithms under test. Starting with the estimated error rates, all four algorithms yielded similar comparison-averaged rates, spanning the range from 10.1% (for TFOCS) to 11.4% (for BG-LR GAMP). Interestingly, the algorithm ranking under the consistency metric was exactly the opposite of that for the error-rate metric: BG-LR GAMP yielded the highest consistency (of 62%) and TFOCS the lowest consistency (of 37%). In terms of sparsity, BG-LR GAMP appears to be the winner, but perhaps a direct comparison to the other algorithms should be avoided due to the differences in the definition of sparsity. For runtime, however, the clear winner is EM-tuned BG-LR GAMP, which runs an order-of-magnitude faster than cross-validated OneBitCS and nearly two orders-of-magnitude faster than cross-validated TFOCS.

A direct comparison between cross-validated TFOCS and EM-tuned L1-LR GAMP is illuminating, since these two algorithms share the L1-LR objective and thus differ mainly

in tuning strategy.<sup>11</sup> For this Haxby data, Table X shows that L1-LR GAMP’s classifiers are uniformly more sparse (and nearly twice as sparse on average) as those generated by TFOCS, while suffering only a small degradation in error-rate. Meanwhile, L1-LR GAMP’s classifiers are uniformly more consistent, and its runtimes are about  $9 \times$  faster on average.

## VII. CONCLUSION

In this work, we presented the first comprehensive study of the *generalized approximate message passing* (GAMP) algorithm [22] in the context of linear binary classification and feature selection. We established that a number of popular discriminative models, including logistic and probit regression, as well as support vector machines (via hinge loss), can be implemented in an efficient manner using the GAMP algorithmic framework, and that GAMP’s state evolution formalism can be used in certain instances to predict the misclassification rate of these models. In addition, we demonstrated that a number of sparsity-promoting weight vector priors can be paired with these activation functions to encourage feature selection. Importantly, GAMP’s message passing framework enables us to learn the hyper-parameters that govern our probabilistic models adaptively from the data using expectation-maximization (EM), a trait which can be advantageous in terms of runtime. The flexibility imparted by the GAMP framework allowed us to consider several modifications to the basic discriminative models, such as robust classification, which can be effectively implemented using existing non-robust modules.

In a numerical study, we confirmed the efficacy of our approach on both synthetic and real-world classification problems. For example, we found that the proposed EM parameter tuning can be both computationally efficient and accurate in the applications of text classification and multi-voxel pattern analysis. We also observed on synthetic data that GAMP can attain nearly optimal error rates in the  $N \gg M$  regime when  $N$  is sufficiently large and the number of discriminatory features,  $K$  is sufficiently small. Furthermore, we observed that the robust classification extension can substantially outperform a non-robust counterpart.

## APPENDIX

### SUM-PRODUCT GAMP HINGE-LOSS COMPUTATIONS

In this appendix, we describe the steps needed to compute the sum-product GAMP nonlinear steps for the hinge-loss acti-

<sup>11</sup>It is known that, if max-sum GAMP converges, then it converges to a critical point of the optimization objective [33], which in the (convex) L1-LR case is unique.

vation function, (28). For convenience, we define the associated *un-normalized* likelihood function

$$\tilde{p}_{y|z}(y|z) \triangleq \exp(-\max(0, 1 - yz)), \quad y \in \{-1, 1\}. \quad (43)$$

Note from (9) that the sum-product  $(\hat{z}, \tau_z)$  can be interpreted as the posterior mean and variance of a random variable,  $z$ , with prior  $\mathcal{N}(\hat{p}, \tau_p)$  and likelihood proportional to  $\tilde{p}_{y|z}(y|z)$ .

To compute the statistics  $\hat{z} \equiv E[z|y = y]$  and  $\tau_z \equiv \text{var}\{z|y = y\}$ , we first write the posterior pdf as

$$p_{z|y}(z|y) = C_y^{-1} \tilde{p}_{y|z}(y|z) p_z(z), \quad (44)$$

where  $C_y$  is an appropriate normalization constant. Defining

$$\alpha_y \triangleq ((1 - \tau_p) - y\hat{p})/\sqrt{\tau_p} \quad (45)$$

$$\beta_y \triangleq (y\hat{p} - 1)/\sqrt{\tau_p} \quad (46)$$

$$\delta_y \triangleq y\hat{p} - 1 + \tau_p/2, \quad (47)$$

it can be shown [36] that

$$C_1 = \int_{-\infty}^1 \exp(z - 1) \mathcal{N}(z; \hat{p}, \tau_p) + \int_1^{\infty} \mathcal{N}(z; \hat{p}, \tau_p) \quad (48)$$

$$= \exp(\delta_1) \Phi(\alpha_1) + \Phi(\beta_1) \quad (49)$$

The posterior mean for  $y = 1$  is therefore given by

$$E[z|y = 1] = \frac{1}{C_1} \int_z z \tilde{p}_{y|z}(y = 1|z) p_z(z) \quad (50)$$

$$= \frac{1}{C_1} \left[ e^{\delta_1} \int_{-\infty}^1 z \mathcal{N}(z; \hat{p} + \tau_p, \tau_p) + \int_1^{\infty} z \mathcal{N}(z; \hat{p}, \tau_p) \right] \quad (51)$$

$$= \frac{e^{\delta_1} \Phi(\alpha_1)}{C_1} \int_{-\infty}^1 z \frac{\mathcal{N}(z; \hat{p} + \tau_p, \tau_p)}{\Phi(\alpha_1)} + \frac{\Phi(\beta_1)}{C_1} \int_1^{\infty} z \frac{\mathcal{N}(z; \hat{p}, \tau_p)}{\Phi(\beta_1)}, \quad (52)$$

where each integral in (52) represents the first moment of a truncated normal random variable. Similar expressions can be derived for  $E[z|y = -1]$ . Then, defining the quantities

$$\gamma_y \triangleq e^{-\delta_y} \Phi(\beta_y) / \Phi(\alpha_y) \quad (53)$$

$$\underline{\mu}_y \triangleq \hat{p} + y(\tau_p - \sqrt{\tau_p} \phi(\alpha_y) / \Phi(\alpha_y)) \quad (54)$$

$$\bar{\mu}_y \triangleq \hat{p} + y\sqrt{\tau_p} \phi(\beta_y) / \Phi(\beta_y), \quad (55)$$

it can be shown [51] that, for  $y \in \{-1, 1\}$ ,

$$\hat{z}(y) = E[z|y = y] = (1 + \gamma_y)^{-1} \underline{\mu}_y + (1 + \gamma_y^{-1})^{-1} \bar{\mu}_y. \quad (56)$$

To compute  $\tau_z \equiv \text{var}\{z|y = y\}$ , it suffices to derive an expression for  $E[z^2|y = y]$ . Following the same line of reasoning that produced (52), we find

$$E[z^2|y = 1] = \frac{e^{\delta_1} \Phi(\alpha_1)}{C_1} \int_{-\infty}^1 z^2 \frac{\mathcal{N}(z; \hat{p} + \tau_p, \tau_p)}{\Phi(\alpha_1)} + \frac{\Phi(\beta_1)}{C_1} \int_1^{\infty} z^2 \frac{\mathcal{N}(z; \hat{p}, \tau_p)}{\Phi(\beta_1)}, \quad (57)$$

where each integral in (57) is the second moment of a truncated normal random variable. A similar expression can be derived for  $E[z^2|y = -1]$ . Defining

$$\underline{v}_y \triangleq \tau_p \left[ 1 - \frac{\phi(\alpha_y)}{\Phi(\alpha_y)} \left( \frac{\phi(\alpha_y)}{\Phi(\alpha_y)} + \alpha_y \right) \right] \quad (58)$$

$$\bar{v}_y \triangleq \tau_p \left[ 1 - \frac{\phi(\beta_y)}{\Phi(\beta_y)} \left( \frac{\phi(\beta_y)}{\Phi(\beta_y)} + \beta_y \right) \right], \quad (59)$$

it can be shown [51] that

$$E[z^2|y = y] = (1 + \gamma_y)^{-1} (\underline{v}_y + \underline{\mu}_y^2) + (1 + \gamma_y^{-1})^{-1} (\bar{v}_y + \bar{\mu}_y^2), \quad (60)$$

allowing us to compute  $\tau_z(y) = E[z^2|y = y] - \hat{z}^2(y)$ .

## REFERENCES

- [1] C. M. Bishop, *Pattern Recognition and Machine Learning*. New York, NY, USA: Springer-Verlag, 2006.
- [2] G. Forman, "An extensive empirical study of feature selection metrics for text classification," *J. Mach. Learn. Res.*, vol. 3, pp. 1289–1305, 2003.
- [3] J. V. Haxby, M. I. Gobbini, M. L. Furey, A. Ishai, J. L. Schouten, and P. Pietrini, "Distributed and overlapping representations of faces and objects in ventral temporal cortex," *Sci.*, vol. 293, pp. 2425–2430, Sep. 2001.
- [4] K. A. Norman, S. M. Polyn, G. J. Detre, and J. V. Haxby, "Beyond mind-reading: Multi-voxel pattern analysis of fMRI data," *Trends Cognit. Sci.*, vol. 10, pp. 424–430, Sep. 2006.
- [5] F. Pereira, T. Mitchell, and M. Botvinick, "Machine learning classifiers and fMRI: A tutorial overview," *NeuroImage*, vol. 45, pp. S199–S209, Mar. 2009.
- [6] A. Gustafsson, A. Hermann, and F. Huber, *Conjoint Measurement: Methods and Applications*. Berlin, Germany: Springer-Verlag, 2007.
- [7] E. P. Xing, M. I. Jordan, and R. M. Karp, "Feature selection for high-dimensional genomic microarray data," in *Proc. Int. Workshop Mach. Learn.*, 2001, pp. 601–608.
- [8] P. T. Boufounos and R. G. Baraniuk, "1-bit compressive sensing," presented at the Conf. Inf. Sci. Sys., Princeton, NJ, USA, Mar. 2008.
- [9] Y. Plan and R. Vershynin, "Robust 1-bit compressed sensing and sparse logistic regression: A convex programming approach," *IEEE Trans. Inf. Theory*, vol. 59, no. 1, pp. 482–494, 2013.
- [10] D. Koller and M. Sahami, L. Saitta, Ed., "Toward optimal feature selection," in *Proc. 13th Int. Conf. Mach. Learn. (ICML)*, Bari, Italy, 1996, pp. 284–292.
- [11] R. Kohavi and G. John, "Wrapper for feature subset selection," *Artif. Intell.*, vol. 97, pp. 273–324, 1997.
- [12] M. E. Tipping, "Sparse Bayesian learning and the relevance vector machine," *J. Mach. Learn. Res.*, vol. 1, pp. 211–244, 2001.
- [13] M. Figueiredo, "Adaptive sparseness using Jeffreys' prior," in *Proc. 14th Conf. Adv. Neural Inf. Process. Syst.*, Cambridge, MA, USA, 2001, pp. 697–704.
- [14] M. Figueiredo, "Adaptive sparseness for supervised learning," *IEEE Trans. Pattern Anal. Mach. Intell. (PAMI)*, vol. 25, no. 9, pp. 1150–1159, 2003.
- [15] A. Kabán, "On Bayesian classification with Laplace priors," *Pattern Recognit. Lett.*, vol. 28, no. 10, pp. 1271–1282, 2007.
- [16] H. Chen, P. Tino, and X. Yao, "Probabilistic classification vector machines," *IEEE Trans. Neural Net.*, vol. 20, no. 6, pp. 901–914, 2009.
- [17] G.-X. Yuan, K.-W. Chang, C.-J. Hsieh, and C.-J. Lin, "A comparison of optimization methods and software for large-scale L1-regularized linear classification," *J. Mach. Learn. Res.*, vol. 11, pp. 3183–3234, 2010.
- [18] A. Gupta, R. Nowak, and B. Recht, "Sample complexity for 1-bit compressed sensing and sparse classification," presented at the Int. Symp. Inf. Theory (ISIT), Austin, TX, 2010.
- [19] J. N. Laska, Z. Wen, W. Yin, and R. G. Baraniuk, "Trust, but verify: Fast and accurate signal recovery from 1-bit compressive measurements," *IEEE Trans. Signal Process.*, vol. 59, no. 11, pp. 5289–5301, 2011.
- [20] U. S. Kamilov, A. Bourquard, A. Amini, and M. Unser, "One-bit measurements with adaptive thresholds," *IEEE Signal Process. Lett.*, vol. 19, pp. 607–610, 2012.

- [21] U. S. Kamilov, V. K. Goyal, and S. Rangan, "Message-passing de-quantization with applications to compressed sensing," *IEEE Trans. Signal Process.*, vol. 60, pp. 6270–6281, Dec. 2012.
- [22] S. Rangan, "Generalized approximate message passing for estimation with random linear mixing," in *Proc. IEEE Int. Symp. Inf. Theory*, St. Petersburg, Russia, Aug. 2011, pp. 2168–2172.
- [23] D. L. Donoho, A. Maleki, and A. Montanari, "Message passing algorithms for compressed sensing," in *Proc. Nat. Acad. Sci.*, Nov. 2009, vol. 106, pp. 18914–18919.
- [24] D. L. Donoho, A. Maleki, and A. Montanari, "Message passing algorithms for compressed sensing: I. Motivation and construction," in *Proc. Inf. Theory Workshop*, Jan. 2010, pp. 1–5.
- [25] A. Javanmard and A. Montanari, "State evolution for general approximate message passing algorithms, with applications to spatial coupling," *Inf. Inference*, vol. 2, no. 2, pp. 115–144, 2013.
- [26] J. Ziniel, S. Rangan, and P. Schniter, "A generalized framework for learning and recovery of structured sparse signals," presented at the IEEE Statist. Signal Process. Workshop, Ann Arbor, MI, USA, Aug. 2012.
- [27] J. P. Vila and P. Schniter, "Expectation-Maximization Gaussian-mixture approximate message passing," *IEEE Trans. Signal Process.*, vol. 61, pp. 4658–4672, Oct. 2013.
- [28] P. Schniter, "Turbo reconstruction of structured sparse signals," in *Proc. Conf. Inf. Sci. Syst. (CISS)*, Princeton, NJ, USA, Mar. 2010, pp. 1–6.
- [29] B. J. Frey and D. J. C. MacKay, "A revolution: Belief propagation in graphs with cycles," *Adv. Neural Inf. Process. Syst.*, pp. 479–485, 1998.
- [30] F. R. Kschischang, B. J. Frey, and H. A. Loeliger, "Factor graphs and the sum-product algorithm," *IEEE Trans. Inf. Theory*, vol. 47, pp. 498–519, Feb. 2001.
- [31] R. J. McEliece, D. J. C. MacKay, and J. Cheng, "Turbo decoding as an instance of Pearl's belief propagation algorithm," *IEEE J. Sel. Areas Commun.*, vol. 16, pp. 140–152, Feb. 1998.
- [32] W. T. Freeman, E. C. Pasztor, and O. T. Carmichael, "Learning low-level vision," *Int. J. Comput. Vis.*, vol. 40, pp. 25–47, Oct. 2000.
- [33] S. Rangan, P. Schniter, E. Riegler, A. K. Fletcher, and V. Cevher, "Fixed points of generalized approximate message passing with arbitrary matrices," in *Proc. IEEE Int. Symp. Inf. Theory*, Istanbul, Jul. 2013, pp. 664–668.
- [34] S. Rangan, P. Schniter, and A. Fletcher, "On the convergence of generalized approximate message passing with arbitrary matrices," in *Proc. IEEE Int. Symp. Inf. Theory*, Honolulu, HI, USA, Jul. 2014, pp. 236–240.
- [35] H.-A. Loeliger, "An introduction to factor graphs," *IEEE Signal Process. Mag.*, vol. 21, pp. 28–41, Jan. 2004.
- [36] J. Ziniel, "Message passing approaches to compressive inference under structured signal priors," Ph.D. dissertation, The Ohio State Univ., Columbus, OH, USA, 2014.
- [37] M. I. Jordan, "Why the logistic function? A tutorial discussion on probabilities and neural networks," MIT Comput. Cognit. Sci., Cambridge, MA, USA, Tech. Rep. 9503, Aug. 13, 1995.
- [38] C. E. Rasmussen and C. K. I. Williams, *Gaussian Processes for Machine Learning*. Cambridge, MA, USA: MIT Press, 2006.
- [39] J. Bi, K. Bennett, M. Embrechts, C. Breneman, and M. Song, "Dimensionality reduction via sparse support vector machines," *J. Mach. Learn. Res.*, vol. 3, pp. 1229–1243, 2003.
- [40] M. Opper and O. Winther, *Gaussian Processes and SVM: Mean Field Results and Leave-One-Out Estimator*. Cambridge, MA, USA: MIT Press, 2000, ch. 17, pp. 311–326.
- [41] H. Zou and T. Hastie, "Regularization and variable selection via the elastic net," *J. Roy. Statist. Soc., B*, vol. 67, no. 2, pp. 301–320, 2005.
- [42] A. K. Nigam, K. McCallum, S. Thrun, and T. Mitchell, "Text classification from labeled and unlabeled documents using EM," *Mach. Learn.*, vol. 39, pp. 103–134, 2000.
- [43] J. P. Vila and P. Schniter, "An empirical-Bayes approach to recovering linearly constrained non-negative sparse signals," *IEEE Trans. Signal Process.*, vol. 62, pp. 4689–4703, Sep. 2014.
- [44] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *J. Roy. Statist. Soc., B*, vol. 39, pp. 1–38, 1977.
- [45] U. S. Kamilov, S. Rangan, A. K. Fletcher, and M. Unser, "Approximate message passing with consistent parameter estimation and applications to sparse learning," in *Proc. Neural Inf. Process. Syst. Conf.*, Lake Tahoe, NV, USA, Dec. 2012, pp. 2447–2455.
- [46] D. D. Lewis, Y. Yang, T. G. Rose, and F. Li, "RCV1: A new benchmark collection for text categorization research," *J. Mach. Learn. Res.*, vol. 5, pp. 361–397, 2004.
- [47] C. Lin, R. C. Weng, and S. S. Keerthi, "Trust region Newton methods for large-scale logistic regression," in *Proc. 24th Int. Conf. Mach. Learn.*, Corvallis, OR, USA, 2007, pp. 561–568.
- [48] S. R. Becker, E. J. Candès, and M. C. Grant, "Templates for convex cone problems with applications to sparse signal recovery," *Math. Prog. Comput.*, vol. 3, no. 3, pp. 165–218, 2011.
- [49] C. J. Lin and J. J. Moré, "Newton's method for large-scale bound constrained problems," *SIAM J. Optim.*, vol. 9, pp. 1100–1127, 1999.
- [50] P. Schniter and S. Rangan, "Compressive phase retrieval via generalized approximate message passing," presented at the Allerton Conf. Commun., Control, Comput., Monticello, IL, USA, Oct. 2012.
- [51] D. R. Barr and E. T. Sherrill, "Mean and variance of truncated normal distributions," *Amer. Statist.*, vol. 53, Nov. 1999.



**Justin Ziniel** (S'11) received the B.S., M.S., and Ph.D. degrees in Electrical and Computer Engineering from The Ohio State University in Columbus, OH, in 2007, 2012, and 2014, respectively.

His research interests include statistical signal processing, machine learning, and interdisciplinary "big data" inference problems.



**Philip Schniter** (F'14) received the B.S. and M.S. degrees in Electrical Engineering from the University of Illinois at Urbana-Champaign in 1992 and 1993, respectively, and the Ph.D. degree in Electrical Engineering from Cornell University in Ithaca, NY, in 2000.

From 1993 to 1996 he was employed by Tektronix Inc. in Beaverton, OR as a systems engineer. After receiving the Ph.D. degree, he joined the Department of Electrical and Computer Engineering at The Ohio State University, Columbus, where he is currently a

Professor and a member of the Information Processing Systems (IPS) Lab. In 2008–2009 he was a visiting professor at Eurecom, Sophia Antipolis, France, and Supélec, Gif-sur-Yvette, France.

In 2003, Dr. Schniter received the National Science Foundation CAREER Award. His areas of interest currently include statistical signal processing, wireless communications and networks, and machine learning.

**Per Sederberg**, photograph and biography not available at the time of publication.