

Sketched Clustering via Approximate Message Passing

Phil Schniter



THE OHIO STATE UNIVERSITY

Joint work with **Evan Byrne** (Ohio State),
Antoine Chatalic, and **Rémi Gribonval** (INRIA)

Supported in part by NSF Grant 1716388 and MIT Lincoln Laboratory

SPARS Workshop (Toulouse, France)

July 3, 2019

Outline

- 1 Sketched Clustering
- 2 Cluster Recovery via EM and AMP
- 3 Numerical Experiments
 - Synthetic Data
 - Spectral MNIST
 - Spike Super-Resolution Recovery from Fourier Samples

Outline

- 1 Sketched Clustering
- 2 Cluster Recovery via EM and AMP
- 3 Numerical Experiments
 - Synthetic Data
 - Spectral MNIST
 - Spike Super-Resolution Recovery from Fourier Samples

Clustering with K-Means

- **Given:** T feature vectors $\{\mathbf{x}_t\}$ with $\mathbf{x}_t \in \mathbb{R}^N$
- **Goal:** Find K centroids $\{\mathbf{c}_k\}$ that minimize sum of squared errors:

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

- Finding the SSE-minimizing centroids is NP-hard
- K-means++ is the standard heuristic approach:
 - Lloyd's algorithm plus a careful random initialization
 - Per-iteration complexity of $O(NKT)$
 - Challenge: Complexity and memory can be prohibitive for large T

Sketched Learning

- **Sketched learning** is an alternative framework:

- 1 Compress data $\mathbf{X} \in \mathbb{R}^{N \times T}$ down to $\mathbf{y} \in \mathbb{C}^M$ (with $M \ll NT$).
- 2 Learn parameters (e.g., centroids) from \mathbf{y} .

- We choose to build the sketch $\mathbf{y} = [y_1, \dots, y_M]^\top$ using¹⁴

$$y_m = \frac{1}{T} \sum_{t=1}^T \exp(j\mathbf{w}_m^\top \mathbf{x}_t) \quad \text{with random } \{\mathbf{w}_m\}_{m=1}^M$$

- Well matched to **distributed** and/or **streaming** scenarios!
 - Complexity & memory of learning are invariant to T !
- Can interpret y_m as samples of the **empirical characteristic function**:

$$y_m = \phi(\mathbf{w}_m) = \int_{\mathbf{R}^N} p(\mathbf{x}) \exp(j\mathbf{w}_m^\top \mathbf{x}) d\mathbf{x} \quad \text{with } p(\mathbf{x}) = \frac{1}{T} \sum_{t=1}^T \delta(\mathbf{x} - \mathbf{x}_t)$$

¹Keriven, Bourrier, Gribonval, Pérez'17, ⁴Keriven, Tremblay, Traonmilin, Gribonval'17

Sketched Clustering

- How do we learn the centroids C from the sketch y ?

- The CL-OMPR algorithm³⁴ aims to solve

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

using a greedy heuristic inspired by OMP.

- In practice, CL-OMPR ...

- recovers accurate centroids with sketch length $M \approx 10KN$
- has a per-iteration complexity of $O(MNK^2)$

- Can we do better in terms of **sample complexity** and **computational complexity**?

³Keriven, Bourrier, Gribonval, Pérez'17, ⁴Keriven, Tremblay, Traonmilin, Gribonval'17

Outline

- 1 Sketched Clustering
- 2 Cluster Recovery via EM and AMP
- 3 Numerical Experiments
 - Synthetic Data
 - Spectral MNIST
 - Spike Super-Resolution Recovery from Fourier Samples

Formulation as a Generalized Linear Model

- Suppose we model the data \mathbf{x}_t using a **Gaussian mixture model** (GMM):

$$\mathbf{x}_t \sim \sum_{k=1}^K \alpha_k \mathcal{N}(\mathbf{c}_k, \Phi_k) \quad \text{with} \quad \sum_{k=1}^K \alpha_k = 1, \quad \alpha_k \geq 0, \quad \Phi_k > 0.$$

- As $T \rightarrow \infty$, have $y_m = \frac{1}{T} \sum_{t=1}^T \exp(\mathbf{j} \mathbf{w}_m^\top \mathbf{x}_t) \rightarrow \mathbb{E} \{ \exp(\mathbf{j} \mathbf{w}_m^\top \mathbf{x}_t) \}$

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

where $g_m \triangleq \|\mathbf{w}_m\|$ and $\mathbf{a}_m \triangleq \mathbf{w}_m / g_m$. $\triangleq z_{mk}$ $\triangleq \tau_{mk}$

- As $N \rightarrow \infty$, with isotropic \mathbf{a}_m , we have $\tau_{mk} \rightarrow \text{tr}(\Phi_k) / N \triangleq \tau_k$.
- Thus for large T and N we have the **generalized linear model** (GLM)

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

with transformed centroids $\mathbf{Z} = \mathbf{A} \mathbf{C}$ & random \mathbf{A} w/ isotropic columns

Sketched Clustering via EM

- Objective: Recover the centroids \mathbf{C} from the sketch \mathbf{y} under the GLM

$$p(\mathbf{y}|\mathbf{Z}; \boldsymbol{\alpha}, \boldsymbol{\tau}) = \prod_{m=1}^M p(y_m|z_m; \boldsymbol{\alpha}, \boldsymbol{\tau}), \quad \mathbf{Z} = \mathbf{A}\mathbf{C}$$

- Challenge: GMM weights $\boldsymbol{\alpha}$ and variances $\boldsymbol{\tau}$ are unknown!

- Approach: **Expectation Maximization** (EM): Iterate ...

$$\begin{aligned} (\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\tau}})^{\text{new}} &= \arg \max_{(\boldsymbol{\alpha}, \boldsymbol{\tau}): \boldsymbol{\alpha}^T \mathbf{1} = 1, \boldsymbol{\alpha} \geq \mathbf{0}, \boldsymbol{\tau} > \mathbf{0}} \mathbb{E} \{ \ln p(\mathbf{y}, \mathbf{Z}; \boldsymbol{\alpha}, \boldsymbol{\tau}) \mid \mathbf{y}, \hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\tau}} \} \\ &= \arg \max_{(\boldsymbol{\alpha}, \boldsymbol{\tau}): \boldsymbol{\alpha}^T \mathbf{1} = 1, \boldsymbol{\alpha} \geq \mathbf{0}, \boldsymbol{\tau} > \mathbf{0}} \sum_{m=1}^M \int_{\mathbb{R}^K} \mathcal{N}(z_m; \hat{z}_m, \mathbf{Q}_m^z) \ln p(y_m|z_m; \boldsymbol{\alpha}, \boldsymbol{\tau}) \end{aligned}$$

with conditional mean $\hat{z}_m = \mathbb{E}\{z_m \mid \mathbf{y}; \hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\tau}}\}$ and conditional covariance \mathbf{Q}_m^z .

- Thus we aim to compute **MMSE centroid estimates** $\hat{\mathbf{C}} = \mathbb{E}\{\mathbf{C} \mid \mathbf{Y}; \hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\tau}}\}$, since 1) they provide $\hat{\mathbf{Z}} = \mathbf{A}\hat{\mathbf{C}}$ for EM and 2) solve our sketched clustering problem.

MMSE Inference for Sketched Clustering

- Objective: Compute MMSE centroid estimate $\hat{\mathbf{C}}$ from \mathbf{y} under GLM

$$p(\mathbf{y}|\mathbf{Z}) = \prod_{m=1}^M p_{y|z}(y_m|z_m; \hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\tau}}), \quad \mathbf{Z} = \mathbf{AC}.$$

- Note that the posterior centroid density is

$$p(\mathbf{C}|\mathbf{y}) \propto \prod_{m=1}^M p_{y|z}(y_m|\mathbf{a}_m^T \mathbf{C}; \hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\tau}}) \prod_{n=1}^N p_{\mathbf{c}}(\mathbf{c}_n)$$

- We assume the trivial centroid prior $p_{\mathbf{c}}(\mathbf{c}_n) \propto 1$, but other priors are possible
- We can approximately compute $\hat{\mathbf{C}}$ using [approximate message passing](#)
 - Due to the form of the likelihood, we use the “HyGAMP” algorithm⁵

⁵Rangan, Fletcher, Goyal, Schniter'12

Lineage of HyGAMP

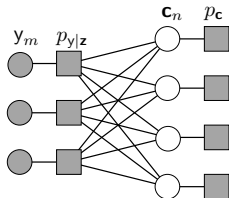
- **Approximate Message Passing (AMP)** [Donoho, Maleki, Montanari'09]
 - Estimate \mathbf{c} under the **standard linear model** $\mathbf{y} = \mathbf{A}\mathbf{c} + \mathbf{w}$ with known iid \mathbf{A}
 - Assumes separable prior $p_{\mathbf{c}}(\mathbf{c}) = \prod_n p_{\mathbf{c}}(c_n)$ and AWGN \mathbf{w}

- **Generalized AMP (GAMP)** [Rangan'11]
 - Estimate \mathbf{c} under **generalized linear model** $\mathbf{y} \sim p(\mathbf{y}|\mathbf{z})$ with $\mathbf{z} = \mathbf{A}\mathbf{c}$
 - Assumes separable prior and likelihood $p(\mathbf{y}|\mathbf{z}) = \prod_m p_{y|z}(y_m|z_m)$

- **Hybrid GAMP (HyGAMP)** [Rangan, Fletcher, Goyal, Schniter'12]
 - **GAMP with vector-valued variables** $\mathbf{z}_m, \mathbf{c}_n \in \mathbb{R}^K$
 - Separable likelihood: $\mathbf{y} \sim p(\mathbf{y}|\mathbf{Z}) = \prod_m p_{y|z}(y_m|\mathbf{z}_m)$ with $\mathbf{Z} = \mathbf{A}\mathbf{C}$
 - Separable prior: $p(\mathbf{C}) = \prod_n p_{\mathbf{c}}(\mathbf{c}_n)$

Message-Passing View of HyGAMP

- HyGAMP can be derived by approximating belief propagation (either sum-product or max-product algorithm) on a factor graph with the form:



- Messages are approximated as K -dimensional Gaussian pdfs assuming $N \rightarrow \infty$
- HyGAMP tackles the (NK) -dimensional inference problem by **iteratively solving $M+N$ inference problems of dimension K**

HyGAMP Inference Steps

- HyGAMP's K -dimensional inference steps compute the **posterior mean and covariance** of the random vectors $\{\mathbf{c}_n\}$ and $\{\mathbf{z}_m\}$ under the posterior pdfs

$$p(\mathbf{c}_n | \mathbf{r}_n; \mathbf{Q}^r) \propto p_{\mathbf{c}}(\mathbf{c}_n) \mathcal{N}(\mathbf{c}_n; \mathbf{r}_n, \mathbf{Q}^r)$$

$$p(\mathbf{z}_m | y_m, \mathbf{p}_m; \mathbf{Q}^p) \propto p_{y|z}(y_m | \mathbf{z}_m) \mathcal{N}(\mathbf{z}_m; \mathbf{p}_m, \mathbf{Q}^p)$$

- The correctness of these posteriors can be argued, under large i.i.d. Gaussian \mathbf{A} , using the analysis in [Javanmard, Montanari'13]
- To reduce computational complexity, we use the Simplified HyGAMP (SHyGAMP) algorithm,⁶ which **approximates covariance matrices as diagonal**
 - The per-iteration complexity of SHyGAMP is only $O(MNK)$.
- For the sketched-clustering likelihood $p_{y|z}(y_m | \mathbf{z}_m)$, the computation of $\hat{\mathbf{z}}_m$ and $\text{diag}(\mathbf{Q}_m^z)$ uses generalized von Mises functions, and is somewhat involved.⁷

⁶Byrne, Schniter'15, ⁷Byrne, Chatalic, Gribonval, Schniter'19

Outline

- 1 Sketched Clustering
- 2 Cluster Recovery via EM and AMP
- 3 Numerical Experiments
 - Synthetic Data
 - Spectral MNIST
 - Spike Super-Resolution Recovery from Fourier Samples

Experiment 1: Synthetic Data

Data generation:

- $\{\mathbf{x}_t\}$ drawn i.i.d. from a GMM with
 - centroids \mathbf{c}_k drawn $\sim \mathcal{N}(\mathbf{0}, 1.5^2 K^{2/N} \mathbf{I}_N)$
 - equal weights $\alpha_k = 1/K$
 - covariances $\mathbf{\Phi}_k = \mathbf{I}$
- $N = 100$ dimensional, $K = 10$ classes, $T = 10^7$ samples

Sketching:⁸

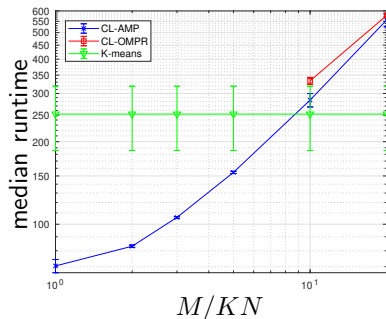
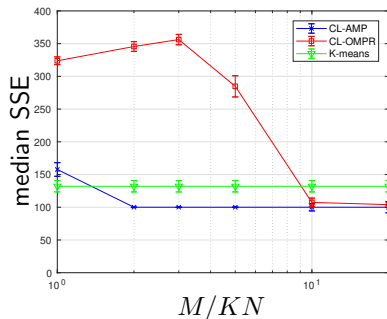
- frequencies $\mathbf{w}_m = g_m \mathbf{a}_m$ with unit-norm isotropic \mathbf{a}_m
- $g_m \sim p(g) = 1_{[0, \infty)} \sqrt{g^2 \sigma^2 + \frac{g^4 \sigma^4}{4}} \exp(-g^2 \sigma^2 / 2)$ with $\sigma^2 = \frac{1}{NT} \|\mathbf{X}\|_F^2$

Accuracy metric:

- median of $\text{SSE}(\hat{\mathbf{C}}) = \sum_{t=1}^T \min_k \|\mathbf{x}_t - \hat{\mathbf{c}}_k\|_2^2$ over 10 trials

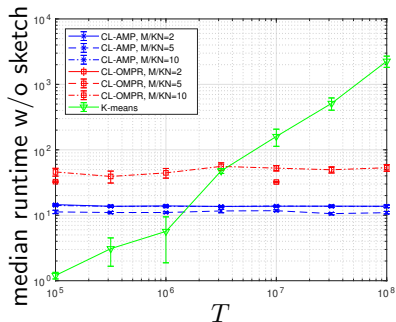
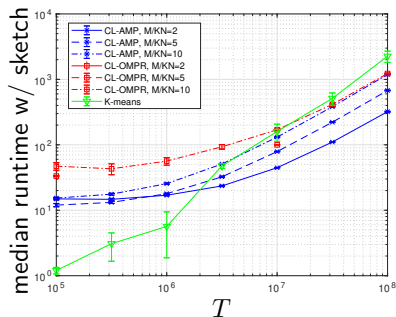
⁸Keriven, Bourier, Gribonval, Perez'17

Accuracy & Runtime vs Sketch Length M



- Sample complexity:
 - CL-AMP needs only $M \approx 2KN$ samples
 - CL-OMPR needs $M \approx 10KN$
- Computational complexity (including sketch):
 - CL-AMP $3\times$ faster than K-means++ for similar accuracy

Runtime vs Data Size T



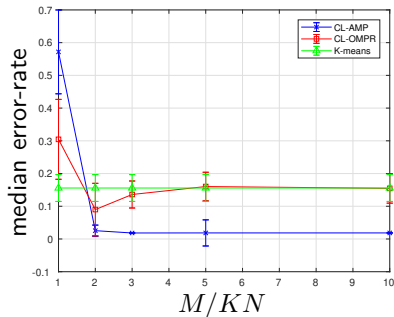
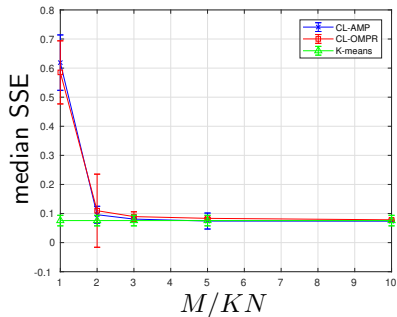
When $T > 2 \times 10^6 \dots$

- sketching+CL-AMP is faster than K-means++
- sketching is more expensive than CL-AMP

Experiment 2: Spectral Clustering of MNIST

- We repeat an experiment from [Keriven, Tremblay, Traonmilin, Gribonval'17]
- Original MNIST data:
 - $T = 70,000$ samples of handwritten digits from $K = 10$ classes
- Preprocessing used to extract features of dimension $N = 10$
 - Compute SIFT descriptors
 - Compute k -NN adjacency matrix (for $k=10$) using FLANN
 - Compute $K=10$ principle eigenvectors of normalized Laplacian matrix
- Dataset partitioned into equal-sized training and test sets (10 trials)
- Kmeans++, CL-OMPR, and CL-AMP estimate $K = 10$ centroids from training set
- Accuracy metrics:
 - 1) SSE on training set
 - 2) error of minimum-distance classifier on test set

Accuracy vs Sketch Length M



For $M \geq 2KN \dots$

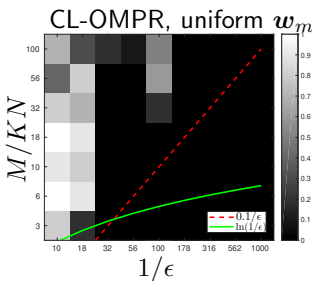
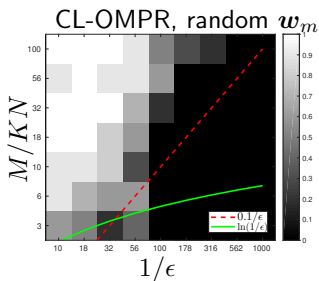
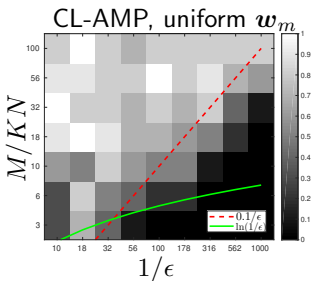
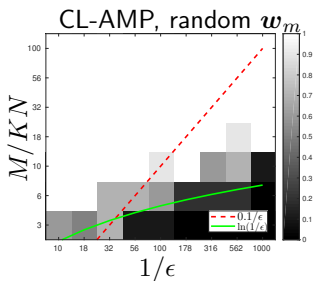
- CL-OMPR and CL-AMP give **SSE similar** to that of k-means++
- CL-AMP gives **error rate much better** than CL-OMPR and k-means++

Experiment 3: Spike Super-Resolution w/ Fourier Samples

- **Sum-of-spikes signal:** $\sum_{k=1}^K \alpha_k \delta(\mathbf{t} - \mathbf{c}_k)$ with time $\mathbf{t} \in \mathbb{R}^N$
- **Fourier transform:** $y(\mathbf{w}) = \sum_{k=1}^K \alpha_k \exp(\mathbf{j}\mathbf{w}^\top \mathbf{c}_k)$ with freq $\mathbf{w} \in \mathbb{R}^N$
- **Goal:** Recover $\{\mathbf{c}_k\}_{k=1}^K$ from Fourier samples $\{y(\mathbf{w}_m)\}_{m=1}^M$
- **Experiment:**
 - Generate frequency pairs $\{(\mathbf{c}_{2i-1}, \mathbf{c}_{2i})\}_{i=1}^{K/2}$ with $\|\mathbf{c}_{2i-1} - \mathbf{c}_{2i}\| = \epsilon \forall i$
 - “Success” if $\max_k \|\hat{\mathbf{c}}_k - \mathbf{c}_{i_k}\| < \epsilon/2$ for some $\{i_1, \dots, i_K\} = \{1, \dots, K\}$
 - Theoretical analysis⁹ says that
 - $M \geq O(\log(1/\epsilon))$ samples suffice for **random** frequencies $\{\mathbf{w}_m\}$
 - $M \geq O(1/\epsilon)$ samples suffice for **uniformly spaced** frequencies $\{\mathbf{w}_m\}$

⁹Traonmilin, Keriven, Gribonval, Blanchard'17

Frequency Estimation Results ($K = 4, N = 2$)



Conclusion

- **Sketched clustering** is an alternative to traditional clustering that
 - 1 compresses the dataset down to a sketch (of generalized moments)
 - 2 extracts centroids from that sketchand is well matched to **distributed** and/or **streamed** scenarios
- We formulated sketched clustering as a **GLM inference** problem, and applied **EM-SHyGAMP**.
- Numerical results suggest that has CL-AMP has good sample & computational complexity
- Ongoing work to analyze the AMP **state evolution** in the large-system limit ($N, M \rightarrow \infty$)

Full paper

E. Byrne, A. Chatalic, R. Gribonval, and P. Schniter, “Sketched Clustering via Hybrid Approximate Message Passing,” *IEEE Trans. Signal Processing*, to appear 2019 (see also <https://arxiv.org/abs/1712.02849>).