

Sketched Clustering via Hybrid Approximate Message Passing

Evan Byrne, Rémi Gribonval, Philip Schniter

► **To cite this version:**

Evan Byrne, Rémi Gribonval, Philip Schniter. Sketched Clustering via Hybrid Approximate Message Passing. Asilomar Conference on Signals, Systems, and Computers, Oct 2017, Pacific Grove, California, United States. hal-01650160

HAL Id: hal-01650160

<https://hal.inria.fr/hal-01650160>

Submitted on 28 Nov 2017

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,^{*} Rémi Gribonval,[†] and Philip Schniter,^{*}

^{*}Dept. of ECE, The Ohio State Univ., Columbus, OH, 43210, USA. (byrne.133@osu.edu, schniter.1@osu.edu)

[†]Univ Rennes, Inria, CNRS, IRISA, France. (remi.gribonval@inria.fr)

Abstract—In sketched clustering, the dataset is first sketched down to a vector of modest size, from which the cluster centers are subsequently extracted. The goal is to perform clustering more efficiently than with methods that operate on the full training data, such as k-means++. For the sketching methodology recently proposed by Keriven, Gribonval, et al., which can be interpreted as a random sampling of the empirical characteristic function, we propose a cluster recovery algorithm based on simplified hybrid generalized approximate message passing (SHyGAMP). Numerical experiments suggest that our approach is more efficient than the state-of-the-art sketched clustering algorithms (in both computational and sample complexity) and more efficient than k-means++ in certain regimes.

I. INTRODUCTION

Given a dataset $\mathbf{X} \triangleq [\mathbf{x}_1, \dots, \mathbf{x}_T] \in \mathbb{R}^{N \times T}$ comprising T feature vectors 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)

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

Finding the optimal \mathbf{C} is an NP-hard problem [1]. Thus, many heuristic approaches have been proposed, with one of the most popular being the *k-means algorithm* [2], [3]. Because k-means can get trapped in bad local minima, many robust variants have been proposed. One of the best known is *k-means++* [4], which uses a careful random initialization procedure to yield solutions with SSE that are on average $\leq 8(\ln K + 2)$ times the minimal SSE. But even with k-means++, many random re-initializations may be required to find a near-optimal clustering. For each initialization, the computational complexity of k-means++ scales as $O(TKNI)$, with I the number of iterations, which can be prohibitive for large T .

A. Sketched Clustering

In *sketched clustering* [5], [6], the dataset is first sketched down to a vector \mathbf{y} with $M \ll TN$ components, from which the cluster centers are subsequently extracted. If the sketch can be performed efficiently, then—since the cluster-extraction complexity will be independent of T —there is a chance that sketched clustering will be more efficient than direct clustering methods like k-means++ when T is large.

^{*}E. Byrne and P. Schniter acknowledge support from NSF grant 1716388 and MIT Lincoln Labs.

The approach proposed in [5], [6] uses a sketch $\mathbf{y} \triangleq [y_1, \dots, y_M]^T$ of the form

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

with randomly 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, i.e.,

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

under the empirical distribution $p_{\mathbf{x}}(\mathbf{x}) = \frac{1}{T} \sum_{t=1}^T \delta(\mathbf{x} - \mathbf{x}_t)$. Note that sketching \mathbf{X} as \mathbf{y} via (2) costs $O(TMN)$ operations, but it is easily parallelized.

For the recovery of cluster centers from \mathbf{y} , the state-of-the-art algorithm is *compressed learning orthogonal matching pursuit with replacement* (CL-OMPR) [5], [6]. It aims to solve

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

using a greedy heuristic inspired by the *orthogonal matching pursuit* algorithm [7] popular in compressed sensing. With $M \approx 10KN$, CL-OMPR often attains SSEs similar to those attained with k-means++, despite the lack of a direct link between the problem formulations (1) and (4). CL-OMPR's computational complexity is $O(MNK^2)$, however, which can be impractical when K is large. Thus, we seek a sketched clustering scheme whose complexity grows linearly in K .

B. Contributions

For the recovery of the cluster centers from a sketch \mathbf{y} of the form given in (2), we propose *compressive learning via approximate message passing* (CL-AMP). As we will see, CL-AMP has a computational complexity of $O(MNK)$ and performs favorably to CL-OMPR in terms of both runtime and sample complexity M . Furthermore, we find that CL-AMP performs favorably to k-means++ in certain operating regimes. CL-AMP can be understood as an application of the *simplified hybrid generalized approximate message passing* (SHyGAMP) framework [8] to sketched clustering. Further details will be provided in the sequel.

II. COMPRESSIVE LEARNING VIA AMP

A. High-Dimensional Inference Framework

CL-AMP formulates cluster recovery as a high-dimensional inference problem rather than as an optimization problem like

(4). In particular, it assumes a Gaussian mixture model (GMM)

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

where the GMM means are the cluster centers \mathbf{c}_k and the GMM weights α_k and covariances \mathbf{R}_k are unknown.

Defining $g_m \triangleq \|\mathbf{w}_m\|$ and $\tilde{\mathbf{w}}_m \triangleq \mathbf{w}_m/g_m$, so $\|\tilde{\mathbf{w}}_m\| = 1$,

$$y_m = \frac{1}{T} \sum_{t=1}^T \exp(\mathbf{j} \mathbf{w}_m^\top \mathbf{x}_t) \approx \mathbb{E}\{\exp(\mathbf{j} \mathbf{w}_m^\top \mathbf{x}_t)\} \quad (6)$$

$$= \sum_{k=1}^K \alpha_k \exp\left(\mathbf{j} g_m \underbrace{\tilde{\mathbf{w}}_m^\top \mathbf{c}_k}_{\triangleq z_{mk}} - g_m^2 \underbrace{\tilde{\mathbf{w}}_m^\top \mathbf{R}_k \tilde{\mathbf{w}}_m / 2}_{\triangleq \tau_{mk}}\right), \quad (7)$$

where (6) holds under large T and (7) follows from the facts that $\mathbf{w}_m^\top \mathbf{x}_t \sim \sum_k \alpha_k \mathcal{N}(g_m z_{mk}, g_m^2 \tau_{mk})$ and that $\mathbb{E}\{e^{\mathbf{j}x}\} = e^{\mathbf{j}z - \tau/2}$ when $x \sim \mathcal{N}(z, \tau)$. For $\tilde{\mathbf{w}}_m$ uniform on the sphere,

$$\tau_{mk} \xrightarrow{P} \mathbb{E}\{\tau_{mk}\} = \text{tr}(\mathbf{R}_k)/N \triangleq \tau_k \quad (8)$$

as $N \rightarrow \infty$, as long as the peak-to-average eigenvalue ratio of \mathbf{R}_k remains bounded [9]. Thus, for $\mathbf{z}_m \triangleq [z_{m1}, \dots, z_{mK}]^\top$,

$$p_{y|\mathbf{z}}(y_m | \mathbf{z}_m; \boldsymbol{\alpha}, \boldsymbol{\tau}) = \delta\left(y_m - \sum_{k=1}^K \alpha_k \exp\left(\mathbf{j} g_m z_{mk} - \frac{g_m^2 \tau_k}{2}\right)\right) \quad (9)$$

with hyperparameters $\boldsymbol{\tau} \triangleq [\tau_1, \dots, \tau_K]^\top$, $\boldsymbol{\alpha} \triangleq [\alpha_1, \dots, \alpha_K]^\top$.

The cluster coordinates $\{c_{nk}\}$ are modeled as i.i.d. $p_{\mathbf{c}}(c; \nu)$, where nominally $p_{\mathbf{c}}(c; \nu) = \mathcal{N}(c; 0, \nu)$ with large ν . Our main objective is then to compute the conditional mean

$$\hat{\mathbf{C}} = \mathbb{E}\{\mathbf{C} | \mathbf{y}\}, \quad (10)$$

where the expectation is taken over

$$p(\mathbf{C} | \mathbf{y}) \propto \prod_{m=1}^M p_{y|\mathbf{z}}(y_m | \mathbf{w}_m^\top \mathbf{C}; \boldsymbol{\alpha}, \boldsymbol{\tau}) \prod_{k=1}^K \prod_{n=1}^N p_{\mathbf{c}}(c_{nk}; \nu), \quad (11)$$

while simultaneously learning the hyperparameters $\boldsymbol{\alpha}, \boldsymbol{\tau}, \nu$. Here and in the sequel, we format random variables in sans-serif font for clarity.

B. Approximate Message Passing

Exact computation of $\hat{\mathbf{C}}$ in (10) is impractical due to the form of $p_{y|\mathbf{z}}$. One might consider approximate inference via the sum-product algorithm (SPA), but even the SPA is intractable due to the form of $p_{y|\mathbf{z}}$. Given the presence of a large random matrix \mathbf{W} in the problem formulation, we instead proposed to tackle approximate inference using *approximate message passing* (AMP) [10]. In particular, we apply the *simplified hybrid generalized AMP* (SHyGAMP) methodology from [8], while simultaneously estimating $\boldsymbol{\alpha}, \boldsymbol{\tau}, \nu$ through expectation maximization (EM). Some background on AMP methods will now be provided to justify our approach.

The original AMP algorithm of Donoho, Maleki, and Montanari [10] was designed to estimate i.i.d. \mathbf{c} under the standard linear model (i.e., $\mathbf{y} = \mathbf{W}\mathbf{c} + \mathbf{n}$ with known $\mathbf{W} \in \mathbb{R}^{M \times N}$

and additive white Gaussian noise \mathbf{n}). The generalized AMP (GAMP) algorithm of Rangan [11] extended AMP to the generalized linear model (i.e., $\mathbf{y} \sim p(\mathbf{y} | \mathbf{z})$ for $\mathbf{z} = \mathbf{W}\mathbf{c}$ 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{W} , while maintaining a computational complexity of only $O(MN)$. Furthermore, both can be rigorously analyzed via the state-evolution framework, which shows that they are Bayes-optimal in certain regimes [12].

A limitation of AMP [10] and GAMP [11] is that they cover only problems with i.i.d. estimand \mathbf{c} and separable likelihood $p(\mathbf{y} | \mathbf{z}) = \prod_{m=1}^M p(y_m | z_m)$. Thus, Hybrid GAMP (HyGAMP) [13] was developed to tackle problems with a structured prior and/or likelihood. HyGAMP could be applied to (10)-(11), but it requires computing and inverting $O(N+M)$ covariance matrices of dimension K at each iteration. The SHyGAMP algorithm [8] is a simplification of HyGAMP that uses diagonal covariance matrices to drastically reduce complexity. As described in [8], SHyGAMP can be readily combined with the EM algorithm for hyperparameter learning.

C. SHyGAMP

The SHyGAMP algorithm is summarized in Algorithm 1 using the language of Section II-A, assuming $\widehat{\mathbf{W}}$ has unit-norm rows. There, with some abuse of notation, we use \mathbf{c}_n^\top to denote the n th row of \mathbf{C} (where previously we used \mathbf{c}_k to denote the k th column of \mathbf{C}). We also use $\widehat{\mathbf{P}} \triangleq [\widehat{p}_1, \dots, \widehat{p}_M]^\top$, $\widehat{\mathbf{Z}} \triangleq [\widehat{z}_1, \dots, \widehat{z}_M]^\top$, $\widehat{\mathbf{R}} \triangleq [\widehat{r}_1, \dots, \widehat{r}_N]^\top$, \odot for componentwise division, and \odot for componentwise multiplication.

At each iteration, lines 10-11 of Algorithm 1 compute an approximation of the posterior mean and variance of $\{\mathbf{c}_{nk}\}$ using the ‘‘pseudo-measurements’’ $\widehat{\mathbf{r}}_n = \mathbf{c}_n + \mathbf{v}_n$, where \mathbf{v}_n is treated as a typical realization of $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}^r)$. Thus, the approximate posterior pdf used in lines 10-11 is

$$p_{\mathbf{c}|\mathbf{r}}(\mathbf{c}_n | \widehat{\mathbf{r}}_n; \mathbf{Q}^r) = \frac{p_{\mathbf{c}}(\mathbf{c}_n) \mathcal{N}(\mathbf{c}_n; \widehat{\mathbf{r}}_n, \mathbf{Q}^r)}{\int p_{\mathbf{c}}(\mathbf{c}'_n) \mathcal{N}(\mathbf{c}'_n; \widehat{\mathbf{r}}_n, \mathbf{Q}^r) d\mathbf{c}'_n}. \quad (12)$$

Similarly, lines 4-5 approximate the posterior mean and covariance of $\mathbf{z}_m \triangleq \mathbf{C}^\top \mathbf{w}_m$, which uses the pseudo-prior $\mathbf{z}_m \sim \mathcal{N}(\widehat{\mathbf{p}}_m, \mathbf{Q}^p)$ and hence the approximate posterior pdf

$$p_{\mathbf{z}|\mathbf{y}, \mathbf{p}}(\mathbf{z}_m | y_m, \widehat{\mathbf{p}}_m; \mathbf{Q}^p) = \frac{p_{y|\mathbf{z}}(y_m | \mathbf{z}_m) \mathcal{N}(\mathbf{z}_m; \widehat{\mathbf{p}}_m, \mathbf{Q}^p)}{\int p_{y|\mathbf{z}}(y_m | \mathbf{z}'_m) \mathcal{N}(\mathbf{z}'_m; \widehat{\mathbf{p}}_m, \mathbf{Q}^p) d\mathbf{z}'_m}. \quad (13)$$

Essentially, the SHyGAMP algorithm breaks an NK -dimensional inference problem into $M+N$ K -dimensional inference problems involving an independent-Gaussian pseudo-prior or pseudo-likelihood, evaluated iteratively. The resulting computational complexity is $O(MNK)$.

D. From SHyGAMP to CL-AMP

The SHyGAMP algorithm can be applied to many different problems via appropriate choice of $p_{y|\mathbf{z}}$ and $p_{\mathbf{c}}$. To apply SHyGAMP to sketched-clustering, we choose $p_{y|\mathbf{z}}$ and $p_{\mathbf{c}}$ as described in Section II-A. The principal remaining challenge is to evaluate lines 4-5 of Algorithm 1.

Algorithm 1 SHyGAMP

Require: Measurements \mathbf{y} , matrix $\widehat{\mathbf{W}}$ with $\|\widehat{\mathbf{W}}\|_F^2 = M$, pdfs $p_{\mathbf{c}|\mathbf{r}}$ and $p_{\mathbf{z}|\mathbf{y},\mathbf{p}}$ from (12)-(13), initializations $\widehat{\mathbf{C}} = \mathbb{E}\{\mathbf{C}\}$, $\mathbf{q}_n^{\mathbf{c}} = \text{diag}(\text{cov}\{\mathbf{c}_n\})$

Ensure: $\widehat{\mathbf{S}} \leftarrow \mathbf{0}$.

```

1: repeat
2:    $\mathbf{q}^{\mathbf{p}} \leftarrow \frac{1}{N} \sum_{n=1}^N \mathbf{q}_n^{\mathbf{c}}$ 
3:    $\widehat{\mathbf{P}} \leftarrow \widehat{\mathbf{W}} \widehat{\mathbf{C}} - \widehat{\mathbf{S}} \text{Diag}(\mathbf{q}^{\mathbf{p}})$ 
4:    $\mathbf{q}_m^{\mathbf{z}} \leftarrow \text{diag}(\text{cov}\{\mathbf{z}_m | y_m, \mathbf{p}_m = \widehat{\mathbf{p}}_m; \text{Diag}(\mathbf{q}^{\mathbf{p}})\}) \quad \forall m$ 
5:    $\widehat{\mathbf{z}}_m \leftarrow \mathbb{E}\{\mathbf{z}_m | y_m, \mathbf{p}_m = \widehat{\mathbf{p}}_m; \text{Diag}(\mathbf{q}^{\mathbf{p}})\} \quad \forall m$ 
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:    $\widehat{\mathbf{R}} \leftarrow \widehat{\mathbf{C}} + \widehat{\mathbf{W}} \widehat{\mathbf{S}}^T \text{Diag}(\mathbf{q}^{\mathbf{r}})$ 
10:   $\mathbf{q}_n^{\mathbf{c}} \leftarrow \text{diag}(\text{cov}\{\mathbf{c}_n | \mathbf{r}_n = \widehat{\mathbf{r}}_n; \text{Diag}(\mathbf{q}^{\mathbf{r}})\}) \quad \forall n$ 
11:   $\widehat{\mathbf{c}}_n \leftarrow \mathbb{E}\{\mathbf{c}_n | \mathbf{r}_n = \widehat{\mathbf{r}}_n; \text{Diag}(\mathbf{q}^{\mathbf{r}})\} \quad \forall n$ 
12: until Terminated

```

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

$$\widehat{z}_{mk} = C_m^{-1} \int_{\mathbb{R}^K} z_{mk} p_{\mathbf{y}|\mathbf{z}}(y_m | \mathbf{z}_m) \mathcal{N}(\mathbf{z}_m; \widehat{\mathbf{p}}_m, \mathbf{Q}^{\mathbf{p}}) d\mathbf{z}_m \quad (14)$$

$$q_{mk}^{\mathbf{z}} = \frac{\int_{\mathbb{R}^K} (z_{mk} - \widehat{z}_{mk})^2 p_{\mathbf{y}|\mathbf{z}}(y_m | \mathbf{z}_m) \mathcal{N}(\mathbf{z}_m; \widehat{\mathbf{p}}_m, \mathbf{Q}^{\mathbf{p}}) d\mathbf{z}_m}{C_m}, \quad (15)$$

where $C_m = \int_{\mathbb{R}^K} p_{\mathbf{y}|\mathbf{z}}(y_m | \mathbf{z}_m) \mathcal{N}(\mathbf{z}_m; \widehat{\mathbf{p}}_m, \mathbf{Q}^{\mathbf{p}}) d\mathbf{z}_m$. We propose approximations of \widehat{z}_{mk} and $q_{mk}^{\mathbf{z}}$ that are summarized below; a full derivation has been omitted due to space limitations. For the remainder of this section, we omit the subscripts m and $\mathbf{y}|\mathbf{z}$ to simplify the notation.

The main idea behind our approximation of \widehat{z}_{mk} and $q_{mk}^{\mathbf{z}}$ is to define $\theta_k \triangleq g z_k$ and then apply the Gaussian approximation (whose accuracy grows with K)

$$p\left(\begin{bmatrix} \text{Re}\{y\} \\ \text{Im}\{y\} \end{bmatrix} \middle| \theta_k\right) \approx \mathcal{N}\left(\begin{bmatrix} \text{Re}\{y\} \\ \text{Im}\{y\} \end{bmatrix}; \beta_k \begin{bmatrix} \cos(\theta_k) \\ \sin(\theta_k) \end{bmatrix} + \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\right) \quad (16)$$

to (9), where

$$\boldsymbol{\mu}_k = \sum_{l \neq k} \alpha_l e^{-g^2(\tau_k + [\mathbf{Q}^{\mathbf{p}}]_{kk})/2} \begin{bmatrix} \cos(g\widehat{p}_l) \\ \sin(g\widehat{p}_l) \end{bmatrix} \quad (17)$$

$$\boldsymbol{\Sigma}_k = \frac{1}{2} \sum_{l \neq k} \beta_l^2 (1 - e^{-g^2[\mathbf{Q}^{\mathbf{p}}]_{ll}}) \times \left(\mathbf{I} - e^{-g^2[\mathbf{Q}^{\mathbf{p}}]_{ll}} \begin{bmatrix} \cos(2g\widehat{p}_l) & \sin(2g\widehat{p}_l) \\ \sin(2g\widehat{p}_l) & -\cos(2g\widehat{p}_l) \end{bmatrix} \right) \quad (18)$$

$$\beta_k = \alpha_k \exp(-g^2 \tau_k / 2). \quad (19)$$

Rewriting (16) as

$$p\left(\beta_k^{-1} \begin{bmatrix} \text{Re}\{y\} \\ \text{Im}\{y\} \end{bmatrix} \middle| \theta_k\right) \approx \mathcal{N}\left(\begin{bmatrix} \cos(\theta_k) \\ \sin(\theta_k) \end{bmatrix}; \beta_k^{-1} \begin{bmatrix} \text{Re}\{y\} \\ \text{Im}\{y\} \end{bmatrix} - \beta_k^{-1} \boldsymbol{\mu}_k, \beta_k^{-2} \boldsymbol{\Sigma}_k\right), \quad (20)$$

the right side of (20) can be recognized as being proportional to the *generalized von Mises* (GvM) density [14] over $\theta_k \in [0, 2\pi)$. Under this GvM approximation, we have [14] that

$$p(y|\theta_k) \propto \exp(\kappa_k \cos(\theta_k - \zeta_k) + \bar{\kappa}_k \cos[2(\theta_k - \bar{\zeta}_k)]) \quad (21)$$

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

$$\kappa_k \cos(\zeta_k) = -\frac{1}{1 - \rho_k^2} \left(\frac{\rho_k \bar{\nu}_k}{\sigma_k \bar{\sigma}_k} - \frac{\nu_k}{\sigma_k^2} \right) \quad (22)$$

$$\kappa_k \sin(\zeta_k) = -\frac{1}{1 - \rho_k^2} \left(\frac{\rho_k \nu_k}{\sigma_k \bar{\sigma}_k} - \frac{\bar{\nu}_k}{\bar{\sigma}_k^2} \right) \quad (23)$$

$$\bar{\kappa}_k \cos(2\bar{\zeta}_k) = -\frac{1}{4(1 - \rho_k^2)} \left(\frac{1}{\sigma_k^2} - \frac{1}{\bar{\sigma}_k^2} \right) \quad (24)$$

$$\bar{\kappa}_k \sin(2\bar{\zeta}_k) = \frac{\rho_k}{2(1 - \rho_k^2) \sigma_k \bar{\sigma}_k}, \quad (25)$$

where

$$\begin{bmatrix} \nu_k \\ \bar{\nu}_k \end{bmatrix} \triangleq \beta_k^{-1} \left(\begin{bmatrix} \text{Re}\{y\} \\ \text{Im}\{y\} \end{bmatrix} - \boldsymbol{\mu}_k \right) \quad (26)$$

and

$$\begin{bmatrix} \sigma_k^2 & \rho_k \sigma_k \bar{\sigma}_k \\ \rho_k \sigma_k \bar{\sigma}_k & \bar{\sigma}_k^2 \end{bmatrix} \triangleq \beta_k^{-2} \boldsymbol{\Sigma}_k. \quad (27)$$

Given the SHyGAMP pseudo-prior $\mathbf{z}_k \sim \mathcal{N}(\widehat{p}_k, [\mathbf{Q}^{\mathbf{p}}]_{kk})$, the posterior on θ_k takes the form

$$p(\theta_k | y) \propto \mathcal{N}(\theta_k; g\widehat{p}_k, g^2[\mathbf{Q}^{\mathbf{p}}]_{kk}) p(y|\theta_k) \quad (28)$$

$$\propto \exp \left[\kappa_k \cos(\theta_k - \zeta_k) + \bar{\kappa}_k \cos[2(\theta_k - \bar{\zeta}_k)] - \frac{(\theta_k - g\widehat{p}_k)^2}{2g^2[\mathbf{Q}^{\mathbf{p}}]_{kk}} \right].$$

We then face the task of computing $\mathbb{E}\{\theta_k | y\}$ and $\mathbb{E}\{\theta_k^2 | y\}$ under (28). Several methods could be applied here, such as numerical integration. The method employed for the experiments in Section III is based on the Laplace approximation [15]. For this, we compute $\widehat{\theta}_{k,\text{MAP}} \triangleq \arg \max_{\theta_k} \ln p(\theta_k | y)$ using bisection and then approximate $\mathbb{E}\{\theta_k | y\} \approx \widehat{\theta}_{k,\text{MAP}}$ and $\text{var}\{\theta_k | y\} \approx -\frac{d^2}{d\theta_k^2} \ln p(\theta_k | y) \big|_{\theta_k = \widehat{\theta}_{k,\text{MAP}}}$. Finally, we compute $\widehat{z}_k = \mathbb{E}\{\theta_k | y\} / g$ and $q_k^{\mathbf{z}} = \text{var}\{\theta_k | y\} / g^2$.

2) *Inference of \mathbf{c}_n* : For lines 10-11 of Algorithm 1, recall that $p_{\mathbf{c}}(\mathbf{c}_n) = \mathcal{N}(\mathbf{c}_n; \mathbf{0}, \nu \mathbf{I})$. Thus $p_{\mathbf{c}|\mathbf{r}}$ is Gaussian and the posterior mean and covariance of \mathbf{c}_n can be computed straightforwardly as

$$\mathbf{Q}_n^{\mathbf{c}} = (\nu^{-1} \mathbf{I} + [\mathbf{Q}^{\mathbf{r}}]^{-1})^{-1} \triangleq \mathbf{Q}^{\mathbf{c}} \quad (29)$$

$$\widehat{\mathbf{c}}_n = \mathbf{Q}^{\mathbf{c}} [\mathbf{Q}^{\mathbf{r}}]^{-1} \widehat{\mathbf{r}}_n \quad (30)$$

Above, $[\mathbf{Q}^{\mathbf{r}}]^{-1}$ is simplified by the fact that $\mathbf{Q}^{\mathbf{r}}$ is diagonal.

E. Hyperparameter Tuning

The likelihood model $p_{\mathbf{y}|\mathbf{z}}$ in (9) depends on the unknown hyperparameters $\boldsymbol{\alpha}$ and $\boldsymbol{\tau}$. Similarly, the prior $p_{\mathbf{c}}$ depends on the unknown variance ν . We propose to estimate these hyperparameters using a combination of *expectation maximization* (EM) and SHyGAMP, as suggested in [8] and detailed—for the simpler case of GAMP—in [16]. Extrapolating [16] to the SHyGAMP case, we estimate $\boldsymbol{\alpha}$ and $\boldsymbol{\tau}$ via

$$\{\widehat{\boldsymbol{\alpha}}, \widehat{\boldsymbol{\tau}}\} = \arg \max_{\boldsymbol{\alpha} \geq \mathbf{0}, \boldsymbol{\alpha}^T \mathbf{1} = 1, \boldsymbol{\tau} > \mathbf{0}} \sum_{m=1}^M \int_{\mathbb{R}^K} \mathcal{N}(\mathbf{z}_m; \widehat{\mathbf{z}}_m, \text{Diag}(\mathbf{q}_m^{\mathbf{z}})) \times \ln p(y_m | \mathbf{z}_m; \boldsymbol{\alpha}, \boldsymbol{\tau}) d\mathbf{z}_m \quad (31)$$

at each SHyGAMP iteration, immediately after line 5 in Algorithm 1. For tractability, we approximate the Dirac delta in (9) by a Gaussian pdf with small variance $\epsilon > 0$, giving

$$\ln p(y_m | \mathbf{z}_m; \boldsymbol{\alpha}, \boldsymbol{\tau}) \approx -\frac{1}{\epsilon} \left| y - \sum_{k=1}^K \alpha_k \exp(jg_m z_{mk} - \frac{g_m^2 \tau_k}{2}) \right|^2 + \text{const.} \quad (32)$$

The resulting optimization problem (31) (which does not depend on ϵ) can be straightforwardly solved using gradient projection, since closed-form expressions for the objective and its gradient exist.

III. NUMERICAL EXPERIMENTS

In this section, we present the results of two numerical 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 [17] and enabled the “++” initialization method. CL-OMPR and CL-AMP used the same sketch \mathbf{y} , whose frequency vectors \mathbf{W} were drawn using the method described in [5]. For both experiments, the clusters were randomly drawn as $\mathbf{c}_k \sim \mathcal{N}(\mathbf{0}_N, 1.5^2 K^{2/N} \mathbf{I}_N)$, after which the training (and test) data were drawn from the GMM (5) with weights $\alpha_k = \frac{1}{K} \forall k$ and covariances $\mathbf{R}_k = \mathbf{I}_N \forall k$. For CL-OMPR and CL-AMP, the runtimes reported include the time of computing the sketch.

A. SSE Minimization

In the first experiment, we test each algorithm’s ability to minimize SSE on the training data, i.e., to solve the problem (1). For each pair of $(K, N) \in \{(5, 100), (10, 50), (10, 100)\}$, 10 trials were performed, where in each trial, a training dataset was randomly generated with $T = 10^4$ samples. For cluster recovery, k-means++ was invoked on this training dataset with 1 replicate (i.e., 1 run from a random initialization), while CL-AMP and CL-OMPR were invoked using a sketch of length M . Several values of M , logarithmically spaced in the interval $[KN, 10KN]$, were evaluated.

For each (K, N) pair under test, Figs. 1a, 1c, and 1e show the median SSE of CL-AMP and CL-OMPR versus M/KN , with the error-bars showing the standard deviation. The median SSE of k-means++ was superimposed on these figures as a reference, although k-means++ has no dependence on M . Likewise, Figs. 1b, 1d, and 1f show the corresponding median runtime for CL-AMP and CL-OMPR vs M/KN , where again the result for k-means++ was superimposed. Because a low runtime is meaningless if the corresponding SSE is very high, the runtime was not shown for CL-AMP or CL-OMPR when its SSE was more than twice that of k-means++.

Figure 1 shows that CL-AMP achieved a low SSE with fewer measurements M than CL-OMPR. In particular, CL-AMP required $M \approx 3KN$ to minimize the SSE, while CL-OMPR required $M \approx 10KN$. Also, the minimum SSE achieved by CL-AMP was in most cases lower than that of k-means++ and CL-OMPR. As we will see in Fig. 2, the SSE

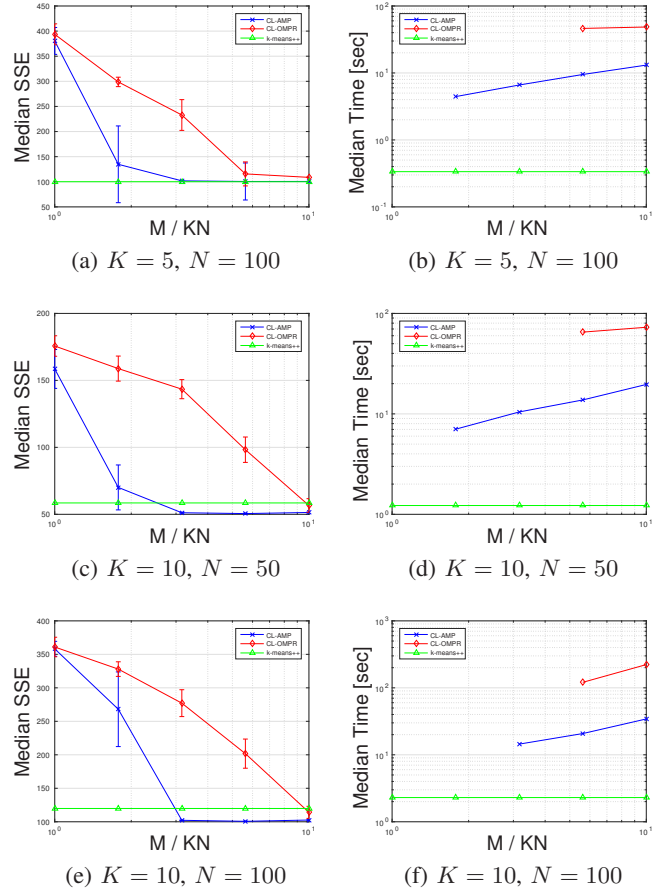


Fig. 1: Median sum-squared error and runtime vs M/KN .

performance of k-means++ can be improved (at the expense of runtime) by increasing the number of replicates. Figure 1 also shows that CL-AMP was an order of magnitude faster than CL-OMPR for all (K, N, M) under test. Meanwhile, CL-AMP was approximately one order-of-magnitude slower than k-means++, although the SSE achieved by CL-AMP was often lower. A direct SSE-vs-runtime comparison is given below.

B. Performance versus Runtime

The previous experiment demonstrated CL-AMP’s ability to minimize SSE faster and with fewer measurements than CL-OMPR. However, the comparison with k-means++ was inconclusive: k-means++ was faster but achieved a worse SSE in many cases. We now describe a different experiment that aimed to evaluate clustering performance versus runtime. For clustering performance, we consider both training SSE and classification error on test data. For the latter, training data is used for cluster recovery and the estimated clusters are used for minimum-distance classification of test data. To control the performance and computational complexity of k-means++, we allowed multiple replicates as well as training-data subsampling. Details are provided in the sequel.

We first drew $K = 30$ random centroids of dimension $N = 20$ under the previously described GMM. Then we

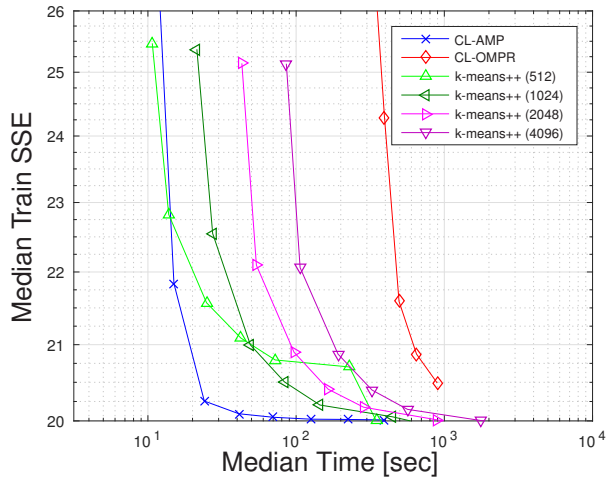


Fig. 2: Training sum-squared error vs runtime. Each k-means++ trace corresponds to a different number of replicates.

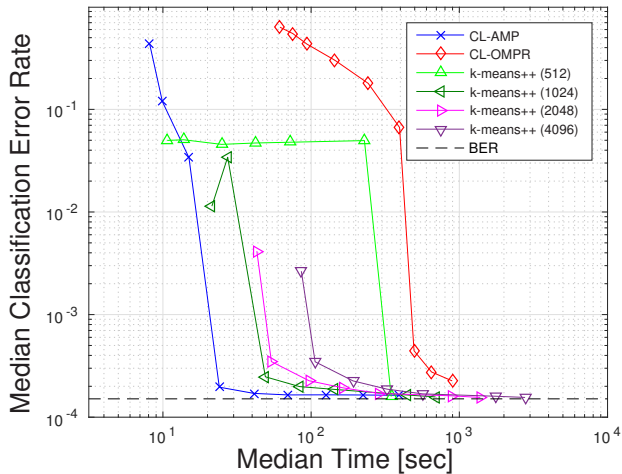


Fig. 3: Classification error rate vs runtime. Each k-means++ trace corresponds to a different number of replicates.

generated $T = 10^4$ random training samples from this GMM. To recover clusters, CL-AMP and CL-OMPR were applied with sketch length M , while k-means++ was applied with random subsampling of the training set and multiple replicates. We tested several sketch lengths $M \in [KN, 100KN]$, k-means++ sampling rates $\in [0.5^6, 1]$, and k-means++ replicates $\in [512, 4096]$, all logarithmically spaced. Finally, the resulting training-data SSE was evaluated using the full training dataset.

The quality of the estimated centroids was also evaluated by computing the error-rate of minimum-distance classification of a test dataset (of size $T_{\text{test}} = 5 \times 10^6$, drawn from the same GMM as the training data). Here, we used the Hungarian algorithm to assign labels to the estimated centroids.

Figure 2 shows median training SSE versus runtime over 10 trials for each algorithm under test, while Fig. 3 shows the corresponding median test error rate versus runtime. In

each CL-AMP and CL-OMPR trace, the different datapoints correspond to increasing values of sketch length M , while in each k-means++ trace, the different datapoints correspond to increasing sampling rates for a fixed number of replicates (specified in the legend). Figure 3 shows the corresponding classification error rate, computed on the test set, as well as the Bayes (i.e., minimum possible) classification error rate.

Figures 2 and 3 tell a similar story: to achieve near-optimal training-SSE or classification error-rate with this GMM, CL-AMP (with properly adjusted M) requires less runtime than CL-OMPR or any variation of k-means++. Note that CL-AMP is easily “tuned” by choosing $M \approx 5KN$, while k-means++ is much more difficult to tune: it is not clear how to choose the best combination of subsampling rate and number-of-replicates to achieve both low SSE and low runtime. Note also that the implementation of k-means++ is highly optimized while that of CL-AMP is not, so further improvements may be possible by optimizing the implementation of CL-AMP.

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, pp. 651–666, June 2010.
- [4] D. Arthur and S. Vassilvitskii, “k-means++: The advantages of careful seeding,” in *Proc. ACM-SIAM Symp. Discrete Alg.*, pp. 1027–1035, 2007.
- [5] N. Keriven, A. Bourrier, R. Gribonval, and P. Perez, “Sketching for large-scale learning of mixture models,” Jun 2016. (found at [arXiv:1606.02838](https://arxiv.org/abs/1606.02838)).
- [6] N. Keriven, N. Tremblay, Y. Traonmilin, and R. Gribonval, “Compressive k-means,” Oct 2016. (found at [arXiv:1610.08738](https://arxiv.org/abs/1610.08738)).
- [7] 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), pp. 40–44, 1993.
- [8] 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.
- [9] M. Rudelson and R. Vershynin, “Hanson-Wright inequality and sub-Gaussian concentration,” *Electron. Commun. Probab.*, vol. 18, no. 82, pp. 1–9, 2013.
- [10] D. L. Donoho, A. Maleki, and A. Montanari, “Message passing algorithms for compressed sensing,” *Proc. Nat. Acad. Sci.*, vol. 106, pp. 18914–18919, Nov. 2009.
- [11] S. Rangan, “Generalized approximate message passing for estimation with random linear mixing,” in *Proc. IEEE Int. Symp. Inform. Thy.*, pp. 2168–2172, Aug. 2011. (full version at [arXiv:1010.5141](https://arxiv.org/abs/1010.5141)).
- [12] M. Bayati and A. Montanari, “The dynamics of message passing on dense graphs, with applications to compressed sensing,” *IEEE Trans. Inform. Theory*, vol. 57, pp. 764–785, Feb. 2011.
- [13] 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.
- [14] R. Gatto and S. R. Jammalamadaka, “The generalized von Mises distribution,” *Stat. Method.*, vol. 4, pp. 341–353, 2007.
- [15] C. M. Bishop, *Pattern Recognition and Machine Learning*. New York: Springer, 2007.
- [16] J. P. Vila and P. Schniter, “Expectation-maximization Gaussian-mixture approximate message passing,” *IEEE Trans. Signal Process.*, vol. 61, pp. 4658–4672, Oct. 2013.
- [17] N. Keriven, N. Tremblay, and R. Gribonval, “SketchMLbox : a Matlab toolbox for large-scale learning of mixture models,” 2016. <http://sketchml.gforge.inria.fr>.