

# Sketched Clustering via Hybrid Approximate Message Passing

Evan Byrne, Antoine Chatalic, Rémi Gribonval, Philip Schniter

► **To cite this version:**

Evan Byrne, Antoine Chatalic, Rémi Gribonval, Philip Schniter. Sketched Clustering via Hybrid Approximate Message Passing. *IEEE Transactions on Signal Processing*, Institute of Electrical and Electronics Engineers, 2019, 67 (17), pp.4556-4569. 10.1109/TSP.2019.2924585 . hal-01991231v2

**HAL Id: hal-01991231**

**<https://hal.inria.fr/hal-01991231v2>**

Submitted on 5 Jul 2019

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Sketched Clustering via Hybrid Approximate Message Passing

Evan Byrne,<sup>\*</sup> Antoine Chatalic,<sup>†</sup> Rémi Gribonval<sup>†</sup>, *IEEE Fellow*, and Philip Schniter<sup>\*</sup>, *IEEE Fellow*

**Abstract**—In sketched clustering, a dataset of  $T$  samples is first sketched down to a vector of modest size, from which the centroids are subsequently extracted. Advantages include i) reduced storage complexity and ii) centroid extraction complexity independent of  $T$ . For the sketching methodology recently proposed by Keriven et al., which can be interpreted as a random sampling of the empirical characteristic function, we propose a sketched clustering algorithm based on approximate message passing. Numerical experiments suggest that our approach is more efficient than the state-of-the-art sketched clustering algorithm “CL-OMPR” (in both computational and sample complexity) and more efficient than k-means++ when  $T$  is large.

**Index Terms**—clustering algorithms, data compression, compressed sensing, approximate message passing

## I. INTRODUCTION

Given a dataset  $\mathbf{X} \triangleq [\mathbf{x}_1, \dots, \mathbf{x}_T] \in \mathbb{R}^{N \times T}$  comprising  $T$  samples of dimension  $N$ , the standard clustering problem is to find  $K$  centroids  $\mathbf{C} \triangleq [\mathbf{c}_1, \dots, \mathbf{c}_K] \in \mathbb{R}^{N \times K}$  that minimize the sum of squared errors (SSE)

$$\text{SSE}(\mathbf{X}, \mathbf{C}) \triangleq \frac{1}{T} \sum_{t=1}^T \min_k \|\mathbf{x}_t - \mathbf{c}_k\|_2^2. \quad (1)$$

Finding the optimal  $\mathbf{C}$  is an NP-hard problem [1]. Thus, many heuristic approaches have been proposed, such as the *k-means* algorithm [2,3]. Because k-means can get trapped in bad local minima, robust variants have been proposed, such as *k-means++* [4], which uses a careful random initialization procedure to yield solutions with SSE that have on average  $\leq 8(\ln K + 2)$  times the minimal SSE. The computational complexity of k-means++ scales as  $O(TKNI)$ , with  $I$  the number of iterations, which is impractical when  $T$  is large.

### A. Sketched Clustering

In *sketched clustering* [5,6,7], the dataset  $\mathbf{X}$  is first sketched down to a vector  $\mathbf{y}$  with  $M = O(KN)$  components, from

<sup>\*</sup>E. Byrne (byrne.133@osu.edu) and P. Schniter (schniter.1@osu.edu) are with the Department of Electrical and Computer Engineering at The Ohio State University, Columbus, OH, USA.

<sup>†</sup>A. Chatalic (antoine.chatalic@irisa.fr) and R. Gribonval (remi.gribonval@inria.fr) are with Univ. Rennes, Inria, CNRS, IRISA, France. A. Chatalic received a travel grant from the French research network GdR MIA to visit The Ohio State University.

Please direct all correspondence to Philip Schniter, Dept. ECE, 2015 Neil Ave., Columbus OH 43210, phone 614.247.6488, fax 614.292.7596.

<sup>\*</sup>E. Byrne and P. Schniter acknowledge support from NSF grant 1716388 and MIT Lincoln Labs.

Portions of this work were presented at the 2017 Asilomar Conference of Signals, Systems, and Computers.

which the centroids  $\mathbf{C}$  are subsequently extracted. In the typical case that  $K \ll T$ , the sketch consumes much less memory than the original dataset. If the sketch can be performed efficiently, then—since the complexity of centroid-extraction is invariant to  $T$ —sketched clustering may be more efficient than direct clustering methods when  $T$  is large. Note, for example, that k-means++ processes the  $T$  data samples in  $\mathbf{X}$  at every iteration, whereas sketched clustering processes the  $T$  data samples in  $\mathbf{X}$  only once, during the sketching step.

In this work, we focus on sketches of the type proposed by Keriven et al. in [5,6], which use  $\mathbf{y} = [y_1, \dots, y_M]^T$  with

$$y_m = \frac{1}{T} \sum_{t=1}^T \exp(j\mathbf{w}_m^T \mathbf{x}_t) \quad (2)$$

and randomly<sup>1</sup> generated  $\mathbf{W} \triangleq [\mathbf{w}_1, \dots, \mathbf{w}_M]^T \in \mathbb{R}^{M \times N}$ . Note that  $y_m$  in (2) can be interpreted as a sample of the empirical characteristic function [8], i.e.,

$$\phi(\mathbf{w}_m) = \int_{\mathbb{R}^N} p(\mathbf{x}) \exp(j\mathbf{w}_m^T \mathbf{x}) d\mathbf{x} \quad (3)$$

under the empirical distribution  $p(\mathbf{x}) = \frac{1}{T} \sum_{t=1}^T \delta(\mathbf{x} - \mathbf{x}_t)$ , with Dirac  $\delta(\cdot)$ . Here, each  $\mathbf{w}_m$  can be interpreted as a multidimensional frequency sample. The process of sketching  $\mathbf{X}$  down to  $\mathbf{y}$  via (2) costs  $O(TMN)$  operations, but it can be performed efficiently in an online and/or distributed manner.

To recover the centroids  $\mathbf{C}$  from  $\mathbf{y}$ , the state-of-the-art algorithm is *compressed learning via orthogonal matching pursuit with replacement* (CL-OMPR) [5,6]. It aims to solve

$$\arg \min_{\mathbf{C}} \min_{\alpha: \mathbf{1}^T \alpha = 1} \sum_{m=1}^M \left| y_m - \sum_{k=1}^K \alpha_k \exp(j\mathbf{w}_m^T \mathbf{c}_k) \right|^2 \quad (4)$$

using a greedy heuristic inspired by the *orthogonal matching pursuit* (OMP) algorithm [9] popular in compressed sensing. With sketch length  $M \geq 10KN$ , CL-OMPR typically recovers centroids of similar or better quality to those attained with k-means++. One may wonder, however, whether it is possible to recover accurate centroids with sketch lengths closer to the counting bound  $M = KN$ . Also, since CL-OMPR’s computational complexity is  $O(MNK^2)$ , one may wonder whether it is possible to recover accurate centroids with computational complexity  $O(MNK)$ .

<sup>1</sup>In [5] it was proposed to generate  $\{\mathbf{w}_m\}$  as independent draws from a distribution for which  $\mathbf{w}_m / \|\mathbf{w}_m\|$  is uniformly distributed on the unit sphere but  $\|\mathbf{w}_m\|$  has a prescribed density. More details are given in Section II-A.

## B. Contributions

To recover the centroids  $\mathbf{C}$  from a sketch  $\mathbf{y}$  of the form in (2), we propose the *compressive learning via approximate message passing* (CL-AMP) algorithm, with computational complexity  $O(MNK)$ . Numerical experiments show that, in most cases, CL-AMP accurately recovers centroids from sketches of length  $M \geq 2KN$ . This is an improvement over CL-OMPR, which typically requires  $M \geq 10KN$ . Our experiments establish these behaviors over many combinations of  $K \in [5, 50]$ ,  $N \in [10, 300]$ , and sample numbers  $T \in [10^5, 10^8]$ . Experiments also show that CL-AMP recovers centroids faster and more accurately than k-means++ when  $T$  is large, e.g.,  $T \geq 10^7$  in our numerical experiments.

We proposed a simple incarnation of the CL-AMP algorithm in the conference paper [10], with derivation details omitted due to space limitations. In this paper, we present the full derivation of CL-AMP with an improved initialization and hyperparameter tuning scheme, and a much more comprehensive set of numerical experiments.

The remainder of the paper is organized as follows. In Section II, we derive CL-AMP after reviewing relevant background on approximate message passing (AMP) algorithms. In Section III, we present numerical experiments using synthetic and MNIST data, and we apply CL-AMP to multidimensional frequency estimation. In Section IV, we conclude.

## II. COMPRESSIVE LEARNING VIA AMP

### A. High-Dimensional Inference Framework

CL-AMP treats centroid recovery as a high-dimensional inference problem rather than an optimization problem like minimizing (1) or (4). In particular, it models the data using a Gaussian mixture model (GMM)

$$\mathbf{x}_t \sim \sum_{k=1}^K \alpha_k \mathcal{N}(\mathbf{c}_k, \Phi_k), \quad (5)$$

where the centroids  $\mathbf{c}_k$  act as the GMM means, and the GMM weights  $\alpha_k$  and covariance matrices  $\Phi_k$  are treated as unknown parameters. That is,  $\{\mathbf{x}_t\}_{t=1}^T$  are assumed to be drawn i.i.d. from the GMM distribution (5). To recover the centroids  $\mathbf{C} \triangleq [\mathbf{c}_1, \dots, \mathbf{c}_K]$  from  $\mathbf{y}$ , CL-AMP computes an approximation to the MMSE estimate

$$\hat{\mathbf{C}} = \mathbb{E}\{\mathbf{C} | \mathbf{y}\}, \quad (6)$$

where the expectation is taken over the posterior density

$$p(\mathbf{C} | \mathbf{y}) \propto p(\mathbf{y} | \mathbf{C}) p(\mathbf{C}). \quad (7)$$

In (7),  $p(\mathbf{y} | \mathbf{C})$  is the likelihood function of  $\mathbf{C}$ , and  $p(\mathbf{C})$  is the prior density on  $\mathbf{C}$ . The dependence of  $p(\mathbf{y} | \mathbf{C})$  on  $\{\alpha_k\}$  and  $\{\Phi_k\}$  will be detailed in the sequel.

As we now establish, the form of the sketch in (2) implies that, conditioned on the centroids  $\mathbf{C}$  and the frequencies  $\mathbf{W}$ , the elements of  $\mathbf{y}$  can be treated as i.i.d. In other words, the sketch  $\mathbf{y}$  follows a generalized linear model (GLM) [11]. To establish this result, let us first define the normalized frequency vectors

$$\mathbf{a}_m \triangleq \mathbf{w}_m / g_m \text{ with } g_m \triangleq \|\mathbf{w}_m\| \quad (8)$$

and the (normalized) transform outputs

$$\mathbf{z}_m^\top \triangleq \mathbf{a}_m^\top \mathbf{C} \in \mathbb{R}^K. \quad (9)$$

Then  $p(\mathbf{y} | \mathbf{C})$  takes the form of a GLM, i.e.,

$$p(\mathbf{y} | \mathbf{C}) = \prod_{m=1}^M p_{y|z}(y_m | \mathbf{a}_m^\top \mathbf{C}), \quad (10)$$

for a conditional pdf  $p_{y|z}$  that will be detailed in the sequel.

From (2) and the definitions of  $\mathbf{a}_m$  and  $g_m$  in (8), we have

$$y_m = \frac{1}{T} \sum_{t=1}^T \exp(j\mathbf{w}_m^\top \mathbf{x}_t) \quad (11)$$

$$\approx \mathbb{E} \left\{ \exp(j\mathbf{w}_m^\top \mathbf{x}_t) \mid \mathbf{w}_m \right\} \quad (12)$$

$$= \sum_{k=1}^K \alpha_k \exp \left( jg_m \underbrace{\mathbf{a}_m^\top \mathbf{c}_k}_{\triangleq z_{mk}} - \frac{g_m^2}{2} \underbrace{\mathbf{a}_m^\top \Phi_k \mathbf{a}_m}_{\triangleq \tau_{mk}} \right), \quad (13)$$

where (12) holds under large  $T$  and (13) follows from the fact

$$\mathbf{w}_m^\top \mathbf{x}_t \mid \mathbf{w}_m \sim \sum_{k=1}^K \alpha_k \mathcal{N}(g_m z_{mk}, g_m^2 \tau_{mk}) \quad (14)$$

under (5), and the following well-known result [12, p.153]:

$$\mathbb{E}\{e^{jx}\} = \exp(j\mu - \sigma^2/2) \text{ when } x \sim \mathcal{N}(\mu, \sigma^2). \quad (15)$$

For  $\mathbf{a}_m$  distributed uniformly on the sphere, the elements  $\{\tau_{mk}\}_{m=1}^M$  in (13) concentrate as  $N \rightarrow \infty$  [13], in that

$$\tau_{mk} \xrightarrow{P} \mathbb{E}\{\tau_{mk}\} = \text{tr}(\Phi_k)/N \triangleq \tau_k, \quad (16)$$

as long as the peak-to-average eigenvalue ratio of  $\Phi_k$  remains bounded. Thus, for large  $T$  and  $N$ , (13) and (16) imply that

$$y_m = \sum_{k=1}^K \alpha_k \exp \left( jg_m z_{mk} - \frac{g_m^2 \tau_k}{2} \right), \quad (17)$$

which implies that the inference problem depends on the covariance matrices  $\{\Phi_k\}$  only through the hyperparameters  $\{\tau_k\}$ . Equation (17) can then be rephrased as

$$p_{y|z}(y_m | z_m; \boldsymbol{\alpha}, \boldsymbol{\tau}) = \delta \left( y_m - \sum_{k=1}^K \alpha_k \exp \left( jg_m z_{mk} - \frac{g_m^2 \tau_k}{2} \right) \right), \quad (18)$$

where  $\boldsymbol{\tau} \triangleq [\tau_1, \dots, \tau_K]^\top$  and  $\boldsymbol{\alpha} \triangleq [\alpha_1, \dots, \alpha_K]^\top$  are hyperparameters of the GLM that will be estimated from  $\mathbf{y}$ .

For the CL-AMP framework, any prior of the form

$$p(\mathbf{C}) = \prod_{n=1}^N p_{\mathbf{c}}(\mathbf{c}_n^\top) \quad (19)$$

is admissible, where (with some abuse of notation)  $\mathbf{c}_n^\top$  denotes the  $n$ th row of  $\mathbf{C}$ . For all experiments in Section III, we used the trivial prior  $p(\mathbf{C}) \propto 1$ .

In summary, CL-AMP aims to compute the MMSE estimate of  $\mathbf{C} \in \mathbb{R}^{N \times K}$  from the sketch  $\mathbf{y} \in \mathbb{C}^M$  under the prior  $\mathbf{C} \sim \prod_{n=1}^N p_{\mathbf{c}}(\mathbf{c}_n)$  from (19) and the likelihood  $\mathbf{y} \sim \prod_{m=1}^M p_{y|z}(y_m | z_m; \boldsymbol{\alpha}, \boldsymbol{\tau})$  from (18), where  $\mathbf{z}_m^\top$  is the  $m$ th row of  $\mathbf{Z} = \mathbf{A}\mathbf{C} \in \mathbb{R}^{M \times K}$  and  $\mathbf{A} \in \mathbb{R}^{M \times N}$  is a large

random matrix with rows  $\{\mathbf{a}_m^\top\}$  distributed uniformly on the unit sphere. CL-AMP estimates the values of  $\alpha$  and  $\tau$  from the sketch prior to estimating  $\mathbf{C}$ , as detailed in the sequel.

As proposed in [5], the row-norms  $\{g_m\}$  from (8) were drawn i.i.d. from the distribution

$$p(g; \sigma^2) \propto 1_{[0, \infty)}(g) \sqrt{g^2 \sigma^2 + \frac{g^4 \sigma^4}{4}} \exp\left(-\frac{1}{2} g^2 \sigma^2\right) \quad (20)$$

with shape parameter  $\sigma^2$ . The authors in [5] suggest using  $\sigma^2 = \frac{1}{NK} \sum_{k=1}^K \text{tr}(\Phi_k)$  and propose a method to estimate  $\sigma^2$  from  $\mathbf{y}$ . However, our numerical experiments suggest that the simpler assignment

$$\sigma^2 = \mathbb{E}\{\|\mathbf{x}\|_2^2\}/N \approx \|\mathbf{X}\|_F^2/NT \quad (21)$$

provides significantly improved performance. Note that the right side of (21) can be computed in an online manner, or approximated using a subset of the data.

### B. Approximate Message Passing

Exactly computing the MMSE estimate of  $\mathbf{C}$  from  $\mathbf{y}$  is impractical due to the form of  $p_{\mathbf{y}|\mathbf{z}}$ . Instead, one might consider approximate inference via the sum-product algorithm (SPA), but even the SPA is intractable due to the form of  $p_{\mathbf{y}|\mathbf{z}}$ . Given the presence of a large random matrix  $\mathbf{A}$  in the problem formulation, we instead leverage *approximate message passing* (AMP) methods. In particular, we propose to apply the *simplified hybrid generalized AMP* (SHyGAMP) methodology from [14], while simultaneously estimating  $\alpha$  and  $\tau$  through expectation maximization (EM). A brief background on AMP methods will now be provided to justify our approach.

The original AMP algorithm of Donoho, Maleki, and Montanari [15] was designed to estimate i.i.d.  $c$  under the standard linear model (i.e.,  $\mathbf{y} = \mathbf{A}\mathbf{c} + \mathbf{n}$  with known  $\mathbf{A} \in \mathbb{R}^{M \times N}$  and additive white Gaussian noise  $\mathbf{n}$ ). The generalized AMP (GAMP) algorithm of Rangan [16] extended AMP to the generalized linear model (i.e.,  $\mathbf{y} \sim p(\mathbf{y}|\mathbf{z})$  for  $\mathbf{z} = \mathbf{A}\mathbf{c}$  and separable  $p(\mathbf{y}|\mathbf{z}) = \prod_{m=1}^M p(y_m|z_m)$ ). Both AMP and GAMP give accurate approximations of the SPA under large i.i.d. sub-Gaussian  $\mathbf{A}$ , while maintaining a computational complexity of only  $O(MN)$ . Furthermore, both can be rigorously analyzed via the state-evolution framework, which proves that they compute MMSE optimal estimates of  $c$  in certain regimes [17].

A limitation of AMP [15] and GAMP [16] is that they treat only problems with i.i.d. estimand  $c$  and separable likelihood  $p(\mathbf{y}|\mathbf{z}) = \prod_{m=1}^M p(y_m|z_m)$ . Thus, *Hybrid GAMP* (HyGAMP) [18] was developed to tackle problems with a structured prior and/or likelihood. HyGAMP could be applied to the compressive learning problem described in Section II-A, but it would require computing and inverting  $O(N+M)$  covariance matrices of dimension  $K$  at each iteration. For this reason, we instead apply the *simplified HyGAMP* (SHyGAMP) algorithm from [14], which uses diagonal covariance matrices in HyGAMP to reduce its computational complexity. As described in [14], SHyGAMP can be readily combined with the EM algorithm to learn the hyperparameters  $\alpha$  and  $\tau$ .

### C. SHyGAMP

The SHyGAMP algorithm was proposed and described in detail in [14]; we provide only a brief review here. Algorithm 1 summarizes the SHyGAMP algorithm using the language of Section II-A. In lines 10-11, with some abuse of notation, we use  $\mathbf{c}_n^\top$  to denote the  $n$ th row of the centroid matrix  $\mathbf{C}$  (where in (5) we used  $\mathbf{c}_k$  to denote the  $k$ th column of  $\mathbf{C}$ ). We also use  $\hat{\mathbf{P}} \triangleq [\hat{\mathbf{p}}_1, \dots, \hat{\mathbf{p}}_M]^\top$ ,  $\hat{\mathbf{Z}} \triangleq [\hat{\mathbf{z}}_1, \dots, \hat{\mathbf{z}}_M]^\top$ ,  $\hat{\mathbf{R}} \triangleq [\hat{\mathbf{r}}_1, \dots, \hat{\mathbf{r}}_N]^\top$ ,  $\oslash$  for componentwise division, and  $\odot$  for componentwise multiplication. In the sequel, covariance matrices will be denoted by (superscripted)  $\mathbf{Q}$  and vectors of their diagonal elements denoted by (superscripted)  $\mathbf{q}$ . A brief interpretation of SHyGAMP is now provided.

At each iteration, lines 4-5 of Algorithm 1 generate the posterior mean and covariance of the transform outputs  $\mathbf{z}_m$  from (9) under a likelihood  $p_{\mathbf{y}|\mathbf{z}}$  like (18) and the ‘‘pseudo’’ prior  $\mathbf{z}_m \sim \mathcal{N}(\hat{\mathbf{p}}_m, \mathbf{Q}^{\mathbf{P}})$ , where  $\hat{\mathbf{p}}_m$  and  $\mathbf{Q}^{\mathbf{P}} = \text{Diag}(\mathbf{q}^{\mathbf{P}})$  are updated at each SHyGAMP iteration. Thus, the pdf used for the covariance and expectation in lines 4-5 is

$$\begin{aligned} p_{\mathbf{z}|\mathbf{y}, \mathbf{p}}(\mathbf{z}_m | y_m, \hat{\mathbf{p}}_m; \mathbf{Q}^{\mathbf{P}}, \alpha, \tau) \\ = \frac{p_{\mathbf{y}|\mathbf{z}}(y_m | \mathbf{z}_m; \alpha, \tau) \mathcal{N}(\mathbf{z}_m; \hat{\mathbf{p}}_m, \mathbf{Q}^{\mathbf{P}})}{\int p_{\mathbf{y}|\mathbf{z}}(y_m | \mathbf{z}'_m; \alpha, \tau) \mathcal{N}(\mathbf{z}'_m; \hat{\mathbf{p}}_m, \mathbf{Q}^{\mathbf{P}}) d\mathbf{z}'_m}. \end{aligned} \quad (22)$$

Similarly, lines 10-11 compute the posterior mean and covariance of  $\mathbf{c}_n$  under a prior  $p_{\mathbf{c}}$  of the form (19) and ‘‘pseudo’’ measurements  $\hat{\mathbf{r}}_n$  that follow the statistical model

$$\hat{\mathbf{r}}_n = \mathbf{c}_n + \mathbf{v}_n, \quad \mathbf{v}_n \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}^{\mathbf{r}}), \quad (23)$$

where  $\hat{\mathbf{r}}_n$  and  $\mathbf{Q}^{\mathbf{r}} = \text{Diag}(\mathbf{q}^{\mathbf{r}})$  are updated at each SHyGAMP iteration. Thus, the pdf used for the covariance and expectation in lines 10-11 is

$$p_{\mathbf{c}|\mathbf{r}}(\mathbf{c}_n | \hat{\mathbf{r}}_n; \mathbf{Q}^{\mathbf{r}}) = \frac{p_{\mathbf{c}}(\mathbf{c}_n) \mathcal{N}(\mathbf{c}_n; \hat{\mathbf{r}}_n, \mathbf{Q}^{\mathbf{r}})}{\int p_{\mathbf{c}}(\mathbf{c}'_n) \mathcal{N}(\mathbf{c}'_n; \hat{\mathbf{r}}_n, \mathbf{Q}^{\mathbf{r}}) d\mathbf{c}'_n}. \quad (24)$$

As the SHyGAMP iterations progress, the output  $[\hat{\mathbf{c}}_1, \dots, \hat{\mathbf{c}}_N]^\top$  of line 11 converges to an approximation of the MMSE estimate  $\mathbb{E}\{\mathbf{C}|\mathbf{y}\}$ , and the output  $[\hat{\mathbf{z}}_1, \dots, \hat{\mathbf{z}}_M]^\top$  of line 5 converges to an approximation of the MMSE estimate  $\mathbb{E}\{\mathbf{Z}|\mathbf{y}\}$ . Essentially, the SHyGAMP algorithm breaks an inference problem of dimension  $NK$  into  $O(M+N)$  inference problems of dimension  $K$  (i.e., lines 4-5 and 10-11 of Algorithm 1), each involving an independent-Gaussian pseudo-prior or pseudo-likelihood, evaluated iteratively. The computational complexity of SHyGAMP is  $O(MNK)$ .

### D. From SHyGAMP to CL-AMP

The SHyGAMP algorithm can be applied to many different problems via appropriate choice of  $p_{\mathbf{y}|\mathbf{z}}$  and  $p_{\mathbf{c}}$ . To apply SHyGAMP to sketched clustering, we choose  $p_{\mathbf{y}|\mathbf{z}}$  and  $p_{\mathbf{c}}$  as described in Section II-A. As we will see, the main challenge is evaluating lines 4-5 of Algorithm 1 for the  $p_{\mathbf{y}|\mathbf{z}}$  in (18).

1) *Inference of  $\mathbf{z}_m$* : For lines 4-5 of Algorithm 1, we would like to compute the mean and variance

$$\hat{\mathbf{z}}_{mk} = \frac{\int_{\mathbb{R}^K} z_{mk} p_{\mathbf{y}|\mathbf{z}}(y_m | \mathbf{z}_m) \mathcal{N}(\mathbf{z}_m; \hat{\mathbf{p}}_m, \mathbf{Q}^{\mathbf{P}}) d\mathbf{z}_m}{C_m} \quad (25)$$

---

**Algorithm 1** SHyGAMP
 

---

**Require:** Measurements  $\mathbf{y} \in \mathbb{C}^M$ , matrix  $\mathbf{A} \in \mathbb{R}^{M \times N}$  with  $\|\mathbf{A}\|_F^2 = M$ , pdfs  $p_{\text{c|r}}(\cdot|\cdot)$  and  $p_{\mathbf{z}|y, \mathbf{p}}(\cdot|\cdot, \cdot; \boldsymbol{\alpha}, \boldsymbol{\tau})$  from (22) and (24), initial  $\widehat{\mathbf{C}}_0 \in \mathbb{R}^{N \times K}$  and  $\mathbf{q}^{\mathbf{p}} = \mathbf{q}_0^{\mathbf{p}} \in \mathbb{R}_+^K$ .

- 1:  $\widehat{\mathbf{S}} \leftarrow \mathbf{0}$ ,  $\widehat{\mathbf{C}} \leftarrow \widehat{\mathbf{C}}_0$ .
- 2: **repeat**
- 3:  $\widehat{\mathbf{P}} \leftarrow \mathbf{A}\widehat{\mathbf{C}} - \widehat{\mathbf{S}} \text{Diag}(\mathbf{q}^{\mathbf{p}})$
- 4:  $\mathbf{q}_m^{\mathbf{z}} \leftarrow \text{diag}(\text{Cov}\{z_m | y_m, \widehat{\mathbf{p}}_m; \text{Diag}(\mathbf{q}^{\mathbf{p}}), \boldsymbol{\alpha}, \boldsymbol{\tau}\})$ ,  $m = 1 \dots M$
- 5:  $\widehat{z}_m \leftarrow \mathbb{E}\{z_m | y_m, \widehat{\mathbf{p}}_m; \text{Diag}(\mathbf{q}^{\mathbf{p}}), \boldsymbol{\alpha}, \boldsymbol{\tau}\}$ ,  $m = 1 \dots M$
- 6:  $\mathbf{q}^{\mathbf{s}} \leftarrow \mathbf{1} \otimes \mathbf{q}^{\mathbf{p}} - (\frac{1}{M} \sum_{m=1}^M \mathbf{q}_m^{\mathbf{z}}) \otimes (\mathbf{q}^{\mathbf{p}} \odot \mathbf{q}^{\mathbf{p}})$
- 7:  $\widehat{\mathbf{S}} \leftarrow (\widehat{\mathbf{Z}} - \widehat{\mathbf{P}}) \text{Diag}(\mathbf{q}^{\mathbf{p}})^{-1}$
- 8:  $\mathbf{q}^{\mathbf{r}} \leftarrow \frac{N}{M} \mathbf{1} \otimes \mathbf{q}^{\mathbf{s}}$
- 9:  $\widehat{\mathbf{R}} \leftarrow \mathbf{C} + \mathbf{A}^T \widehat{\mathbf{S}} \text{Diag}(\mathbf{q}^{\mathbf{r}})$
- 10:  $\mathbf{q}_n^{\mathbf{c}} \leftarrow \text{diag}(\text{Cov}\{c_n | \widehat{\mathbf{r}}_n; \text{Diag}(\mathbf{q}^{\mathbf{r}})\})$ ,  $n = 1 \dots N$
- 11:  $\widehat{c}_n \leftarrow \mathbb{E}\{c_n | \widehat{\mathbf{r}}_n; \text{Diag}(\mathbf{q}^{\mathbf{r}})\}$ ,  $n = 1 \dots N$
- 12:  $\mathbf{q}^{\mathbf{p}} \leftarrow \frac{1}{N} \sum_{n=1}^N \mathbf{q}_n^{\mathbf{c}}$
- 13: **until** convergence
- 14: **return**  $\widehat{\mathbf{C}}$

---

$$q_{mk}^z = \frac{\int_{\mathbb{R}^K} (z_{mk} - \widehat{z}_{mk})^2 p_{y|\mathbf{z}}(y_m | z_m) \mathcal{N}(z_m; \widehat{\mathbf{p}}_m, \mathbf{Q}^{\mathbf{p}}) dz_m}{C_m}, \quad (26)$$

where  $q_{mk}^z$  is the  $k$ th element of  $\mathbf{q}_m^z$  and

$$C_m = \int_{\mathbb{R}^K} p_{y|\mathbf{z}}(y_m | z_m) \mathcal{N}(z_m; \widehat{\mathbf{p}}_m, \mathbf{Q}^{\mathbf{p}}) dz_m. \quad (27)$$

However, due to the form of  $p_{y|\mathbf{z}}$  in (18), we are not able to find closed-form expressions for  $\widehat{z}_{mk}$  or  $q_{mk}^z$ . Thus, we propose to approximate  $\widehat{z}_{mk}$  and  $q_{mk}^z$  by writing (17) as

$$y_m = \alpha_k \exp(-g_m^2 \tau_k / 2) \exp(jg_m z_{mk}) + \sum_{l \neq k} \alpha_l \exp(-g_m^2 \tau_l / 2) \exp(jg_m (z_{ml})) \quad (28)$$

and treating the sum over  $l$  as complex Gaussian. For the remainder of this section, we suppress the subscripts “ $m$ ” and “ $y|\mathbf{z}$ ” to simplify the notation.

We now give a brief sketch of the derivation. First, we write (28) as

$$y = \underbrace{\alpha_k \exp(-g^2 \tau_k / 2)}_{\triangleq \beta_k} \underbrace{\exp(jg(z_k + n_k))}_{\triangleq \theta_k} + \sum_{l \neq k} \underbrace{\alpha_l \exp(-g^2 \tau_l / 2)}_{= \beta_l} \underbrace{\exp(jg(z_l + n_l))}_{\triangleq v_l}. \quad (29)$$

Here we introduce i.i.d.  $n_k \sim \mathcal{N}(0, q^n)$ , which will allow us to leverage the Gaussian multiplication rule (see footnote 2) to bypass tedious linear algebra. Eventually we will take  $q^n \rightarrow 0$ , so that (29) matches (28). Next we derive expressions (42) and (48), which state  $\widehat{z}_k$  and  $q_k^z$  in terms of the posterior mean and variance on the  $2\pi$ -periodic quantity  $\theta_k$  in (29). By approximating the second term in (29) as Gaussian, the posterior of  $\theta_k$  takes the form of a generalized von Mises distribution, as summarized in (68). Because the posterior mean and variance of  $\theta_k$  are not computable in closed-form, we approximate them using numerical integration. Finally, we relate the posterior mean and variance of  $\theta_k$  back to  $\widehat{z}_k$  and  $q_k^z$ .

We now begin the derivation. First, we derive an expression

for the marginal posterior  $p(z_k | y)$  under the pseudo-prior  $z_k \sim \mathcal{N}(\widehat{p}_k, q_k^{\mathbf{p}}) \forall k$ . To start,

$$p(z_k | y) = \int_{\mathbb{R}^K} p(\mathbf{z}, \theta_k | y) d\theta_k d\mathbf{z}_{\setminus k} \quad (30)$$

$$= \frac{1}{p(y)} \int_{\mathbb{R}^K} p(y | \mathbf{z}, \theta_k) p(\theta_k | \mathbf{z}) p(\mathbf{z}) d\theta_k d\mathbf{z}_{\setminus k} \quad (31)$$

$$= \frac{1}{p(y)} \int_{\mathbb{R}^K} p(y | \mathbf{z}_{\setminus k}, \theta_k) \mathcal{N}(\theta_k; g z_k, g^2 q^n) \times \prod_{l=1}^K \mathcal{N}(z_l; \widehat{p}_l, q_l^{\mathbf{p}}) d\theta_k d\mathbf{z}_{\setminus k}, \quad (32)$$

where  $\mathbf{z}_{\setminus k} \triangleq [z_1, \dots, z_{k-1}, z_{k+1}, \dots, z_K]^T$ . A change-of-variables from  $z_l$  to  $\widetilde{z}_l \triangleq z_l - \widehat{p}_l$  for all  $l \neq k$  gives

$$p(z_k | y) = \frac{\mathcal{N}(z_k; \widehat{p}_k, q_k^{\mathbf{p}})}{p(y)} \int_{\mathbb{R}} \mathcal{N}(\theta_k; g z_k, g^2 q^n) \times \left[ \int_{\mathbb{R}^{K-1}} p(y | \widetilde{\mathbf{z}}_{\setminus k}, \theta_k) \prod_{l \neq k} \mathcal{N}(\widetilde{z}_l; 0, q_l^{\mathbf{p}}) d\widetilde{\mathbf{z}}_{\setminus k} \right] d\theta_k, \quad (33)$$

where  $p(y | \widetilde{\mathbf{z}}_{\setminus k}, \theta_k)$  is associated with the generative model

$$y = \beta_k \exp(j\theta_k) + \sum_{l \neq k} \beta_l \exp(jg(\widehat{p}_l + \widetilde{z}_l + n_l)) \quad (34)$$

with i.i.d.  $n_l \sim \mathcal{N}(0, q^n)$ . Now, because  $\widetilde{z}_l$  and  $n_l$  are (apriori) mutually independent zero-mean Gaussian variables, we can work directly with the sum  $\widetilde{n}_l \triangleq \widetilde{z}_l + n_l \sim \mathcal{N}(0, q_l^{\mathbf{p}} + q^n)$  and thus bypass the inner integral in (33). This allows us to write

$$p(z_k | y) = \frac{\mathcal{N}(z_k; \widehat{p}_k, q_k^{\mathbf{p}})}{p(y)} \int_{\mathbb{R}} \mathcal{N}(\theta_k; g z_k, g^2 q^n) p(y | \theta_k) d\theta_k, \quad (35)$$

where  $p(y | \theta_k)$  is associated with the generative model

$$y = \beta_k \exp(j\theta_k) + \sum_{l \neq k} \beta_l \underbrace{\exp(jg(\widehat{p}_l + \widetilde{n}_l))}_{= v_l} \quad (36)$$

with i.i.d.  $\widetilde{n}_l \sim \mathcal{N}(0, q_l^{\mathbf{p}} + q^n)$ . Recalling that  $y \in \mathbb{C}$ , it will sometimes be useful to write (36) as

$$\begin{aligned} \begin{bmatrix} \text{Re}\{y\} \\ \text{Im}\{y\} \end{bmatrix} &\sim \mathcal{N}\left(\beta_k \begin{bmatrix} \cos(\theta_k) \\ \sin(\theta_k) \end{bmatrix} + \sum_{l \neq k} \beta_l \mathbb{E}\left\{ \begin{bmatrix} \text{Re}\{v_l\} \\ \text{Im}\{v_l\} \end{bmatrix} \right\}, \right. \\ &\quad \left. \sum_{l \neq k} \beta_l^2 \text{Cov}\left\{ \begin{bmatrix} \text{Re}\{v_l\} \\ \text{Im}\{v_l\} \end{bmatrix} \right\} \right). \end{aligned} \quad (37)$$

To compute the posterior mean of  $z_k$ , (35) implies

$$\widehat{z}_k \triangleq \mathbb{E}\{z_k | y\} = \int_{\mathbb{R}} z_k p(z_k | y) dz_k \quad (38)$$

$$\begin{aligned} &= \frac{1}{p(y)} \int_{\mathbb{R}} \left[ \int_{\mathbb{R}} z_k \mathcal{N}(g z_k; \theta_k, g^2 q^n) \mathcal{N}(z_k; \widehat{p}_k, q_k^{\mathbf{p}}) dz_k \right] \\ &\quad \times p(y | \theta_k) d\theta_k \\ &= \int_{\mathbb{R}} \left[ \int_{\mathbb{R}} z_k \mathcal{N}\left(z_k; \frac{\theta_k/g + \widehat{p}_k}{q^n + \frac{1}{q_k^{\mathbf{p}}}}, \frac{1}{q^n + \frac{1}{q_k^{\mathbf{p}}}}\right) dz_k \right] \end{aligned} \quad (39)$$

$$\times \underbrace{\frac{\mathcal{N}(\theta_k; g\hat{p}_k, g^2(q^n + q_k^p)) p(y|\theta_k)}{p(y)}}_{= p(\theta_k|y)} d\theta_k \quad (40)$$

$$= \int_{\mathbb{R}} \frac{\frac{\theta_k/g}{q^n} + \frac{\hat{p}_k}{q_k^p}}{\frac{1}{q^n} + \frac{1}{q_k^p}} p(\theta_k|y) d\theta_k \quad (41)$$

$$= \frac{\hat{p}_k}{q_k^p/q^n + 1} + \frac{\hat{\theta}_k/g}{1 + q^n/q_k^p} \text{ for } \hat{\theta}_k \triangleq \int_{\mathbb{R}} \theta_k p(\theta_k|y) d\theta_k, \quad (42)$$

where the Gaussian pdf multiplication rule<sup>2</sup> was used in (40) and where  $\hat{\theta}_k$  denotes the posterior mean of  $\theta_k$ .

For the posterior variance of  $z_k$ , a similar approach gives

$$q_k^z \triangleq \text{var}\{z_k|y\} = \int_{\mathbb{R}} (z_k - \hat{z}_k)^2 p(z_k|y) dz_k \quad (43)$$

$$= \frac{1}{p(y)} \int_{\mathbb{R}} \left[ \int_{\mathbb{R}} (z_k - \hat{z}_k)^2 \mathcal{N}(gz_k; \theta_k, g^2 q^n) \times \mathcal{N}(z_k; \hat{p}_k, q_k^p) dz_k \right] p(y|\theta_k) d\theta_k \quad (44)$$

$$= \int_{\mathbb{R}} \left[ \int_{\mathbb{R}} (z_k - \hat{z}_k)^2 \mathcal{N}\left(z_k; \frac{\theta_k/g + \hat{p}_k}{\frac{1}{q^n} + \frac{1}{q_k^p}}, \frac{1}{\frac{1}{q^n} + \frac{1}{q_k^p}}\right) dz_k \right] \times p(\theta_k|y) d\theta_k. \quad (45)$$

Using a change-of-variables from  $z_k$  to  $\tilde{z}_k \triangleq z_k - \hat{z}_k$ , we get

$$q_k^z = \int_{\mathbb{R}} \left[ \int_{\mathbb{R}} \tilde{z}_k^2 \mathcal{N}\left(\tilde{z}_k; \frac{\theta_k/g - \hat{\theta}_k/g}{\frac{1}{q^n} + \frac{1}{q_k^p}}, \frac{1}{\frac{1}{q^n} + \frac{1}{q_k^p}}\right) d\tilde{z}_k \right] \times p(\theta_k|y) d\theta_k \quad (46)$$

$$= \int_{\mathbb{R}} \left[ \left( \frac{(\theta_k - \hat{\theta}_k)/g}{1 + q^n/q_k^p} \right)^2 + \frac{q^n}{1 + q^n/q_k^p} \right] p(\theta_k|y) d\theta_k \quad (47)$$

$$= \frac{q^n}{1 + q^n/q_k^p} + \frac{1}{g^2} \underbrace{\left( \frac{1}{1 + q^n/q_k^p} \right)^2 \int_{\mathbb{R}} (\theta_k - \hat{\theta}_k)^2 p(\theta_k|y) d\theta_k}_{\triangleq q_k^\theta = \text{var}\{\theta_k|y\}} \quad (48)$$

The computation of  $\hat{z}_k$  and  $q_k^z$  is still complicated by the form of the posterior  $p(\theta_k|y)$  implied by (36). To circumvent this problem, we propose to apply a Gaussian approximation to the sum in (36). Because  $\{\tilde{v}_l\}_{v_l \neq k}$  are mutually independent, the mean and covariance of the sum in (36) are simply the sum of the means and covariances (respectively) of the  $K-1$  terms making up the sum. Recalling (37), this implies that

$$p\left(\begin{bmatrix} \text{Re}\{y\} \\ \text{Im}\{y\} \end{bmatrix} \middle| \theta_k\right) \approx \mathcal{N}\left(\begin{bmatrix} \text{Re}\{y\} \\ \text{Im}\{y\} \end{bmatrix}; \beta_k \begin{bmatrix} \cos(\theta_k) \\ \sin(\theta_k) \end{bmatrix} + \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\right) \quad (49)$$

with

$$\boldsymbol{\mu}_k = \sum_{l \neq k} \alpha_l e^{-g^2(\tau_l + q_l^p)/2} \begin{bmatrix} \cos(g\hat{p}_l) \\ \sin(g\hat{p}_l) \end{bmatrix} \quad (50)$$

$$\boldsymbol{\Sigma}_k = \frac{1}{2} \sum_{l \neq k} \beta_l^2 (1 - e^{-g^2 q_l^p}) \times \left( \mathbf{I} - e^{-g^2 q_l^p} \begin{bmatrix} \cos(2g\hat{p}_l) & \sin(2g\hat{p}_l) \\ \sin(2g\hat{p}_l) & -\cos(2g\hat{p}_l) \end{bmatrix} \right). \quad (51)$$

We note that (50) and (51) were obtained using

$$\mathbb{E}\{\text{Re}\{v_l\}\} = \exp(-g^2 q_l^p/2) \cos(g\hat{p}_l) \quad (52)$$

$$\mathbb{E}\{\text{Im}\{v_l\}\} = \exp(-g^2 q_l^p/2) \sin(g\hat{p}_l) \quad (53)$$

$$2\mathbb{E}\{\text{Re}\{v_l\}^2\} = 1 + \exp(-g^2 q_l^p) \cos(2g\hat{p}_l) \quad (54)$$

$$2\mathbb{E}\{\text{Im}\{v_l\}^2\} = 1 - \exp(-g^2 q_l^p) \cos(2g\hat{p}_l) \quad (55)$$

$$2\mathbb{E}\{\text{Re}\{v_l\} \text{Im}\{v_l\}\} = \exp(-g^2 q_l^p) \sin(2g\hat{p}_l), \quad (56)$$

which use the fact that, after letting  $q^n \rightarrow 0$ ,

$$\mathbb{E}\{v_l\} = \int_{\mathbb{R}} \mathcal{N}(z_l; \hat{p}_l, q_l^p) \exp(jgz_l) dz_l \quad (57)$$

$$= \exp(jg\hat{p}_l - g^2 q_l^p/2). \quad (58)$$

Rewriting (49) as

$$p\left(\beta_k^{-1} \begin{bmatrix} \text{Re}\{y\} \\ \text{Im}\{y\} \end{bmatrix} \middle| \theta_k\right) \approx \mathcal{N}\left(\begin{bmatrix} \cos(\theta_k) \\ \sin(\theta_k) \end{bmatrix}; \beta_k^{-1} \begin{bmatrix} \text{Re}\{y\} \\ \text{Im}\{y\} \end{bmatrix} - \beta_k^{-1} \boldsymbol{\mu}_k, \beta_k^{-2} \boldsymbol{\Sigma}_k\right), \quad (59)$$

the right side of (59) can be recognized as being proportional to the generalized von Mises (GvM) density over  $\theta_k \in [0, 2\pi)$  from [19]. Under this GvM approximation, we have [19] that

$$p(y|\theta_k) \propto \exp(\kappa_k \cos(\theta_k - \zeta_k) + \bar{\kappa}_k \cos[2(\theta_k - \bar{\zeta}_k)]) \quad (60)$$

for parameters  $\kappa_k, \bar{\kappa}_k > 0$  and  $\zeta_k, \bar{\zeta}_k \in [0, 2\pi)$  defined from  $\beta_k^{-1} y$ ,  $\beta_k^{-1} \boldsymbol{\mu}_k$ , and  $\beta_k^{-2} \boldsymbol{\Sigma}_k$ . In particular,

$$\kappa_k \cos(\zeta_k) = -\frac{1}{1 - \rho_k^2} \left( \frac{\rho_k \bar{\nu}_k}{\sigma_k \bar{\sigma}_k} - \frac{\nu_k}{\sigma_k^2} \right) \quad (61)$$

$$\kappa_k \sin(\zeta_k) = -\frac{1}{1 - \rho_k^2} \left( \frac{\rho_k \nu_k}{\sigma_k \bar{\sigma}_k} - \frac{\bar{\nu}_k}{\bar{\sigma}_k^2} \right) \quad (62)$$

$$\bar{\kappa}_k \cos(2\bar{\zeta}_k) = -\frac{1}{4(1 - \rho_k^2)} \left( \frac{1}{\sigma_k^2} - \frac{1}{\bar{\sigma}_k^2} \right) \quad (63)$$

$$\bar{\kappa}_k \sin(2\bar{\zeta}_k) = \frac{\rho_k}{2(1 - \rho_k^2) \sigma_k \bar{\sigma}_k}, \quad (64)$$

where

$$\begin{bmatrix} \nu_k \\ \bar{\nu}_k \end{bmatrix} \triangleq \beta_k^{-1} \left( \begin{bmatrix} \text{Re}\{y\} \\ \text{Im}\{y\} \end{bmatrix} - \boldsymbol{\mu}_k \right) \quad (65)$$

$$\begin{bmatrix} \sigma_k^2 & \rho_k \sigma_k \bar{\sigma}_k \\ \rho_k \sigma_k \bar{\sigma}_k & \bar{\sigma}_k^2 \end{bmatrix} \triangleq \beta_k^{-2} \boldsymbol{\Sigma}_k. \quad (66)$$

From (60) and the SHyGAMP pseudo-prior  $z_k \sim \mathcal{N}(\hat{p}_k, q_k^p)$ , we see that the posterior on  $\theta_k$  takes the form

$$p(\theta_k|y) \propto \mathcal{N}(\theta_k; g\hat{p}_k, g^2 q_k^p) p(y|\theta_k) \quad (67)$$

$$\propto \exp\left[\kappa_k \cos(\theta_k - \zeta_k) + \bar{\kappa}_k \cos[2(\theta_k - \bar{\zeta}_k)] - \frac{(\theta_k - g\hat{p}_k)^2}{2g^2 q_k^p}\right]. \quad (68)$$

<sup>2</sup>According to the Gaussian multiplication rule, we have  $\mathcal{N}(\mathbf{x}; \mathbf{a}, \mathbf{A})\mathcal{N}(\mathbf{x}; \mathbf{b}, \mathbf{B}) = \mathcal{N}(\mathbf{0}; \mathbf{a} - \mathbf{b}, \mathbf{A} + \mathbf{B})\mathcal{N}(\mathbf{x}; (\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1}(\mathbf{A}^{-1}\mathbf{a} + \mathbf{B}^{-1}\mathbf{b}), (\mathbf{A}^{-1} + \mathbf{B}^{-1})^{-1})$ .

We now face the task of computing  $\hat{\theta}_k = \mathbb{E}\{\theta_k|y\}$  and  $q_k^\theta = \text{var}\{\theta_k|y\}$  under (68). Since these quantities do not appear to be computable in closed form, we settle for an approximation, such as that based on the Laplace approximation [20] or numerical integration. For the Laplace approximation, we would first compute  $\hat{\theta}_{k,\text{MAP}} \triangleq \arg \max_{\theta_k} \ln p(\theta_k|y)$  and then approximate  $\hat{\theta}_k \approx \hat{\theta}_{k,\text{MAP}}$  and  $q_k^\theta \approx -\frac{d^2}{d\theta_k^2} \ln p(\theta_k|y)|_{\theta_k=\hat{\theta}_{k,\text{MAP}}}$ . However, since computing  $\arg \max_{\theta_k} \ln p(\theta_k|y)$  is complicated due to the presence of multiple local maxima, we instead use numerical integration. For this, we suggest a grid of  $N_{\text{pts}}N_{\text{per}} + 1$  uniformly-spaced points centered at  $g\hat{p}_k$  with width  $2\pi N_{\text{per}}$ , where  $N_{\text{per}} = \left\lceil \frac{N_{\text{std}}}{\pi} \sqrt{g^2 q_k^\theta} \right\rceil$ . This choice of grid ensures that the sampling points cover at least  $N_{\text{std}}$  standard deviations of the prior on  $\theta_k$ . We used  $N_{\text{std}} = 4$  and  $N_{\text{pts}} = 7$  in the numerical experiments in Section III.

Finally, after approximating  $\hat{\theta}_k$  and  $q_k^\theta$  via numerical integration, we set  $\hat{z}_k = \hat{\theta}_k/g$  and  $q_k^z = q_k^\theta/g^2$ .

2) *Inference of  $c_n$* : Recall that lines 10-11 of Algorithm 1 support an arbitrary prior  $p_{\mathbf{c}}$  on  $c_n$ . For the experiments in Section III, we used the trivial non-informative prior  $p_{\mathbf{c}}(c_n) \propto 1$ , after which lines 10-11 reduce to

$$q_n^{\mathbf{c}} = \mathbf{q}^{\mathbf{r}} \quad \forall n \quad \text{and} \quad \hat{c}_n = \hat{\mathbf{r}}_n \quad \forall n. \quad (69)$$

### E. Initialization

We recommend initializing CL-AMP with  $\hat{\mathbf{C}} = \hat{\mathbf{C}}_0$  and  $\mathbf{q}^{\mathbf{P}} = \mathbf{q}_0^{\mathbf{P}}$ , where  $\hat{\mathbf{C}}_0$  is drawn i.i.d.  $\mathcal{N}(0, \sigma^2)$  and where  $\mathbf{q}_0^{\mathbf{P}} = \sigma^2 \mathbf{1}$ , with  $\sigma^2$  from (21) (as described in Section II-A).

In some cases, running CL-AMP from  $R > 1$  different random initializations can help to avoid spurious solutions. Here, CL-AMP is run from a different random initialization  $\hat{\mathbf{C}}_{0,r}$ , for  $r = 1, \dots, R$ , and then the quality of the recovered solution  $\hat{\mathbf{C}}_r$  is evaluated by constructing the “estimated sketch”  $\hat{\mathbf{y}}_r$  via

$$\hat{\mathbf{y}}_{mr} = \sum_{k=1}^K \alpha_k \exp(-g_m^2 \tau_k) \exp(jg_m \mathbf{a}_m^{\mathbf{T}} \hat{\mathbf{c}}_{kr}) \quad (70)$$

recalling (9) and (17), and then measuring its distance to the true sketch  $\mathbf{y}$ . The initialization index is then selected as

$$r_* = \arg \min_r \|\mathbf{y} - \hat{\mathbf{y}}_r\|, \quad (71)$$

and the centroids saved as  $\hat{\mathbf{C}} = \hat{\mathbf{C}}_{r_*}$ . In Section III, we used  $R = 2$  for all experiments.

### F. Hyperparameter Tuning

The likelihood model  $p_{y|\mathbf{z}}$  in (18) depends on the unknown hyperparameters  $\alpha$  and  $\tau$ . We propose to estimate these hyperparameters using a combination of *expectation maximization* (EM) and SHyGAMP, as suggested in [14] and detailed—for the simpler case of GAMP—in [21]. The idea is to run SHyGAMP using an estimate of  $\alpha$  and  $\tau$ , update  $\alpha$  and  $\tau$  from the SHyGAMP outputs, and repeat until convergence. For the first estimate, we suggest to use  $\alpha_k = \frac{1}{K}$  and  $\tau_k = 0 \quad \forall k$ .

Extrapolating [21, eq. (23)] to the SHyGAMP case, the EM update of  $(\alpha, \tau)$  takes the form

$$(\hat{\alpha}, \hat{\tau}) = \arg \max_{\alpha \geq 0, \alpha^{\mathbf{T}} \mathbf{1} = 1, \tau > 0} \sum_{m=1}^M \int_{\mathbb{R}^K} \mathcal{N}(z_m; \hat{z}_m, \mathbf{Q}_m^z) \times \ln p_{y|\mathbf{z}}(y_m | z_m; \alpha, \tau) dz_m, \quad (72)$$

where  $\hat{z}_m$  and  $\mathbf{Q}_m^z = \text{Diag}\{q_m^z\}$  are obtained by running SHyGAMP to convergence under  $(\alpha, \tau)$ . To proceed, we model the Dirac delta in (18) using a circular Gaussian pdf with vanishingly small variance  $\epsilon > 0$ , in which case

$$\begin{aligned} \ln p_{y|\mathbf{z}}(y_m | z_m; \alpha, \tau) \\ = -\frac{1}{\epsilon} \left| y_m - \sum_{k=1}^K \alpha_k \exp\left(jg_m z_{mk} - \frac{g_m^2 \tau_k}{2}\right) \right|^2 + \text{const}. \end{aligned} \quad (73)$$

Plugging (73) back into (72), we see that the constant and the  $1/\epsilon$ -scaling play no role in the optimization, and so we can discard them to obtain

$$\begin{aligned} (\hat{\alpha}, \hat{\tau}) = \arg \min_{\alpha \geq 0, \alpha^{\mathbf{T}} \mathbf{1} = 1, \tau > 0} \sum_{m=1}^M \int_{\mathbb{R}^K} \mathcal{N}(z_m; \hat{z}_m, \mathbf{Q}_m^z) \\ \times \left| y_m - \sum_{k=1}^K \alpha_k \exp\left(jg_m z_{mk} - \frac{g_m^2 \tau_k}{2}\right) \right|^2 dz_m. \end{aligned} \quad (74)$$

A closed-form solution to the optimization problem in (74) seems out of reach. Also, the optimization objective is convex in  $\alpha$  for fixed  $\tau$ , and convex in  $\tau$  for fixed  $\alpha$ , but not jointly convex in  $[\alpha^{\mathbf{T}}, \tau^{\mathbf{T}}]$ . Although the optimization problem (74) is difficult to solve, the solutions obtained by gradient projection (GP) [22] seem to work well in practice. Also, GP is made practical by closed-form gradient expressions. In particular, let

$$q_{mk} \triangleq \exp\left(-\frac{g_m^2 \tau_k}{2}\right) \quad (75)$$

$$\rho_{mk} \triangleq \exp\left(jg_m \hat{z}_{mk} - \frac{g_m^2 \tau_k}{2}\right), \quad (76)$$

and recall that  $v_{mk} = \exp(jg_m z_{mk})$  from (29) (although there the  $m$  subscript was suppressed). Then the  $m$ th term of the sum in the objective in (74) becomes

$$\begin{aligned} \int_{\mathbb{R}^K} \mathcal{N}(z_m; \hat{z}_m, \mathbf{Q}_m^z) \left| y_m - \sum_{k=1}^K \alpha_k q_{mk} v_{mk} \right|^2 dz_m \\ = |y_m|^2 - 2 \sum_{k=1}^K \alpha_k q_{mk} \text{Re}\{y_m^* \rho_{mk}\} \\ + \sum_{k=1}^K \alpha_k q_{mk} \rho_{mk}^* \sum_{l \neq k}^K \alpha_l q_{ml} \rho_{ml} + \sum_{k=1}^K \alpha_k^2 q_{mk}^2, \end{aligned} \quad (77)$$

where we used the fact that  $\int_{\mathbb{R}} \mathcal{N}(z_{mk}; \hat{z}_{mk}, q_{mk}^z) v_{mk} dz_{mk} = \rho_{mk}$ . After reapplying the sum over  $m$ , we get

$$\begin{aligned} \frac{\partial}{\partial \alpha_k} \sum_{m=1}^M \int_{\mathbb{R}^K} \mathcal{N}(z_m; \hat{z}_m, \mathbf{Q}_m^z) \left| y_m - \sum_{k=1}^K \alpha_k q_{mk} v_{mk} \right|^2 dz_m \\ = -2 \sum_{m=1}^M q_{mk} \gamma_{mk} \end{aligned} \quad (78)$$

---

**Algorithm 2** CL-AMP with hyperparameter tuning and multiple random initializations
 

---

**Require:** Measurements  $\mathbf{y} \in \mathbb{C}^M$ , gains  $\{g_m\}_{m=1}^M$ , number of initializations  $R \geq 1$ , initializations  $\{\widehat{\mathbf{C}}_{0,r}\}_{r=1}^R$ ,  $\mathbf{q}_0^{\mathbf{p}}$ ,  $\alpha_0$ ,  $\tau_0$ .

- 1:  $i = 0$
- 2: **repeat**
- 3:   **if**  $i = 0$  **then**
- 4:     **for**  $r = 1 : R$  **do**
- 5:       Run CL-AMP with fixed  $(\alpha_0, \tau_0)$  from initialization  $(\widehat{\mathbf{C}}_{0,r}, \mathbf{q}_0^{\mathbf{p}})$ , yielding output  $\widehat{\mathbf{C}}_{1,r}$ ,  $\widehat{\mathbf{Z}}_r$ , and  $\{\mathbf{q}_{mr}^{\mathbf{z}}\}_{m=1}^M$ .
- 6:     **end for**
- 7:     Compute  $\widehat{\mathbf{y}}_{mr} \triangleq \sum_{k=1}^K \alpha_{0k} \exp(-g_m^2 \tau_{0k}) \exp(jg_m \widehat{z}_{mkr}) \forall mr$
- 8:     Find  $r_* = \arg \min_r \|\mathbf{y} - \widehat{\mathbf{y}}_r\|$ .
- 9:     Set  $\widehat{\mathbf{C}}_1 = \widehat{\mathbf{C}}_{1,r_*}$ ,  $\widehat{\mathbf{Z}} = \widehat{\mathbf{Z}}_{r_*}$  and  $\{\mathbf{q}_m^{\mathbf{z}}\}_{m=1}^M = \{\mathbf{q}_{mr_*}^{\mathbf{z}}\}_{m=1}^M$ .
- 10:   **else**
- 11:     Run CL-AMP with fixed  $(\alpha_i, \tau_i)$  from initialization  $(\widehat{\mathbf{C}}_i, \mathbf{q}_0^{\mathbf{p}})$ , yielding output  $\widehat{\mathbf{C}}_{i+1}$ ,  $\widehat{\mathbf{Z}}$ , and  $\{\mathbf{q}_m^{\mathbf{z}}\}_{m=1}^M$ .
- 12:   **end if**
- 13:   Compute  $(\alpha_{i+1}, \tau_{i+1})$  via (74) using  $\widehat{\mathbf{Z}}$  and  $\{\mathbf{q}_m^{\mathbf{z}}\}_{m=1}^M$ .
- 14:    $i \leftarrow i + 1$ .
- 15: **until** convergence

---

$$\begin{aligned} & \frac{\partial}{\partial \tau_k} \sum_{m=1}^M \int_{\mathbb{R}^K} \mathcal{N}(\mathbf{z}_m; \widehat{\mathbf{z}}_m, \mathbf{Q}_m^{\mathbf{z}}) \left| y_m - \sum_{k=1}^K \alpha_k q_{mk} v_{mk} \right|^2 d\mathbf{z}_m \\ & = \alpha_k \sum_{m=1}^M g_m^2 q_{mk} \gamma_{mk} \end{aligned} \quad (79)$$

for

$$\gamma_{mk} \triangleq \operatorname{Re} \{ y_m^* \rho_{mk} \} - \alpha_k q_{mk} - \sum_{l \neq k} \alpha_l q_{ml} \operatorname{Re} \{ \rho_{mk}^* \rho_{ml} \}. \quad (80)$$

We found that complexity of hyperparameter tuning can be substantially reduced, without much loss in accuracy, by using only a subset of the terms in the sum in (74), as well as in the corresponding gradient expressions (78)-(79). For the experiments in Section III, we used a fixed random subset of  $\min(M, 20K)$  terms.

### G. Algorithm Summary

Algorithm 2 summarizes the CL-AMP algorithm with  $R$  random initializations and tuning of the hyperparameters  $(\alpha, \tau)$ . Note that the random initializations  $\{\widehat{\mathbf{C}}_{0,r}\}$  are used only for the first EM iteration, i.e.,  $i = 0$ . Subsequent EM iterations (i.e.,  $i \geq 1$ ) are initialized using the output  $\widehat{\mathbf{C}}_i$  of the previous EM iteration.

## III. NUMERICAL EXPERIMENTS

In this section, we present the results of several experiments used to test the performance of the CL-AMP, CL-OMPR, and k-means++ algorithms. For k-means++, we used the implementation provided by MATLAB and, for CL-OMPR, we downloaded the MATLAB implementation from [23]. CL-OMPR and CL-AMP used the same sketch  $\mathbf{y}$ , whose frequency vectors  $\mathbf{W}$  were drawn using the method described in Section II-A, with the scaling parameter  $\sigma^2$  set via (21). For CL-OMPR and CL-AMP, the reported runtimes include the time of computing the sketch, unless otherwise noted.

All experiments were run on a Dell PowerEdge C6320 two-socket server with Intel Xeon E5-2680 v4 processors (14 cores, 2.40GHz) and 128GB RAM.

### A. Experiments with Synthetic Data

1) *Performance vs. sketch length  $M$* : In the first experiment, we test each algorithm's ability to minimize SSE on a set of training data, i.e., to solve the problem (1). In addition, we test how well the recovered centroids work in minimum-distance classification.

The experiment was conducted as follows. Fixing the number of classes at  $K = 10$  and the data dimension at  $N = 100$ , ten Monte Carlo trials were performed. In each trial, the true centroids were randomly drawn<sup>3</sup> as  $\mathbf{c}_k \sim \mathcal{N}(\mathbf{0}_N, 1.5^2 K^{2/N} \mathbf{I}_N)$ . Then, using these centroids, a training dataset  $\{\mathbf{x}_t\}_{t=1}^T$  with  $T = 10^7$  samples was drawn from the GMM (5) with weights  $\alpha_k = 1/K$  and covariances  $\Phi_k = \mathbf{I}_N \forall k$ . Additionally, a test dataset  $\{\bar{\mathbf{x}}_t\}$  of  $10^6$  samples was independently generated.

For centroid recovery, k-means++ was invoked on the training dataset, and both CL-AMP and CL-OMPR were invoked after sketching the training data with  $M$  samples as in (2). Sketch lengths  $M/KN \in \{1, 2, 3, 5, 10, 20\}$  were investigated. CL-AMP used two random initializations, i.e.,  $R = 2$  as defined in Algorithm 2.

For each algorithm, the SSE of its estimated centroids  $\{\widehat{\mathbf{c}}_k\}_{k=1}^K$  was calculated using the training data  $\{\mathbf{x}_t\}_{t=1}^T$  via (1). Additionally, the performance of the estimated centroids in minimum-distance classification was evaluated as follows. First, labels  $\{j_k\}_{k=1}^K$  were assigned to the estimated centroids by solving the linear assignment problem [24] without replacement, given by

$$\arg \min_{\{j_1, \dots, j_K\} = \{1, \dots, K\}} \sum_{k=1}^K \|\mathbf{c}_k - \widehat{\mathbf{c}}_{j_k}\|_2^2. \quad (81)$$

Next, each test sample  $\bar{\mathbf{x}}_t$  was classified using minimum-distance classification, producing the estimated label

$$\widehat{k}_t = \arg \min_{k \in \{1, \dots, K\}} \|\bar{\mathbf{x}}_t - \widehat{\mathbf{c}}_{j_k}\|. \quad (82)$$

The classification error rate (CER) was then calculated as the proportion of estimated labels  $\widehat{k}_t$  that do not equal the true label  $k_t$  from which the test sample  $\bar{\mathbf{x}}_t$  was generated.<sup>4</sup>

Figures 1a, 1b, and 1c show the median SSE, CER, and runtime (including sketching), respectively, for CL-AMP and CL-OMPR versus  $M/KN$ . Also shown is the median SSE, CER, and runtime of k-means++, as a baseline, where k-means++ has no dependence on  $M$ . Because a low runtime is meaningless if the corresponding SSE is very high, the runtime was not shown for CL-AMP and CL-OMPR whenever its SSE

<sup>3</sup>This data-generation model was chosen to match that from [5], and is intended to have a relatively constant Bayes error rate w.r.t.  $N$  and  $K$ . For the chosen parameters, the Bayes error rate is extremely small:  $10^{-24}$ . Thus, when the centroids are accurately recovered, the classification error rate should be essentially zero.

<sup>4</sup>Note that the true label  $k_t$  was assigned when the test sample  $\bar{\mathbf{x}}_t$  was generated. The true label  $k_t$  does not necessarily indicate which of the true centroids  $\{\mathbf{c}_k\}$  is closest to  $\bar{\mathbf{x}}_t$ .



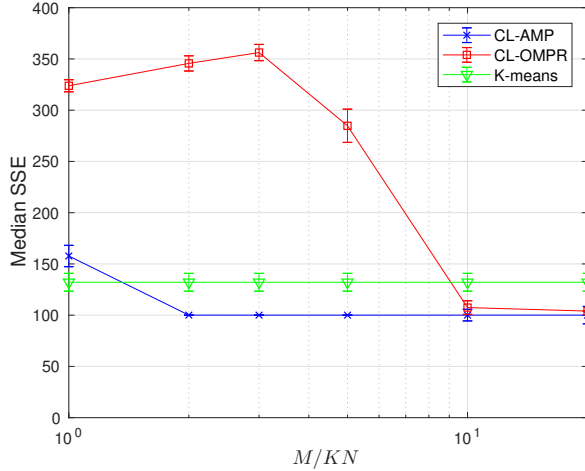
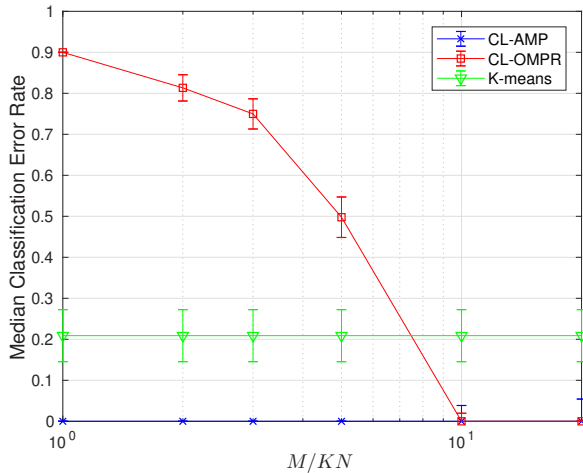
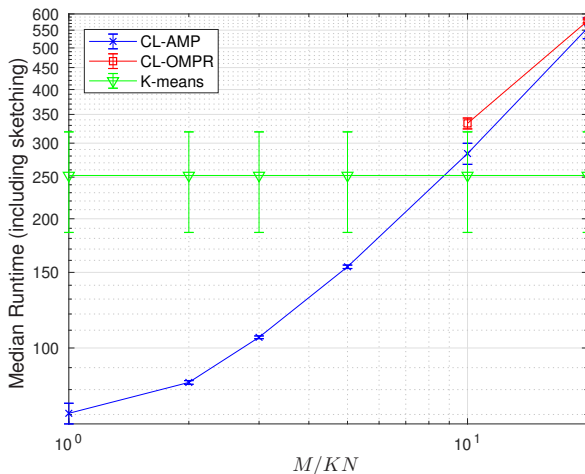
(a) SSE vs.  $M$ (b) Classification Error Rate vs.  $M$ (c) Runtime (including sketching) vs.  $M$ 

Fig. 1: Performance vs. sketch length  $M$  for  $K = 10$  clusters, dimension  $N = 100$ , and  $T = 10^7$  training samples.

was more than 1.5 times that of k-means++. The error bars show the standard deviation of the estimates.

Figure 1a shows that, among the methods tested, CL-AMP achieved the lowest SSE when  $M \geq 2KN$ . Also, CL-AMP supported the use of smaller sketch sizes  $M$  than CL-OMPR. In particular, CL-AMP required  $M \geq 2KN$  to yield a low SSE, while CL-OMPR required  $M \geq 10KN$ . This behavior mirrors the behavior of AMP and OMP in the classical compressive sensing context, where AMP usually requires fewer measurements to accurately recover signals of a given sparsity (see, e.g., [21, Figs. 8-10]). Also, with sufficiently large  $M$ , the SSE achieved by CL-AMP and CL-OMPR was lower than that achieved by k-means++.

Figure 1b shows that CL-AMP achieved a low CER with sketch size  $M \geq KN$ , while again CL-OMPR required  $M \geq 10KN$ . Also, with sufficiently large  $M$ , CL-AMP and CL-OMPR achieved near-zero CER, whereas k-means++ achieved an error rate of  $\approx 0.2$ .

Finally, Fig. 1c shows that, for  $M/KN \in \{10, 20\}$ , k-means++ ran slightly faster than CL-AMP, which ran slightly faster than CL-OMPR. However, for  $M/KN \in \{1, 2, 3, 5\}$ , CL-AMP ran significantly faster than k-means++. For  $M/KN \in \{1, 2, 3, 5\}$ , the runtime of CL-OMPR was not shown because it generated centroids of significantly worse SSE than those of k-means++.

2) *Performance vs. number of classes  $K$* : In a second experiment, we evaluated each algorithm's performance versus the number of classes  $K \in \{5, 10, 15, 20, 25, 30, 40, 50\}$  and sketch sizes  $M/KN \in \{2, 5, 10\}$  for fixed data dimension  $N = 50$ . The data was generated in exactly the same way as the previous experiment, and the same performance metrics were evaluated. Figures 2a, 2b, and 2c show the median SSE, CER, and runtime (including sketching) versus  $K$ , for CL-AMP, CL-OMPR, and k-means++.

Figure 2a shows that, as  $K$  increases, the SSE of k-means++ remained roughly constant, as expected based on the generation of the true centers  $c_k$ . For  $K \leq 20$ , CL-AMP yielded the best SSE for all tested values of  $M$ . For  $K > 20$ , CL-AMP yielded the best SSE with sketch sizes  $M \in \{5KN, 10KN\}$ , but performed poorly with  $M = 2KN$ . Meanwhile, CL-OMPR performed reasonably well with sketch size  $M = 10KN$ , but poorly with  $M \in \{2KN, 5KN\}$ .

Figure 2b shows similar trends. With sketch size  $M \in \{5KN, 10KN\}$ , CL-AMP had the lowest CER of any algorithm for all tested values of  $K$ . With sketch size  $M = 10KN$ , CL-OMPR gave CER better than k-means++ for all tested  $K$ , but with  $M \in \{2KN, 5KN\}$  CL-OMPR gave CER worse than k-means++ for all tested  $K$ .

Finally, Fig. 2c shows that CL-AMP ran faster than CL-OMPR at all tested  $K$  due to its ability to work with a smaller sketch size  $M$ . For large  $K$ , Fig. 2c suggests that the runtime of both CL-AMP and CL-OMPR grow as  $O(K^2)$ . The  $O(K^2)$  complexity scaling is expected for CL-AMP, since its complexity is  $O(MNK)$  and we set  $M = O(K)$ . But the  $O(K^2)$  complexity scaling is somewhat surprising for CL-OMPR, since its complexity is  $O(MNK^2)$  and we set  $M = 10NK$ . Also, Fig. 2c shows that CL-AMP ran faster than k-means++ for most values of  $K$ ; for the smallest tested value of  $K$  (i.e.,  $K = 5$ ), the median runtime of k-means++ was lower than CL-AMP (but the error-bar suggests that the

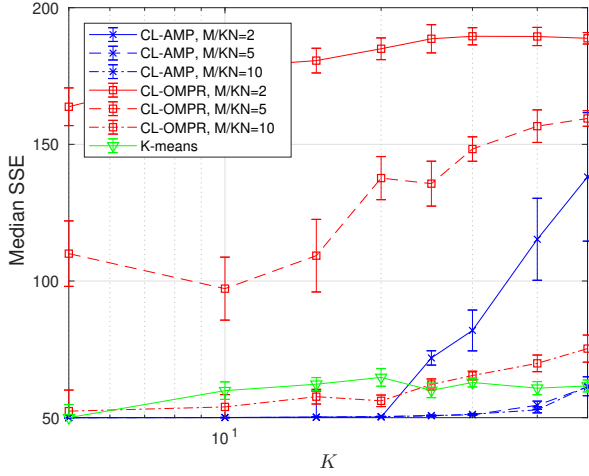
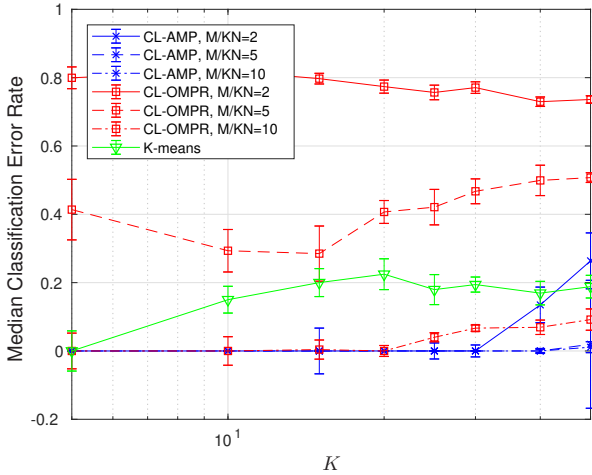
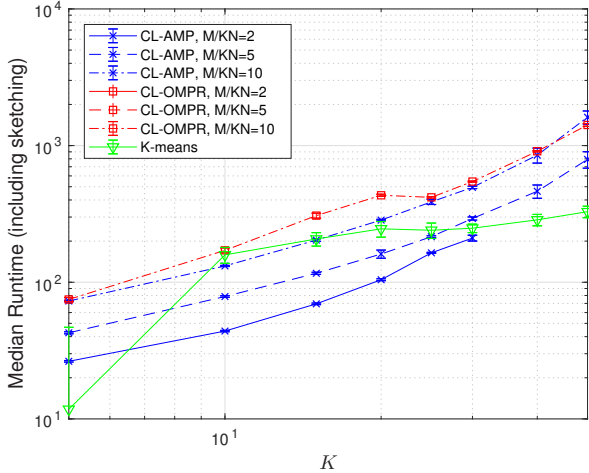
(a) SSE vs.  $K$ (b) Classification Error Rate vs.  $K$ (c) Runtime (including sketching) vs.  $K$ 

Fig. 2: Performance vs. number of clusters  $K$  for dimension  $N = 50$ , sketch size  $M \in \{2, 5, 10\} \times KN$ , and  $T = 10^7$  training samples.

runtime of k-means++ was highly variable at this  $K$ ). For the largest tested value of  $K$ , k-means++ was again faster than CL-AMP, because the runtime of k-means++ is expected to grow linearly with  $K$ , whereas that of CL-AMP is expected to grow quadratically with  $K$  when  $M/KN$  is fixed.

3) *Performance vs. dimension  $N$* : In a third experiment, we evaluated each algorithm's performance versus the dimension  $N$  (logarithmically spaced between 10 and 316) for  $K = 10$  classes and sketch size  $M \in \{2, 5, 10\} \times KN$ . The data was generated in exactly the same way as the previous two experiments, and the same performance metrics were evaluated. Figures 3a, 3b, and 3c show the median SSE/ $N$ , the CER, and the runtime (including sketching) versus  $N$ , for CL-AMP, CL-OMPR, and k-means++.

Figure 3a shows that, among all algorithms, CL-AMP achieved the lowest SSE for all tested values of  $N$  and  $M$ . Meanwhile, both CL-OMPR under sketch size  $M = 10KN$  and k-means++ achieved reasonably good SSE, but CL-OMPR under smaller sketches gave much higher SSE.

Figure 3b shows that, among all algorithms, CL-AMP achieved the lowest CER for all tested values of  $N$  and  $M$ . Meanwhile, CL-OMPR under sketch size  $M = 10KN$  gave similar CER to CL-AMP for most  $N$ , k-means++ gave significantly worse CER compared to CL-AMP for all  $N$ , and CL-OMPR under sketch size  $M = 5KN$  or  $2KN$  gave even worse CER for all  $N$ .

Finally, Fig. 3c shows that, among all algorithms, CL-AMP with sketch size  $M = 2KN$  ran the fastest for all tested values of  $N$ . Meanwhile, CL-OMPR with sketch size  $M = 10KN$  ran at a similar speed to CL-AMP with sketch size  $M = 10KN$ , for all  $N$ . The runtimes for CL-OMPR with smaller sketches are not shown because it achieved significantly worse SSE than k-means++. Figure 3c suggests that, if  $N$  is increased beyond 316, then eventually k-means++ will be faster than CL-AMP under fixed  $M/KN$ .

4) *Performance vs. training size  $T$* : In a final synthetic-data experiment, we evaluated each algorithm's performance versus the number of training samples  $T$  (logarithmically spaced between  $10^5$  and  $10^8$ ) for  $K = 10$  classes, dimension  $N = 50$ , and sketch size  $M \in \{2, 5, 10\}KN$ . The data was generated in exactly the same way as the previous three experiments, and the same performance metrics were evaluated.

Figures 4a and 4b show the median SSE and CER versus  $T$ , for CL-AMP, CL-OMPR, and k-means++. From these figures, we observe that the SSE and CER for each algorithm (and sketch length  $M$ ) were approximately invariant to  $T$ . CL-AMP (under any tested  $M$ ) yielded the lowest values of SSE and CER. Both CL-OMPR under sketch size  $M = 10KN$  and k-means++ gave reasonably good SSE and CER, but CL-OMPR under smaller sketches gave worse SSE and CER.

Figures 4c and 4d show the median runtime with and without sketching, respectively, for the algorithms under test. Figure 4c shows that, if sketching time is included in runtime, then all runtimes increased linearly with training size  $T$ . However, for large  $T$ , CL-AMP ran faster than k-means++ and CL-OMPR (while also achieving lower SSE and CER). Meanwhile, Fig. 4d shows that, if sketching time is not included in runtime, then the runtimes of both CL-AMP and

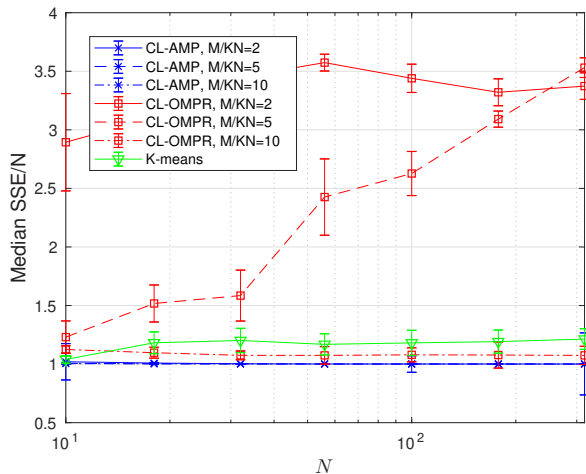
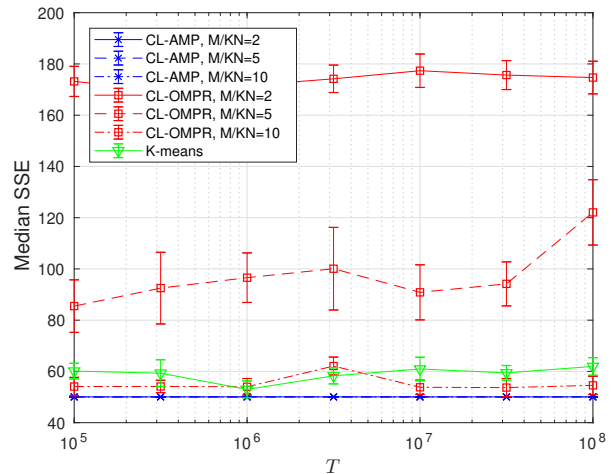
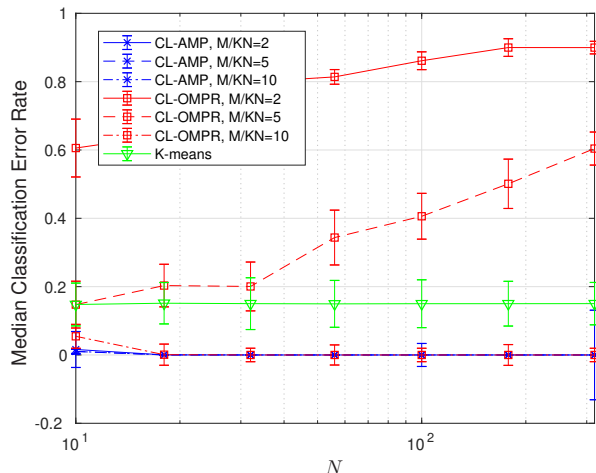
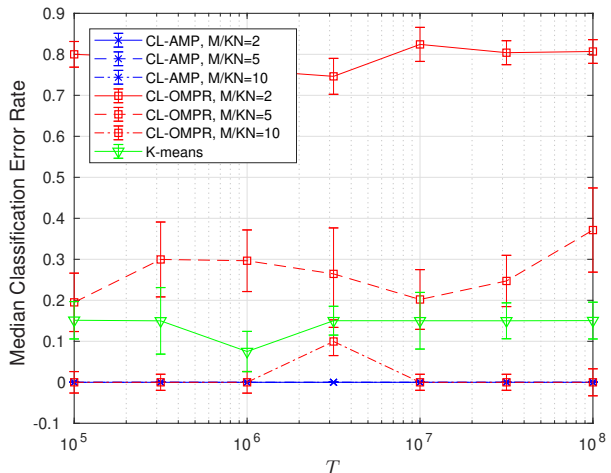
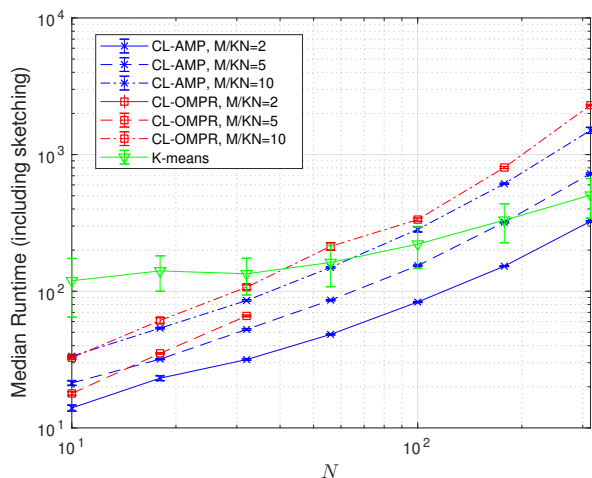
(a) SSE/ $N$  vs.  $N$ (a) SSE vs.  $T$ (b) Classification Error Rate vs.  $N$ (b) Classification Error Rate vs.  $T$ (c) Runtime (including sketching) vs.  $N$ 

Fig. 3: Performance vs. dimension  $N$  for  $K = 10$  classes,  $T = 10^7$  samples, and sketch size  $M \in \{2, 5, 10\} \times KN$ .

CL-OMPR were relatively invariant to  $T$ . Also, Figures 4c and 4d together show that, for  $T > 10^6$ , the sketching time was

Fig. 4: Performance vs. training size  $T$  for  $K = 10$  classes, dimension  $N = 50$ , and sketch size  $M \in \{2, 5, 10\} \times KN$ .

the dominant contributor to the overall runtime.

### B. Spectral Clustering of MNIST

Next we evaluated the algorithms on the task of spectral clustering [25] of the MNIST dataset. This task was previously investigated for CL-OMPR and k-means++ in [6], and we used the same data preprocessing steps: extract SIFT descriptors [26] of each image, compute the  $K$ -nearest-neighbors adjacency matrix (for  $K = 10$ ) using FLANN [27], and compute the 10 principal eigenvectors of the associated normalized Laplacian matrix (since we know  $K = 10$ ), yielding features of dimension  $N = 10$ . We applied this process to the original MNIST dataset, which includes  $T = 7 \times 10^4$  samples, as well as an augmented one with  $T = 3 \times 10^5$  samples constructed as described in [6].

The experiment was conducted as follows. In each of 10 trials, we randomly partitioned each sub-dataset into equally-sized training and testing portions. Then, we invoked CL-AMP, CL-OMPR, and k-means++ on the training portion of

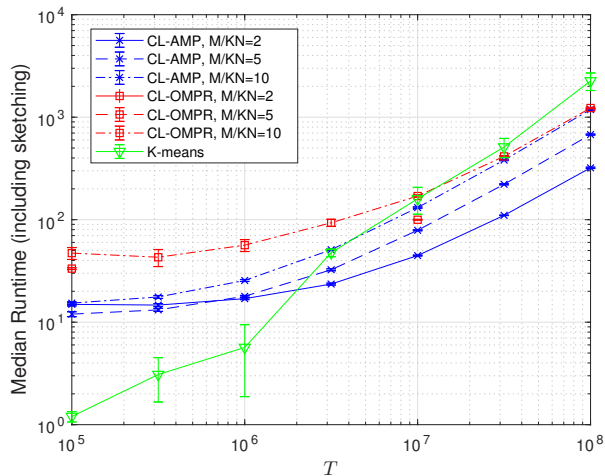
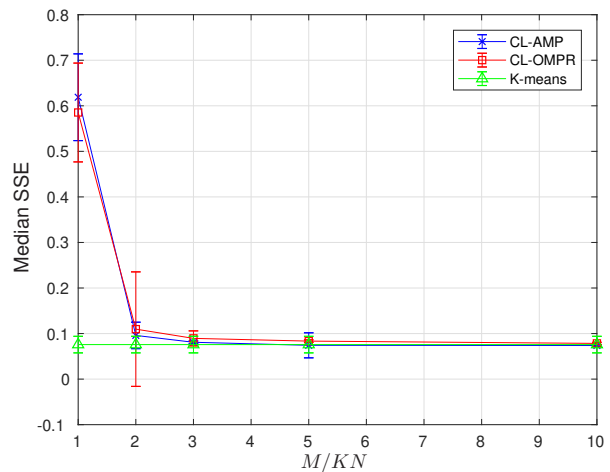
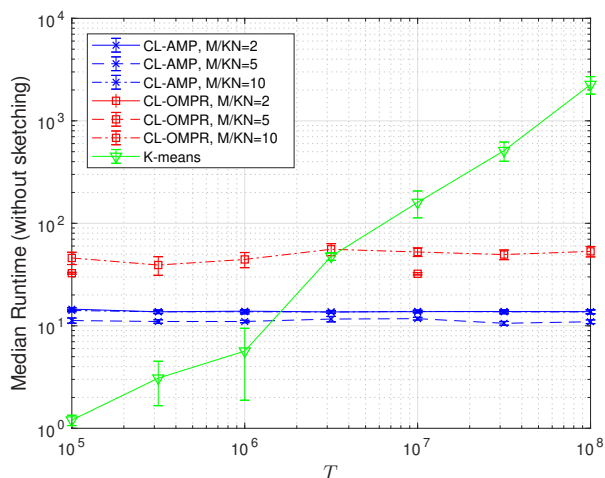
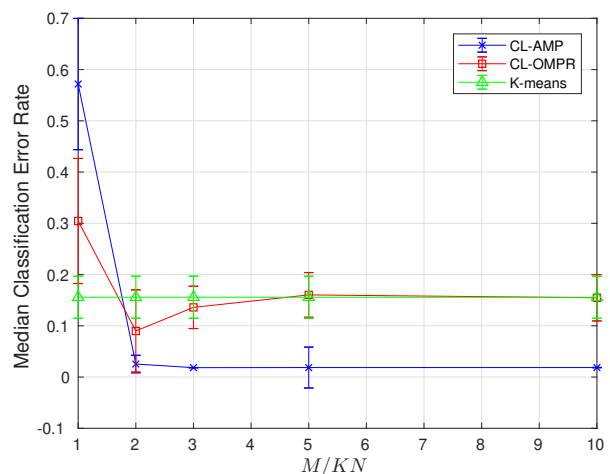
(c) Runtime (including sketching) vs.  $T$ (a) SSE vs.  $M$ (d) Runtime (without sketching) vs.  $T$ (b) Classification Error Rate vs.  $M$ 

Fig. 4: Performance vs. training size  $T$  for  $K = 10$  classes, dimension  $N = 50$ , and sketch size  $M \in \{2, 5, 10\} \times KN$ .

the dataset, using sketch sizes  $M \in \{1, 2, 3, 5, 10\} \times KN$  for CL-AMP and CL-OMPR. The algorithm parameters were the same as in Section III-A. Finally, the estimated centroids produced by each algorithm were evaluated using the same two metrics as in Section III-A: SSE on the training data, and classification error rate (CER) when the centroids were used for minimum-distance classification of the test data samples.

The median SSE, CER, and runtime, versus sketch length  $M$ , are shown for CL-AMP and CL-OMPR in Fig. 5 for the  $T = 7 \times 10^4$ -sample MNIST sub-dataset. As before, k-means++ is shown, as a baseline, although it does not use the sketch and thus its performance is invariant to  $M$ . From this figure, we observe that CL-AMP and CL-OMPR gave respectable results for sketch lengths  $M \geq 2KN$ , and SSE nearly identical to k-means++ for  $M \geq 5KN$ . For  $M \geq 2KN$ , however, CL-AMP yielded significantly lower CER than both CL-OMPR and k-means++, at the cost of a slower runtime. We attribute CL-AMP's slower runtime to its use of many iterations  $i$  in Algorithm 2 for hyperparameter tuning.

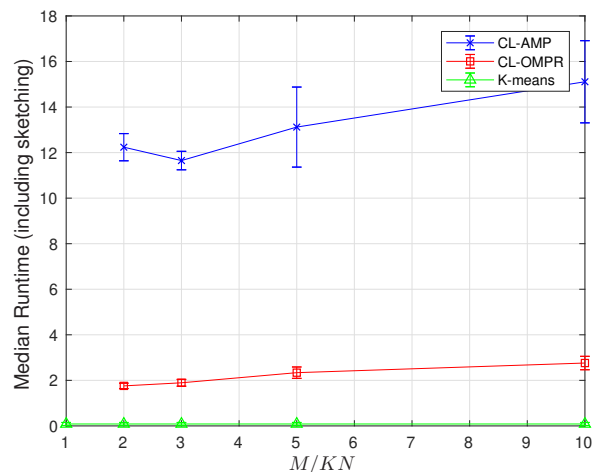
(c) Runtime (including sketching) vs.  $M$ 

Fig. 5: Performance vs.  $M$  for the  $T = 70\,000$ -sample spectral MNIST dataset, with  $K = 10$  clusters and dimension  $N = 10$ .

### C. Examination of Computational Complexity

To better understand the computational bottlenecks of the proposed approach, Fig. 6 shows—for several problem dimensions and data types—the runtime contributions of the “sketch,” i.e., equation (2); the “tuning” steps, i.e., line 13 of Algorithm 2; the “estimation” steps, i.e., lines 4-5 of Algorithm 1; and all other lines from Algorithm 1, which we refer to as the “linear” steps, since their complexity is dominated by the matrix multiplications in lines 3 and 9 of Algorithm 1.

Figure 6 suggests that CL-AMP’s estimation steps require the most computation, followed by its tuning steps, and finally its linear steps. These results motivate additional work to reduce the computational complexity of CL-AMP’s estimation steps. The cost of sketching itself depends on the number of training samples,  $T$ , and the degree to which the sketching operation is distributed over multiple processors. When  $T$  becomes large enough that the sketching time becomes computationally significant (as in Fig. 6b), the simplest remedy is to parallelize the sketch.

### D. Frequency Estimation

Our final experiment concerns multi-dimensional frequency estimation. Consider a sum-of-sinusoids signal of the form

$$y(\mathbf{t}) = \sum_{k=1}^K \alpha_k \exp(j\mathbf{t}^T \mathbf{c}_k), \quad (83)$$

where  $\mathbf{c}_k \in \mathbb{R}^N$  is the frequency of the  $k$ th sinusoid,  $\alpha_k > 0$  is the amplitude of the  $k$ th sinusoid, and  $\mathbf{t} \in \mathbb{R}^N$  denotes time. Given measurements of the signal  $y(\mathbf{t})$  at a collection of random times  $\mathbf{t} \in \{\mathbf{t}_m\}_{m=1}^M$ , i.e.,

$$y_m = y(\mathbf{t}_m) \text{ for } m = 1, \dots, M, \quad (84)$$

we seek to recover the frequencies  $\{\mathbf{c}_k\}_{k=1}^K$ . We are particularly interested in the case where the frequencies  $\{\mathbf{c}_k\}$  are closely spaced, i.e., the “super-resolution” problem.

Note that the model in (83) matches that in (13) with  $g_m \mathbf{a}_m = \mathbf{t}_m \forall m$  and  $\Phi_k = \mathbf{0} \forall k$ , so that we can apply CL-AMP to this frequency estimation problem. The model in (83) also matches (4) with  $\mathbf{w}_m = \mathbf{t}_m \forall m$ , and so we can also apply CL-OMPR. But we cannot apply k-means++.

For frequency pairs  $\{\mathbf{c}_1, \mathbf{c}_2\}$  with  $\|\mathbf{c}_1 - \mathbf{c}_2\|_2 \geq \epsilon$ , [28] claims that, with  $\{\mathbf{w}_m\}$  drawn randomly from an appropriate distribution, one can resolve the frequencies with  $M \geq O(\ln(1/\epsilon))$  measurements. However, choosing  $\mathbf{w}_m$  uniformly spaced on a grid would require  $M \geq O(1/\epsilon)$  measurements. Thus, for a final experiment, similar to those performed in [28], we did the following. For a particular  $N$  and  $K$  (where  $K$  is even for simplicity), we generated  $K/2$  pairs of frequencies  $\{\mathbf{c}_{2k-1}, \mathbf{c}_{2k}\}$ , where  $\|\mathbf{c}_{2k-1} - \mathbf{c}_{2k}\|_2 = \epsilon$  for  $k = 1, \dots, K/2$ . Then, for a particular realization of  $\{\mathbf{c}_k\}_{k=1}^K$  and  $\{\mathbf{w}_m\}_{m=1}^M$ , CL-AMP and CL-OMPR were invoked to estimate  $\{\hat{\mathbf{c}}_k\}_{k=1}^K$ . Recovery was declared successful if

$$\max_k \|\mathbf{c}_{j_k} - \hat{\mathbf{c}}_k\|_2 < \epsilon/2, \quad (85)$$

where  $\{j_k\}_{k=1}^K$  solves the linear assignment problem (81).

For our experiment, we tested  $K=4$  frequency components of dimension  $N=2$  and varied  $M$  from  $3KN$  to  $100KN$  while also varying  $\epsilon$  from  $10^{-1}$  to  $10^{-3}$ . For each combination, 10 trials were performed. The empirical probability of successful recovery is shown in Figures 7-8. In Fig. 7,  $\mathbf{a}_m$  were drawn uniformly on the unit sphere and  $g_m = |g'_m|$  with  $g'_m \sim \mathcal{N}(0, 4\epsilon^2 \log_{10}^2(\epsilon))$ , while in Fig. 8,  $\mathbf{w}_m$  are uniformly spaced on a  $N$ -dimensional hyper-grid with per-dimension spacing  $\pi/2$ . Superimposed on the figures are curves showing  $M/KN = 0.1/\epsilon$  and  $M/KN = \ln(1/\epsilon)$ . From the figures, we see that CL-AMP had a higher empirical probability of recovery than CL-OMPR, especially for small  $\epsilon$ . We also see that the empirical phase transition of CL-AMP is close to the  $\ln(1/\epsilon)$  curve with random frequency samples (i.e., Fig. 7a) and the  $0.1/\epsilon$  curve with uniform frequency samples (i.e., Fig. 8a).

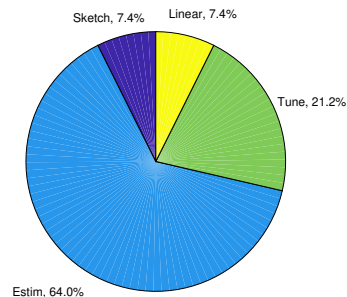
## IV. CONCLUSION

In sketched clustering, the original dataset is sketched down to a relatively short vector, from which the centroids are extracted. For the sketch proposed by [5,6], we proposed the “CL-AMP” centroid-extraction method. Our method assumes that the original data follows a GMM, and exploits the recently proposed simplified hybrid generalized approximate message passing (SHYGAMP) algorithm [14]. Numerical experiments suggest that CL-AMP exhibits better sample complexity (i.e., extracts accurate clusters with fewer compressed samples) than the state-of-the-art sketched-clustering algorithm, CL-OMPR, from [5,6]. In many cases, CL-AMP also exhibits better computational complexity than CL-OMPR. Furthermore, for datasets with many samples, CL-AMP exhibits lower computational complexity than the widely used k-means++ algorithm. As future work, it would be worthwhile to investigate ways to reduce the computational complexity of CL-AMP’s estimation steps, and to analyze the theoretical behavior of CL-AMP using a state-evolution approach. Finally, as new variations of the sketch (2) are proposed (e.g., [29,30]) it would be interesting to modify CL-AMP accordingly.

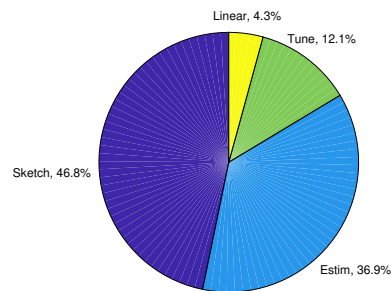
## REFERENCES

- [1] P. Drineas, A. Frieze, R. Kannan, S. Vempala, and V. Vinay, “Clustering large graphs via the singular value decomposition,” *Mach. Learn.*, vol. 56, no. 1-3, pp. 9–33, 2004.
- [2] H. Steinhaus, “Sur la division des corps matériels en parties,” *Bull. Acad. Polon. Sci.*, vol. 4, no. 12, pp. 801–804, 1956.
- [3] A. K. Jain, “Data clustering: 50 years beyond K-means,” *Pattern Recognition Letters*, vol. 31, no. 8, pp. 651–666, Jun. 2010.
- [4] D. Arthur and S. Vassilvitskii, “k-means++: The advantages of careful seeding,” in *Proc. Symp. Discrete Alg. (SODA)*, 2007, pp. 1027–1035.
- [5] N. Keriven, A. Bourrier, R. Gribonval, and P. Pérez, “Sketching for large-scale learning of mixture models,” *Inform. Inference*, vol. 7, no. 3, pp. 447–508, 2017.
- [6] N. Keriven, N. Tremblay, Y. Traonmilin, and R. Gribonval, “Compressive K-means,” in *Proc. IEEE Int. Conf. Acoust. Speech & Signal Process.*, 2017, pp. 6369–6373.
- [7] R. Gribonval, G. Blanchard, N. Keriven, and Y. Traonmilin, “Compressive statistical learning with random feature moments,” *arXiv:1706.07180*, 2017.
- [8] A. Feuerverger and R. A. Mureika, “The empirical characteristic function and its applications,” *Ann. Statist.*, vol. 5, no. 1, pp. 88–97, 1977.
- [9] Y. C. Pati, R. Rezaifar, and P. S. Krishnaprasad, “Orthogonal matching pursuit: Recursive function approximation with applications to wavelet decomposition,” in *Proc. Asilomar Conf. Signals Syst. Comput.*, Pacific Grove, CA, 1993, pp. 40–44.

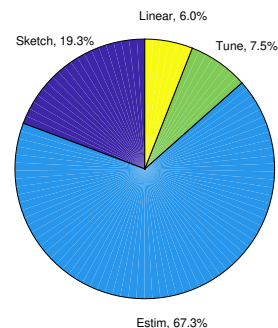
- [10] E. Byrne, R. Gribonval, and P. Schniter, “Sketched clustering via hybrid approximate message passing,” in *Proc. Asilomar Conf. Signals Syst. Comput.*, 2017, pp. 410–414.
- [11] P. McCullagh and J. A. Nelder, *Generalized Linear Models*, 2nd ed. London: Chapman & Hall/CRC, 1989.
- [12] A. Papoulis, *Probability, Random Variables, and Stochastic Processes*, 3rd ed. New York: McGraw-Hill, 1991.
- [13] M. Rudelson and R. Vershynin, “Hanson-Wright inequality and sub-Gaussian concentration,” *Electron. Commun. Probab.*, vol. 18, no. 82, pp. 1–9, 2013.
- [14] E. M. Byrne and P. Schniter, “Sparse multinomial logistic regression via approximate message passing,” *IEEE Trans. Signal Process.*, vol. 64, no. 21, pp. 5485–5498, 2016.
- [15] D. L. Donoho, A. Maleki, and A. Montanari, “Message passing algorithms for compressed sensing,” *Proc. Nat. Acad. Sci.*, vol. 106, no. 45, pp. 18 914–18 919, Nov. 2009.
- [16] S. Rangan, “Generalized approximate message passing for estimation with random linear mixing,” in *Proc. IEEE Int. Symp. Inform. Thy.*, Aug. 2011, pp. 2168–2172. (full version at *arXiv:1010.5141*).
- [17] M. Bayati and A. Montanari, “The dynamics of message passing on dense graphs, with applications to compressed sensing,” *IEEE Trans. Inform. Theory*, vol. 57, no. 2, pp. 764–785, Feb. 2011.
- [18] S. Rangan, A. K. Fletcher, V. K. Goyal, E. Byrne, and P. Schniter, “Hybrid approximate message passing,” *IEEE Trans. Signal Process.*, vol. 65, no. 17, pp. 4577–4592, 2017.
- [19] R. Gatto and S. R. Jammalamadaka, “The generalized von Mises distribution,” *Stat. Method.*, vol. 4, pp. 341–353, 2007.
- [20] C. M. Bishop, *Pattern Recognition and Machine Learning*. New York: Springer, 2007.
- [21] J. P. Vila and P. Schniter, “Expectation-maximization Gaussian-mixture approximate message passing,” *IEEE Trans. Signal Process.*, vol. 61, no. 19, pp. 4658–4672, Oct. 2013.
- [22] D. Bertsekas, *Nonlinear Programming*, 2nd ed. Athena Scientific, 1999.
- [23] N. Keriven, N. Tremblay, and R. Gribonval, “SketchMLbox : a Matlab toolbox for large-scale learning of mixture models,” 2016.
- [24] H. W. Kuhn, “The Hungarian Method for the Assignment Problem,” *Naval Research Logistics Quarterly*, pp. 83–97, 1955.
- [25] A. Y. Ng, M. I. Jordan, and Y. Weiss, “On spectral clustering: Analysis and an algorithm,” in *Proc. Neural Inform. Process. Syst. Conf.*, 2001, pp. 849–856.
- [26] A. Vedaldi and B. Fulkerson, “VLFeat: An open and portable library of computer vision algorithms,” in *Proc. ACM Intl. Conf. Multimedia*, 2010, pp. 1469–1472.
- [27] M. Muja and D. G. Lowe, “Fast approximate nearest neighbors with automatic algorithm configuration,” in *Proc. Intl. Conf. Comp. Vision Thy. Appl. (VISAPP)*, 2009, pp. 331–340.
- [28] Y. Traonmilin, N. Keriven, R. Gribonval, and G. Blanchard, “Spikes super-resolution with random Fourier sampling,” in *Proc. Workshop Signal Process. Adapt. Sparse Struct. Repr. (SPARS)*, 2017, pp. 1–2.
- [29] R. G. A. Chatalic and N. Keriven, “Large-scale high-dimensional clustering with fast sketching,” in *Proc. IEEE Int. Conf. Acoust. Speech & Signal Process.*, 2018, pp. 4714–4718.
- [30] V. Schellekens and L. Jacques, “Quantized compressive k-means,” *IEEE Signal Process. Lett.*, vol. 25, no. 8, pp. 1211–1215, 2018.



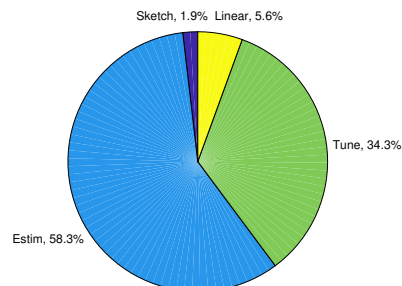
(a) Synthetic:  $N = 100$ ,  $K = 10$ ,  $M = 2KN$ , and  $T = 10^5$



(b) Synthetic:  $N = 100$ ,  $K = 10$ ,  $M = 2KN$ , and  $T = 10^6$

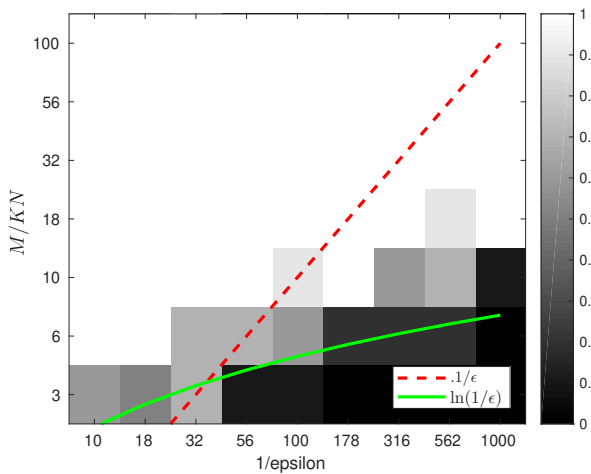


(c) Synthetic:  $N = 100$ ,  $K = 10$ ,  $M = 10KN$ , and  $T = 10^5$

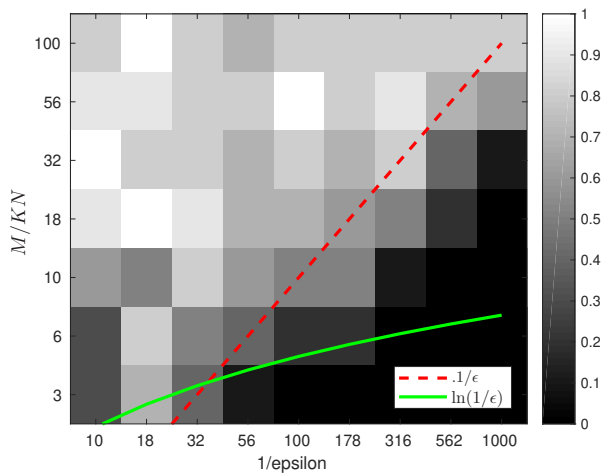


(d) MNIST:  $N = 10$ ,  $K = 10$ ,  $M = 5KN$ , and  $T = 3.5 \times 10^4$

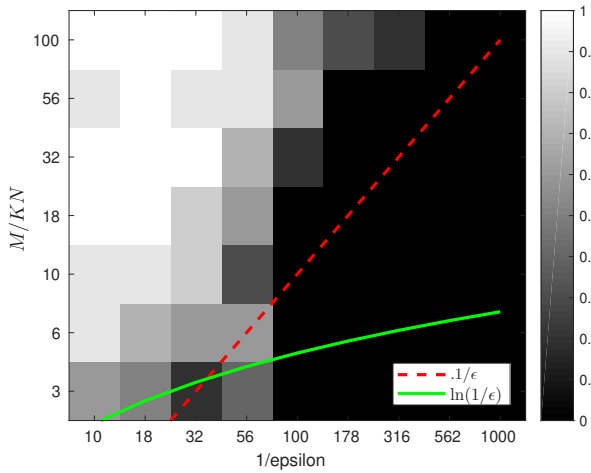
Fig. 6: Proportion of total runtime of the different sections of the CL-AMP algorithm applied to different datasets.



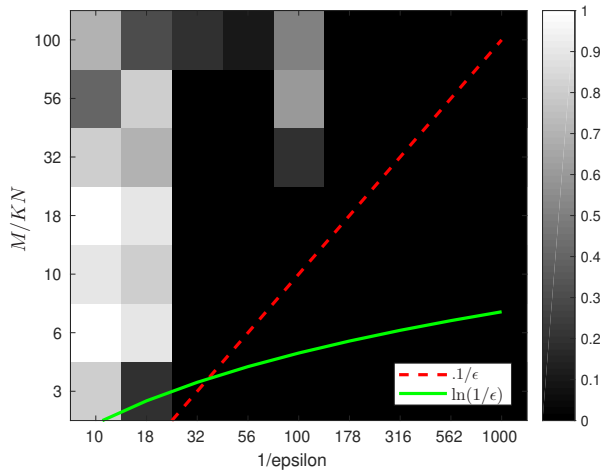
(a) CL-AMP



(a) CL-AMP



(b) CL-OMPR



(b) CL-OMPR

Fig. 7: Frequency estimation for  $K = 4$  and  $N = 2$  with random time samples.

Fig. 8: Frequency estimation for  $K = 4$  and  $N = 2$  with uniformly spaced time samples.