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

9 z w y ] H m ] m ] V̌ e C f ] e ] D Q [ m ] Y V N m R V B̌ V R V t P V t Z t

*Dept. of ECE, The Ohio State University, Columbus, OH 43210. (Email: zinieli@ece.osu.edu, schniter@ece.osu.edu)
†Dept. of Psychology, The Ohio State University, Columbus, OH 43210. (Email: sederberg.1@osu.edu)

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 [3]–[5], conjoint analysis [6], and microarray gene expression [7]. In the $N \gg M$ regime, the model (1) also coincides with that of *noisy one-bit compressed sensing* (CS) [8], [9], where $w$ is the sparse vector of interest and $a_m^T$ is the $m$th row of the $M \times N$ measurement matrix. ... study to confirm the efficacy, flexibility, and speed afforded by our GAMP-based approaches to binary classification.

## II. GAMP

We now formulate the binary linear classifier design problem from both minimum error-rate and maximum a posteriori (MAP) perspectives, and show that each case can be tackled by a variant of the GAMP algorithm.

### A. Sum-Product GAMP

Suppose that we are given $M$ labeled training examples $\{y_m, a_m\}$ and $T$ test feature vectors $\{a_t\}$ associated with unknown test labels $\{y_t\}$, all obeying the noisy linear model (1) under some known error pdf $p_e$ and thus

$$ \text{(1)} $$

## I. INTRODUCTION

In this work we consider binary linear classification
is the signum function and $n$ is a random perturbation accounting for model inaccuracies. For the purpose of learning $w$, we assume the availability of $M$ labeled training examples generated independently according to this model:

$$ y_m = \mathrm{sgn}(a_m^T \cdot w) \oplus n_m, \quad m = 1, \ldots, M, \tag{1} $$

with $n_m$ i.i.d. It is common to express the relationship between the label $y_m$ and the score $a_m^T \cdot w$ via the conditional pdf $p_{y_m|z_m}$ — known as the "activation function," which can be related to the perturbation pdf $p_e$ via

$$ \text{(2)} $$

We are particularly interested in classification problems in which the number of potentially discriminatory features $N$

$$ p_{y_m|z_m}(1|z_m) = \int_{-\infty}^{z_m} p_e(e)\, de = 1 - p_{y_m|z_m}(-1|z_m). $$

dimensional case of interest, we turn to a recently developed approximation: the sum-product variant of GAMP [10], as specified in Algorithm There, the expectation and variance in lines 5-6 and 16-17 are taken elementwise w.r.t. the GAMP-approximated marginal posterior pdfs (with superscript $k$ denoting the iteration)

$$\tag{7}$$

$$\tag{8}$$

with appropriate normalizations 1 and 1 1 and the vector-vector multiplications and divisions in lines 3, 9, 11, 12, 14, 13, 20 are performed elementwise. Due to space limitations, we refer the interested reader to [10] for an overview and derivation of GAMP, to [13] for rigorous analysis under large i.i.d. sub-Gaussian ff1 and to [14] for fixed-point analysis under arbitrary ff 4

---

**Generalized Approximate Message Passing**

Matrix ff1 priors $V_{x}^{3}$. activation functions $V_{fr_w}^t$ w 1 and mode $\in \{\mathsf{SumProduct}, \mathsf{MaxSum}\}$

*(Algorithm listing — content illegible)*

Terminated

---

Applying GAMP to the classification factor graph in Fig. 1a and examining the resulting form of lines 5-6 in Algorithm 1, it becomes evident that the test-label nodes do not affect the GAMP weight estimates and thus the factor graph can effectively be simplified to the form shown in Fig. 1b, after which the (approximated) posterior test-label pdfs are computed via

$$\tag{9}$$

where $\hat{z}_t^\infty$ and ] denote the P element of the GAMP vectors

---



*(figure caption — content illegible)*

known We then consider the problem of computing the classification-error-rate minimizing hypotheses

$$\tag{3}$$

with h and ff5 Note that we treat the labels as random but the features as deterministic parameters. The probabilities in (3) can be computed via the marginalization

$$\tag{4}$$

$$\tag{5}$$

with scaling constant 1 label vector h and constraint set which fixes the $t$th element of h at the value fi and the first $M$ elements of h at the values of the corresponding training labels. The joint pdf in (5) factors as

$$\tag{6}$$

due to the linear model and assuming a separable weight-vector prior, i.e.,

The factorization (6) is illustrated using the *factor graph* in Fig. 1a, 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 (5) is computationally intractable due to the high-dimensional summation and integration, the factor graph in Fig. 1a suggests the use of loopy belief propagation (LBP), and in particular the *sum-product algorithm* (SPA) [11], 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 loopy graphs [12].

Because a direct application of the SPA to the factor graph in Fig. 1a is itself computationally infeasible in the high-

$z^k$ and $a$ respectively, at the final iteration "$k = \check{s}$".

### *B. Max-Sum GAMP*

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

$$\hat{w} = \arg\min_{w \in \mathbb{R}^N} \sum_{m=1}^{M} f_{z_m}(x_m^T w) + \sum_{n=1}^{L\check{s}} f_{w_n}(w_n) \quad (10)$$

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

The solution to (10) can be recognized as the *maximum a posteriori* (MAP) estimate of random vector $w$ under a separable prior $P_w()$ and likelihood corresponding to (1), i.e.,

$$P_{y|w}(y|w; X) = \prod_{m=1}^{M} P_{y_m|z_m}(y_m | x_m^T w), \quad (11)$$

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

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 [15]. Since max-sum LBP is itself intractable for the high-dimensional problems of interest, we turn to the max-sum variant of GAMP [10], which is also specified in Algorithm 1. There, lines 8-9 are to be interpreted as

$$\hat{z}_m^k = \operatorname{prox}_{\tau_{p_m}^k f_{z_m}}(\hat{p}_m^k), \quad m = 1, \ldots, M, \quad (12)$$

$$\tau_{z_m}^k = \tau_{p_m}^k \operatorname{prox}'_{\tau_{p_m}^k f_{z_m}}(\hat{p}_m^k), \quad m = 1, \ldots, M, \quad (13)$$

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

$$\operatorname{prox}_{\tau f}(v) = \arg\min_{u \in \mathbb{R}} \left[ f(u) + \frac{1}{2\tau}(u - v)^2 \right] \quad (14)$$

$$\operatorname{prox}'_{\tau f}(v) = \left(1 + \tau f''(\operatorname{prox}_{\tau f}(v))\right)^{-1} \quad (15)$$

and lines 19–20 are to be interpreted similarly. It is known [14] that, for *arbitrary* $X$, the fixed points of GAMP correspond to the critical points of the optimization objective (10).

### III. ERROR-RATE AND MSE PREDICTION VIA STATE EVOLUTION

One of the hallmarks of GAMP is that its behavior in the large-system limit (i.e., $M, N \to \check{s}$ with fixed ratio $\beta = M/N$) under i.i.d. sub-Gaussian $X$ is characterized by a scalar state evolution [10], [13]. 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 $(y, X, w)$ of the random triple $(y, X, w)$. Recall that, for a given problem realization $(y, X, w)$, the GAMP iterations in Algorithm 1 yield the sequence of estimates $\{\hat{w}^k\}$ of the true weight vector $w$. Then, according to the state evolution, $P_{w,\hat{w}^k}(w, \hat{w}^k)$.

$\prod_n P_{w_n,\hat{w}_n}(w_n, \hat{w}_n)$ and the first two moments of the joint pdf $P_{w,\hat{w}}$ can be computed using [10, Algorithm 3].

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

$$\mathcal{E}^k = E\{\mathbb{1}\{y \neq \operatorname{sgn}(\hat{z}^k)\}\} \quad (16)$$

$$= \sum_{y \in \{-1,1\}} \iint \mathbb{1}\{y \neq \operatorname{sgn}(\hat{z}^k)\} P_{y,z,\hat{z}^k}(y, z, \hat{z}^k) \, dz \, d\hat{z}^k \quad (17)$$

$$= \sum_{y \in \{-1,1\}} \iint \mathbb{1}\{y \neq \operatorname{sgn}(\hat{z}^k)\} P_{y|z}(y|z) p_{z,\hat{z}^k}(z, \hat{z}^k) \, dz \, d\hat{z}^k. \quad (18)$$

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

$$\begin{bmatrix} z \\ \hat{z}^k \end{bmatrix} \xrightarrow{d} \mathcal{N}\left( 0 = \mathcal{N}\left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma^k \right), \begin{bmatrix} \Sigma_{11}^k & \Sigma_{12}^k \\ \Sigma_{21}^k & \Sigma_{22}^k \end{bmatrix} \right), \quad (19)$$

where $\xrightarrow{d}$ indicates convergence in distribution. In [16], it is shown that the above matrix components are

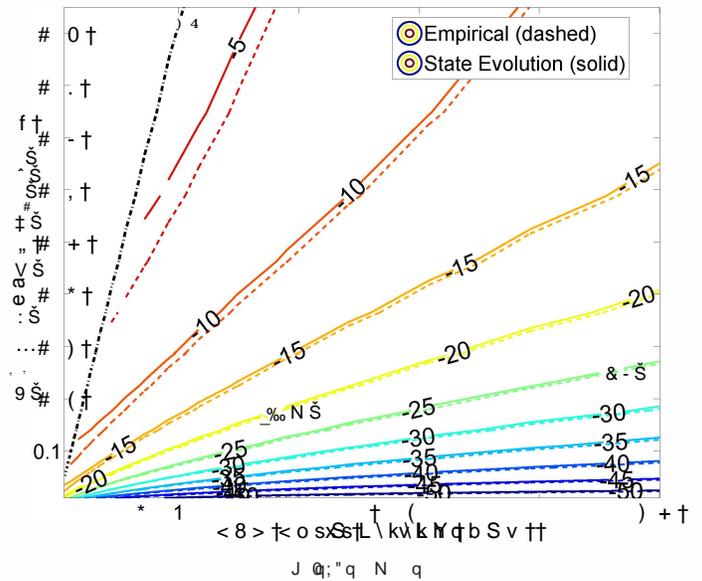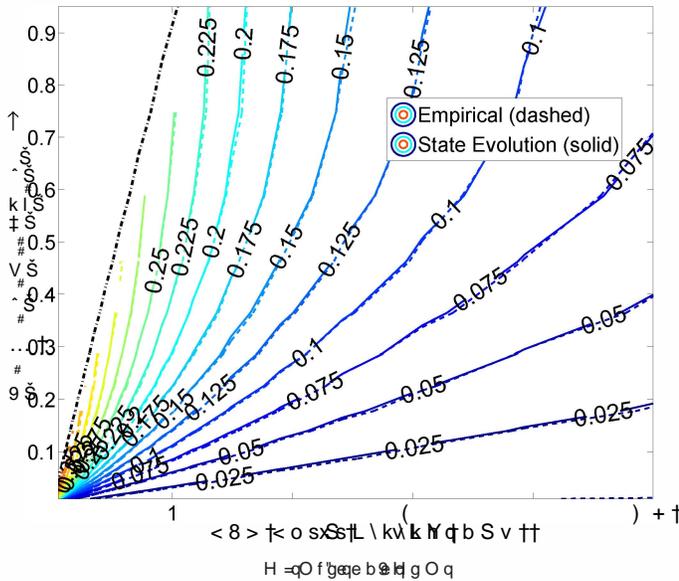$$\Sigma_{11}^k = \delta^{-1}\left(\operatorname{var}\{w_n\} + E[w_n]^2\right), \quad (20)$$

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

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

for label-to-feature ratio $\beta$. As described earlier, the above moments can be computed using [10, Algorithm 3]. The integral in (18) 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 $k^{th}$ 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 100 problem realizations, a true weight vector $w \in \mathbb{R}^{L\check{s}}$ was drawn i.i.d. zero-mean Bernoulli-Gaussian and a feature matrix $X$ was drawn i.i.d. Gaussian, yielding true scores $z = Xw$, from which the true labels $y$ were randomly drawn using a probit activation function $P_{y|z}$. A GAMP weight-vector estimate $\hat{w}^\infty$ was then computed using the training data $(y, X)$ from which the test-label estimates $\{\hat{y}_t\}$ with $\hat{y}_t = \operatorname{sgn}(x_t^T \hat{w}^\infty)$ were computed and compared to the true test-labels in order to calculate the test-error rate for that realization. Figure 2a 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 $\beta$ and discriminative-feature ratio $\rho$, where $K = \|w\|_0$ and $N = 1024$. Similarly, Fig. 2b plots average empirical MSE versus state-evolution predicted MSE, where $MSE = \frac{1}{N} E\{\|\hat{w}^\infty - w\|_2^2\}$.

---

In both Fig. 2a and Fig. 2b, 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. 2a and Fig. 2b show very close agreements between empirical averages and state-evolution predictions.

## IV. GAMP for Binary Classification

We now recall some of the weight-vector priors $p_x(\cdot)$ and activation functions $p_{y|z}(y_m | \cdot)$ (equivalently, regularizations $f_x(\cdot)$ and losses $f_z(\cdot)$) commonly used for binary linear classification, and we describe how GAMP handles them via the nonlinear steps 5-6, 8-9, 16-17, and 19-20 in Algorithm For sum-product GAMP, we recall that the mean and variance computations in lines 5-6 and 16-17 are computed based on the pdfs in (7) and (8), respectively, and for max-sum GAMP the prox steps in 8-9 are computed using equations (12)-(13) and those in 19-20 are computed similarly.

### A. Activation Functions

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

$$p_{y|z}(y|z;e) = \frac{1}{1 + \exp(-yez)}, \quad y \in \{-1, 1\} \qquad (23)$$

where $e > 0$ controls the steepness of the transition. For logistic sum-product GAMP, the mean and variance $(\hat{z}, \tau_z)$ of the marginal posterior approximation can be computed using a variational approach described in [16]. For logistic max-sum GAMP, $\hat{z}$ from (12) solves the scalar minimization problem (14) with $f(u) = -\log p_{y|z}(y|u;e)$ from (23), which is convex. This $\hat{z}$ can be found via bisection search, by locating

the root of $\frac{d}{du}[f(u) + \frac{1}{\tau}(\cdot)]$. The $\tau_z$ from (13) can then be computed in closed-form using $\hat{z}$ and $f''(\cdot)$ via

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

$$p_{y|z}(1|z;v) = \int_{-\infty}^{z} \mathcal{N}(\tau;0,v)d\tau = \Phi\left(\frac{z}{\sqrt{v}}\right) \qquad (24)$$

where $p_{y|z}(-1|z) = 1 - p_{y|z}(1|z)$ and where $v > 0$ controls the steepness of the sigmoid. In the sum-product case, the density (7) corresponds to the posterior pdf of a random variable $z$ with prior $\mathcal{N}(\hat{r}, \tau_p)$ from an observation $y$ measured under the model (24), and a derivation in [17, §3.9] provides expressions for the moments used by GAMP in lines 5-6 of Algorithm In the max-sum case, $(\hat{z}, \tau_z)$ can be computed using bisection search, akin to that described above.

The *hinge loss* $f_z(z) = \max(0, 1 - ymz)$, commonly used in the *support vector machine* (SVM) approach to maximum-margin classification [1, §7.1], yields the activation function

$$p_{y|z}(y|z) \propto \exp(-\max(0, 1 - yz)), \quad y \in \{-1, 1\}. \quad (25)$$

For sum-product GAMP, the corresponding $(\hat{z}, \tau_z)$ from (7) can be computed in closed-form using a procedure described in [16], and for max-sum GAMP, the proximal steps (12)-(13) can be efficiently computed using bisection search.

### B. Weight-Vector Priors

For linear classification and feature selection with $N > M$, it is customary to choose a prior $p_x(\cdot)$ that leads to sparse (or approximately sparse) weight vectors $x$, as discussed below.

In the sum-product case, this is typical accomplished using a Bernoulli-$\tilde{p}$ prior, i.e.,

$$p_x(x) = (1 - \tilde{p})\delta(x) + \tilde{p} f_x(x) \qquad (26)$$

where $\delta(\cdot)$ is the Dirac delta function, $\pi_n \in [0,1]$ is the prior probability that $w_n \neq 0$, and $\tilde{p}_{w_n}(\cdot)$ is the prior pdf of $w_n$ when non-zero. With GAMP it is common to use Bernoulli-Gaussian [18], Bernoulli-Gaussian-mixture [19], and Bernoulli-Laplacian [16], all of which yield (8) with closed-form first and second moments, $g_{w_n}(\cdot, \tau_w)$.

In the max-sum case, the GAMP nonlinear outputs $g_{w_n}(\cdot, \tau_w)$ are computed just as $(\hat{z}, \tau_z)$ are in (12)-(15), but with $f''_{w_n}(\cdot)$ used in place of $p_{y|z}(y_m|\cdot)$. Common examples include $f_{w_n}(w) = \lambda_1 |w|$ for $\ell_1$ regularization [20], $f_{w_n}(w) = \lambda_2 w^2$ for $\ell_2$ regularization [10], and $f_{w_n}(w) = \lambda_1 |w| + \lambda_2 w^2$ for the "elastic net" [21], all of which yield closed-form implementations of $g_{w_n}(\hat{w}, \tau)$. (See [16] for details.)

## 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 (23) depends on $\alpha$; the probit depends on $v$; $\ell_1$ regularization depends on $\lambda$; and the Bernoulli-Gaussian-mixture prior depends on $\pi$ and $\{\omega_l, \theta_l, \phi_l\}_{l=1}^{L}$ where $\omega_l$ parameterizes the weight, $\theta_l$ the mean, and $\phi_l$ 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 design 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., [22]).

As an alternative to CV, we consider *online learning* of the unknown model parameters $\theta$ using the method proposed in [19], [23]. Here, the goal is to compute the maximum-likelihood estimate $\hat{\theta}_{ML} = \arg\max_\theta p(y; \theta)$, where our data model implies a likelihood function of the form

$$p(y; \theta) = \int \prod_m p_{y|z}(y_m | z_m^\top w; \theta) \prod_n p_{w_n}(w_n; \theta) \, (27)$$

Because it is computationally infeasible to evaluate and/or maximize (27) directly, approximate-ML estimation is performed using expectation-maximization (EM) [24], treating $w$ as the "hidden" data, which yields the iteration-$j$ estimate

$$\hat{\theta}^{j+1} = \arg\max_\theta \mathbb{E}\{\ln p(y, w; \theta) | y; \hat{\theta}^j\} \qquad (28)$$

$$= \arg\max_\theta \sum_m \mathbb{E}\{\ln p_{y|z}(y_m | z_m^\top w; \theta) | y; \hat{\theta}^j\}$$

$$+ \sum_n \mathbb{E}\{\ln p_{w_n}(w_n; \theta) | y; \hat{\theta}^j\}. \qquad (29)$$

Furthermore, to evaluate the conditional expectations in (29), GAMP's posterior approximations from (7)-(8) are used. It was shown in [25] that, in the large-system limit, the estimates generated by this procedure are asymptotically consistent (as $j \to \infty$ and under certain identifiability conditions). Moreover, it was shown in [19], [23] 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 (29) are implicitly computed by GAMP, making this approach computationally attractive. Still, this EM procedure runs GAMP several times, once for each EM iteration (although not necessarily to convergence), thus potentially increasing the runtime over GAMP without EM.

In this work, we export the EM-GAMP methods from [19], [23] to the classification setting, which requires deriving the EM updates for the activation-function parameters, such as $\alpha$ in the logistic model (23) and $v$ in the probit model (24). Due to lack of space, we leave the algorithmic details to [16], but we demonstrate the performance of this approach in Section VI.

## VI. NUMERICAL STUDY

In this section we apply GAMP[4] to a text classification problem. This experiment was conducted on a workstation running Red Hat Enterprise Linux (r2.4), with an Intel Core i7-2600 CPU (3.4 GHz, 8MB cache) and 8GB DDR3 RAM. Other experiments can be found in [16], but they are not shown here for reasons of space.

We first consider a binary text classification problem based on the Reuter's Corpus Volume I (RCV1) dataset [26]. As in [27], [28], 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.[5] 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 [27], [28] and swapped training and testing sets in order to test computational efficiency on a large training dataset (and thus $M = 677\,399$). As in [27], 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 $X$ are non-zero. Finally, we trained linear classifiers (i.e., weight vectors) using four GAMP-based methods and three existing state-of-the-art methods: TRON [29], CDN [27], and TFOCS [30] in L1-LR mode. For EM learning we used $5$ EM iterations, and for cross-validation we used 2 folds and a logarithmically spaced grid size of $10$.

Table I 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),[6] 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.

Table I shows all 7 classifiers achieving nearly identical test-set classification accuracy. However, if feature selection is

[4] We used the GAMP implementation from the open-source "GAMPmatlab" suite at http://sourceforge.net/projects/gampmatlab/, and the activation-function and parameter-tuning extensions are available in the same repository.

[5] Our experimental setup matches that used for the experiments reported in [27], [28].

[6] Although L1-regularized approaches are known to produce sparse weight vectors, the non-zero weight density also depends on the regularization parameter selected during tuning.

| Classifier | Tuning | Accuracy | Runtime (s) | Density |
|---|---|---|---|---|
| spGAMP: BG-PR | EM | 97.6% | )$/ ./š | 11.1% |
| spGAMP: BG-HL | EM | 3/ / š | 468 / 93 | 1 š |
| msGAMP: L1-LR | EM | 97.6% | 684 / 123 | 9.8% |
| msGAMP: L1-LR | xval | 97.6% | 3068 / 278 | 19.6% |
| CDN | xval | 3/ / š | 1298 / 112 | 10.9% |
| TRON | xval | 3/ / š | 1682 / 133 | 10.8% |
| TFOCS: L1-LR | xval | 97.6% | 1086 / 94 | 19.2% |

the goal, then weight-vector density is also important, and for this we see significant differences among algorithms. Notably, the two most discriminative (i.e., sparse) classifiers are the result of EM-GAMP methods. We attribute this sparsity in large part to EM-tuning, since Table I shows that when cross-validation is used in place of EM-tuning with L1-LR GAMP, the weight vector is twice as dense.

Table I also shows a wide range of runtimes. Among the post-tuning runtimes, the two fastest are EM-GAMP based. Moreover, among the total runtimes, the three fastest are EM-GAMP based, with the best (at 3 $T$ seconds) beating the fastest non-GAMP algorithm (at 1 0 8 seconds) by more than a factor of 3 . That said, some caution must be used when comparing runtimes. For example, while all algorithms were given a "stopping tolerance" of 1 0 $^-$ 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).

Finally, we note that, although GAMP was derived under the assumption that the elements of $X$ are realizations of a an i.i.d. sub-Gaussian distribution, it worked well even with the $X$ of this experiment, which was far from i.i.d. sub-Gaussian. We attribute the robust performance of GAMP to the adaptive damping mechanism included in the GAMPmatlab implementation (described in [31]).

[1] C. M. Bishop, *Pattern Recognition and Machine Learning.* Springer, 2006.

[2] *J. Mach. Learn. Res.*, jb^ q ɑɕc q q

[3] *Science*, jb^ q ɑɕc q q

[4] *NeuroImage*, jb^ q ɑɕc q q

[5] *NeuroImage*, jb^ q ɑɕc q q

[6] *Conjoint Measurement: Methods and Applications.* Springer-Verlag, q q

[7] *Int'l Wkshp. Mach. Learn.*, cc q q q

[8] *Proc. Conf. Inform. Science Sys.*, 6e[aLO ba0 e q q

[9] *arXiv: 1202.1212*, $OJ q q q

[10] *Proc. IEEE Int'l Symp. Inform. Theory*, arXiv:1010.5141).

[11] *IEEE Trans. Inform. Theory*, jb^ q ɑɕc q q $OJ q q

[12] *Adv. Neural Info. Process. Sys.*, cc q q q

[13] *arXiv:1211.5164*, 2bj q q

[14] *Proc. IEEE Int'l Symp. Inform. Theory*, arXiv:1301.6295).

[15] *IEEE Signal Process. Mag.*, jb^ q ɑɕc q q Ha q q

[16] *Message Passing Approaches to Compressive Inference Under Structured Signal Priors.* q q

[17] *Gaussian Processes for Machine Learning.* q q

[18] *Conf. on Information Sciences and Systems (CISS)*, 6e[aLO ba q q q 0He q q

[19] *IEEE Trans. Signal Process.*, jb^ q q cc q q

[20] *Proceedings of the National Academy of Sciences*, jb^ q ɑɕc q q 2bj q q

[21] *R. Stat. Statist. Soc.*, ? š b^ q qɑb q ɑɕc q q q

[22] *Machine Learning*, jb^ q ɑɕc q q q q

[23] *Proc. IEEE Workshop Comp. Adv. Multi-Sensor Adaptive Process.*, arXiv:1310.2806).

[24] *R. Stat. Statist. Soc.*, ? š jb^ q ɑɕc q q

[25] *Proc. Neural Inform. Process. Syst. Conf.*, arXiv:1207.3859).

[26] *J. Mach. Learn. Res.*, jb^ q q cc q q q

[27] *J. Mach. Learn. Res.*, jb^ q 'q ɑɕc q q q

[28] *Proc. 24th Int'l Conf. Mach. Learn.*, bej H^^ q q c q q

[29] *SIAM J. Optim.*, jb^ q qɑɕc q q q

[30] *Math. Prog. Comp.*, jb^ qɑb q ɑɕc q q

[31] *Allerton Conf. Commun., Control, Comput.*, q q q