

# Sketched Clustering via Hybrid Approximate Message Passing

Evan Byrne, Rémi Gribonval, and Philip Schniter



(Supported by NSF Grant 1716388 and MIT Lincoln Labs)

## Traditional Clustering Problem Statement

- Given a dataset of  $T$   $N$ -dimensional feature vectors  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_T] \in \mathbb{R}^{N \times T}$ , estimate  $K$   $N$ -dimensional cluster centers  $\mathbf{C} = [\mathbf{c}_1, \dots, \mathbf{c}_K] \in \mathbb{R}^{N \times K}$  that minimize sum of squared errors (SSE):

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

- However, finding  $\mathbf{C}$  to minimize the SSE in (1) is NP-hard.
- K-means is a commonly applied heuristic approach.
- K-means generally works well wrt minimizing the SSE, except its complexity is  $\mathcal{O}(NKTl)$ , where  $l$  is the number of iterations, which is prohibitive for large  $T$ .

## Sketched Clustering

- Sketched clustering [Kerivan 16] is an alternate approach possibly more efficient than K-means.
- Let  $\mathbf{y} \in \mathbb{C}^M$  be the "sketch" of  $\mathbf{X}$ , where

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

for some set of  $N$ -dimensional frequency vectors  $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_M]$ .

- The sketch in (2) can be interpreted as the empirical characteristic function of the dataset  $\mathbf{X}$ .
- CLOMPR [Kerivan 17] is the state-of-the-art Sketched Clustering algorithm, which solves

$$\{\hat{\mathbf{C}}, \hat{\boldsymbol{\alpha}}\} = \arg \min_{\mathbf{C}, \boldsymbol{\alpha}} \sum_{m=1}^M \left| \mathbf{y}_m - \sum_{k=1}^K \alpha_k \exp(j\mathbf{w}_m^T \mathbf{c}_k) \right|^2 \quad (3)$$

via a greedy optimization approach.

- In practice,  $\hat{\mathbf{C}}_{\text{CLOMPR}}$  works well wrt SSE compared to  $\hat{\mathbf{C}}_{\text{K-means}}$ , despite no link between (3) and (1).
- CLOMPR's complexity is  $\mathcal{O}(MNK^2l + MNT)$ , which includes the cost of computing  $\mathbf{y}$ .
- Note that once  $\mathbf{y}$  is computed,  $\mathbf{X}$  is not stored during CLOMPR, so the memory requirement is significantly reduced.
- CLOMPR's authors have developed several approaches for randomly generating the frequencies  $\mathbf{w}_m$  and have observed around  $M \approx 10KN$  frequencies necessary for accurate performance.

## Sketched Clustering via Approximate Message Passing

- We choose to model the feature vectors  $\mathbf{x}_t$  with a Gaussian Mixture where the mixture centers are the "true" cluster centers, i.e.,

$$\mathbf{x}_t \sim \sum_{k=1}^K \alpha_k \mathcal{N}(\mathbf{c}_k, \boldsymbol{\Sigma}_k). \quad (4)$$

- Then, for large  $T$ ,

$$\mathbf{y}_m = \frac{1}{T} \sum_{t=1}^T \exp(j\mathbf{w}_m^T \mathbf{x}_t) \approx \mathbb{E}\{\exp(j\mathbf{w}_m^T \mathbf{x}_t)\} = \sum_{k=1}^K \alpha_k \exp(j\mathbf{w}_m^T \mathbf{c}_k - \frac{\mathbf{w}_m^T \boldsymbol{\Sigma}_k \mathbf{w}_m}{2}), \quad (5)$$

and so

$$p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}_m | \mathbf{z}_m) = \delta\left(\mathbf{y}_m - \sum_{k=1}^K \alpha_k \exp(j\mathbf{z}_{mk} - \tau_{mk}/2)\right), \quad (6)$$

where  $\{\tau_{mk}\}$  and  $\{\alpha_k\}$  are treated as hyperparameters.

- If we assume  $p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}_m | \mathbf{z}_m)$  are independent across  $m$  and assume  $p_{\mathbf{C}}(\mathbf{C}) = \prod_{n=1}^N p_{\mathbf{C}}(\mathbf{c}_n)$ , we obtain

$$p_{\mathbf{y}, \mathbf{C}}(\mathbf{y}, \mathbf{C}) = \prod_{m=1}^M p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}_m | \mathbf{w}_m^T \mathbf{C}) \prod_{n=1}^N p_{\mathbf{C}}(\mathbf{c}_n). \quad (7)$$

- With (7), we treat sketched clustering as an inference problem rather than an optimization problem.
- In particular, we approximate

$$\hat{\mathbf{C}} = \mathbb{E}\{p_{\mathbf{C}|\mathbf{y}}(\mathbf{C} | \mathbf{y})\}, \quad (8)$$

using the Simplified-Hybrid-GAMP (SHyGAMP) algorithm [Byrne 16].

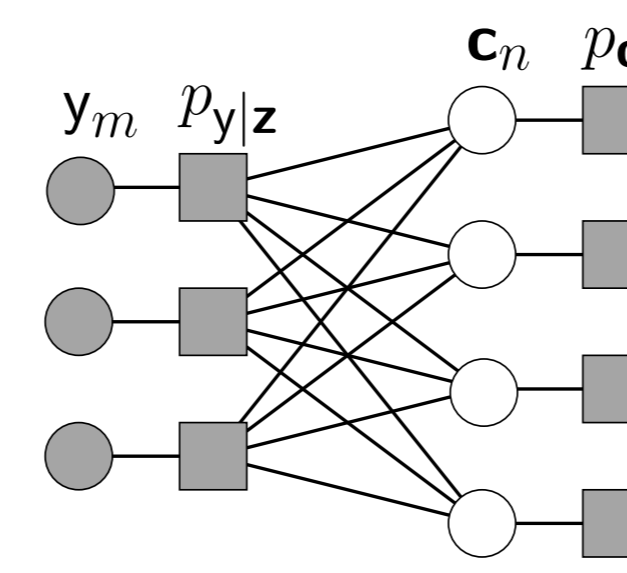
- The SHyGAMP algorithm is based on the more general HyGAMP algorithm [Rangan 17]. The only difference between the two is SHyGAMP restricts the messages that are passed to have diagonal covariance matrices, which drastically reduces computational complexity.

## References

- N. Kerivan, A. Bourrier, R. Gribonval, and P. P  rez, "Sketching for Large-Scale Learning of Mixture Models," in *Proc. IEEE Intl. Conf. on Acoustic, Speech and Signal Processing (ICASSP)*, March 2016.
- N. Kerivan, N. Tremblay, Y. Traonmilin, and R. Gribonval, "Compressive K-means," in *Proc. IEEE Intl. Conf. on Acoustic, Speech and Signal Processing (ICASSP)*, March 2017.
- E. Byrne and P. Schniter, "Sparse Multinomial Logistic Regression via Approximate Message Passing," in *IEEE Trans. on Signal Processing*, vol. 64, no. 21, pp. 5485-5498, Nov. 2016.
- S. Rangan, A. K. Fletcher, V. K. Goyal, E. Byrne, and P. Schniter, "Hybrid Approximate Message Passing," in *IEEE Trans. on Signal Processing*, vol. 65, no. 17, pp. 4577-4592, Sep. 2017.
- P. Schniter and S. Rangan, "Compressive Phase Retrieval via Generalized Approximate Message Passing," *IEEE Trans. on Signal Processing*, vol. 63, no. 4, pp. 1043-1055, Feb. 2015.

## Description of SHyGAMP

- SHyGAMP approximates sum-product loopy belief propagation on factor graphs of the form:



- SHyGAMP iteratively passes messages back and forth between the  $p_{\mathbf{C}}$  and  $p_{\mathbf{y}|\mathbf{z}}$  nodes until convergence.
- Messages are approximated as  $K$ -dimensional Gaussian pdfs with diagonal covariance structure.
- This iterative message passing allows an  $NK$ -dimensional inference problem is broken into many  $K$ -dimensional inference problems.
- SHyGAMP's complexity for sketched clustering is  $\mathcal{O}(K(M+N)l + MNT)$ .
- The SHyGAMP algorithm can be divided into "linear" and "non-linear" steps.
- At each iteration the non-linear steps require computing the mean and covariance of the estimands using the following approximate posterior distributions:

$$p_{\mathbf{C}|\mathbf{r}}(\mathbf{c}_n | \hat{\mathbf{r}}_n; \mathbf{Q}_n^{\mathbf{r}}) \propto p_{\mathbf{C}}(\mathbf{c}_n) \mathcal{N}(\mathbf{c}_n; \hat{\mathbf{r}}_n, \mathbf{Q}_n^{\mathbf{r}}) \quad (9)$$

and

$$p_{\mathbf{z}|\mathbf{y}, \mathbf{p}}(\mathbf{z}_m | \mathbf{y}_m, \hat{\mathbf{p}}_m; \mathbf{Q}_m^{\mathbf{p}}) \propto p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}_m | \mathbf{z}_m) \mathcal{N}(\mathbf{z}_m; \hat{\mathbf{p}}_m, \mathbf{Q}_m^{\mathbf{p}}), \quad (10)$$

where the quantities  $\hat{\mathbf{p}}_m$ ,  $\mathbf{Q}_m^{\mathbf{p}}$ ,  $\hat{\mathbf{r}}_n$ , and  $\mathbf{Q}_n^{\mathbf{r}}$  are computed during the linear steps.

## The SHyGAMP Algorithm

**Require:** frequency matrix  $\mathbf{W}$ , sketch  $\mathbf{y}$ , pdfs  $p_{\mathbf{C}|\mathbf{r}}$  and  $p_{\mathbf{z}|\mathbf{y}, \mathbf{p}}$  from (9)-(10), initializations  $\hat{\mathbf{r}}_n(0)$ ,  $\mathbf{Q}_n^{\mathbf{r}}(0)$ .  
**Ensure:**  $t \leftarrow 0$ ;  $\hat{\mathbf{s}}_m(0) \leftarrow \mathbf{0}$ .

- repeat
- $\forall n: \hat{\mathbf{c}}_n(t) \leftarrow \mathbb{E}\{\mathbf{c}_n | \mathbf{r}_n = \hat{\mathbf{r}}_n(t-1); \mathbf{Q}_n^{\mathbf{r}}(t-1)\}$
- $\forall n: \mathbf{Q}_n^{\mathbf{c}}(t) \leftarrow \text{cov}\{\mathbf{c}_n | \mathbf{r}_n = \hat{\mathbf{r}}_n(t-1); \mathbf{Q}_n^{\mathbf{r}}(t-1)\}$
- $\forall m: \mathbf{Q}_m^{\mathbf{p}}(t) \leftarrow \sum_{n=1}^N W_{nm}^2 \mathbf{Q}_n^{\mathbf{c}}(t)$
- $\forall m: \hat{\mathbf{p}}_m(t) \leftarrow \sum_{n=1}^N W_{nm} \hat{\mathbf{c}}_n(t) - \mathbf{Q}_m^{\mathbf{p}}(t) \hat{\mathbf{s}}_m(t-1)$
- $\forall m: \hat{\mathbf{z}}_m(t) \leftarrow \mathbb{E}\{\mathbf{z}_m | \mathbf{y}_m, \mathbf{p}_m = \hat{\mathbf{p}}_m(t); \mathbf{Q}_m^{\mathbf{p}}(t)\}$
- $\forall m: \mathbf{Q}_m^{\mathbf{z}}(t) \leftarrow \text{cov}\{\mathbf{z}_m | \mathbf{y}_m, \mathbf{p}_m = \hat{\mathbf{p}}_m(t); \mathbf{Q}_m^{\mathbf{p}}(t)\}$
- $\forall m: \mathbf{Q}_m^{\mathbf{s}}(t) \leftarrow [\mathbf{Q}_m^{\mathbf{p}}(t)]^{-1} - [\mathbf{Q}_m^{\mathbf{p}}(t)]^{-1} \mathbf{Q}_m^{\mathbf{z}}(t) [\mathbf{Q}_m^{\mathbf{p}}(t)]^{-1}$
- $\forall m: \hat{\mathbf{s}}_m(t) \leftarrow [\mathbf{Q}_m^{\mathbf{p}}(t)]^{-1} (\hat{\mathbf{z}}_m(t) - \hat{\mathbf{p}}_m(t))$
- $\forall n: \mathbf{Q}_n^{\mathbf{r}}(t) \leftarrow [\sum_{m=1}^M W_{nm}^2 \mathbf{Q}_m^{\mathbf{s}}(t)]^{-1}$
- $\forall n: \hat{\mathbf{r}}_n(t) \leftarrow \hat{\mathbf{c}}_n(t) + \mathbf{Q}_n^{\mathbf{r}}(t) \sum_{m=1}^M W_{nm} \hat{\mathbf{s}}_m(t)$
- $t \leftarrow t + 1$
- until Terminated

## Computation of SHyGAMP Non-linear Steps

- The key technical challenge in applying SHyGAMP to sketched clustering is computing Lines 6-7 of the SHyGAMP algorithm when  $p_{\mathbf{y}|\mathbf{z}}$  has the form in (6).
- We have developed a method based on approximating  $p_{\mathbf{y}|\mathbf{z}}(\mathbf{y}_m | \mathbf{z}_m)$  with a Generalized von Mises distribution and evaluating the necessary integrals with the Laplace Approximation.

## Parameter Tuning

- Our Gaussian Mixture model in (4) requires properly selecting  $\alpha_k$  and  $\tau_{mk}$  in (6).
- Currently, we assume  $\tau_{mk}$  is invariant to  $m$ .
- Allowing  $\tau_{mk}$  to vary with  $m$  increases the generalizability of the model, but is more difficult to learn. Exploring this is one avenue for future work.
- One approach to tuning  $\alpha_k$  and  $\tau_k$  is via approximate EM:

$$\{\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\tau}}\} = \arg \max_{\boldsymbol{\alpha} \geq \mathbf{0}, \boldsymbol{\alpha}^T \mathbf{1} = 1, \boldsymbol{\tau}^w \geq \mathbf{0}} \sum_{m=1}^M \int_{\mathbb{R}^K} \mathcal{N}(\mathbf{z}_m; \hat{\mathbf{z}}_m, \mathbf{Q}_m^{\mathbf{z}}) \log p(\mathbf{y}_m | \mathbf{z}_m) d\mathbf{z}_m, \quad (11)$$

which can be optimized at every SHyGAMP iteration (immediately after Line 7) using gradient-projection.

- An alternate approach based on Bethe Free Energy Minimization [Schniter 15] is currently in development.

## Comparison Between SHyGAMP, CLOMPR++ and K-means++

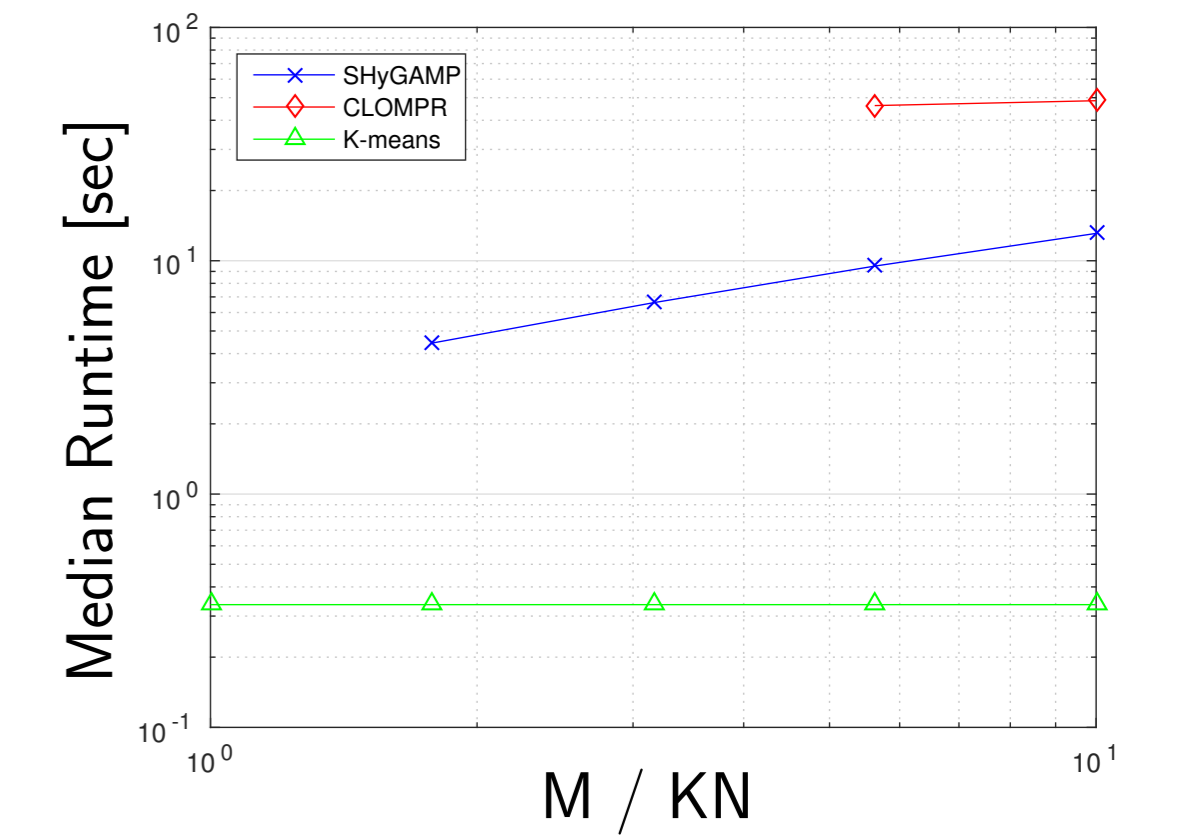
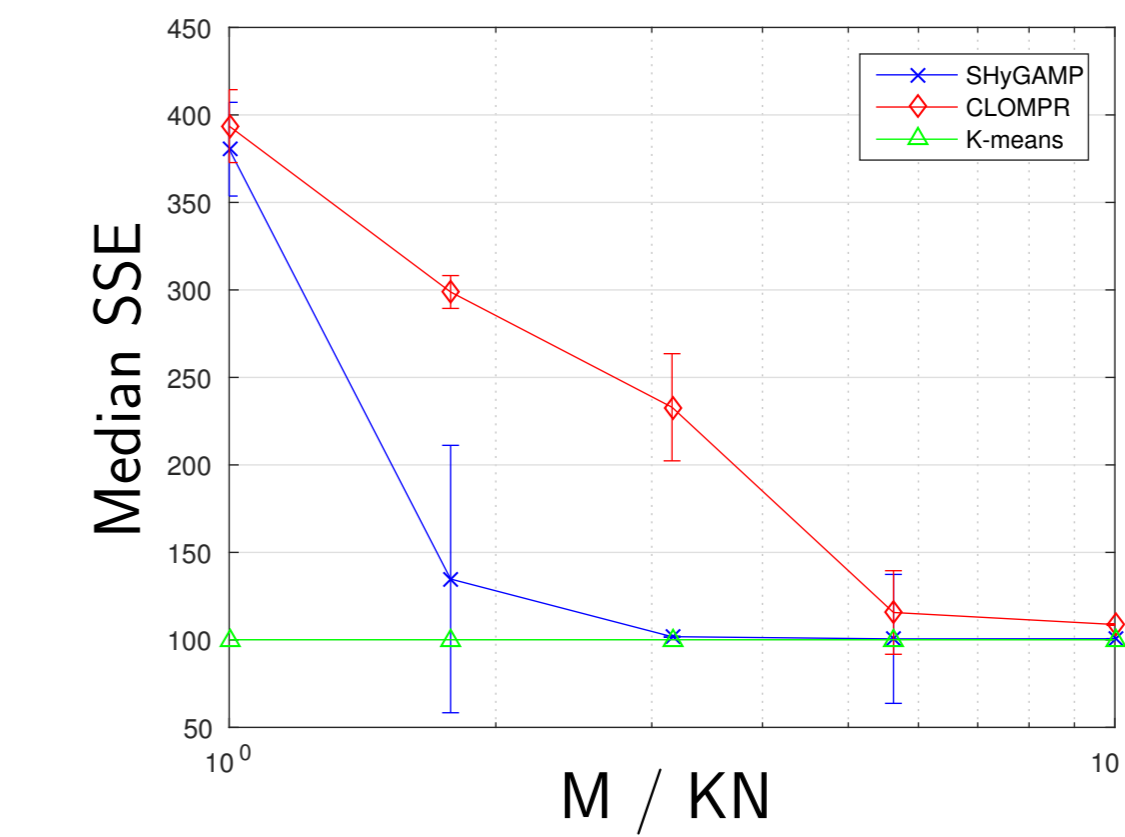
### Data generation model

- True  $\mathbf{x}_t \sim \sum_{k=1}^K \alpha_k \mathcal{N}(\mathbf{c}_k, \mathbf{I}_N)$  for  $\mathbf{c}_k \sim \mathcal{N}(\mathbf{0}_N, (1.5\sqrt{K})^2 \mathbf{I}_N)$ , and  $\alpha_k = \frac{1}{K} \forall k$ .

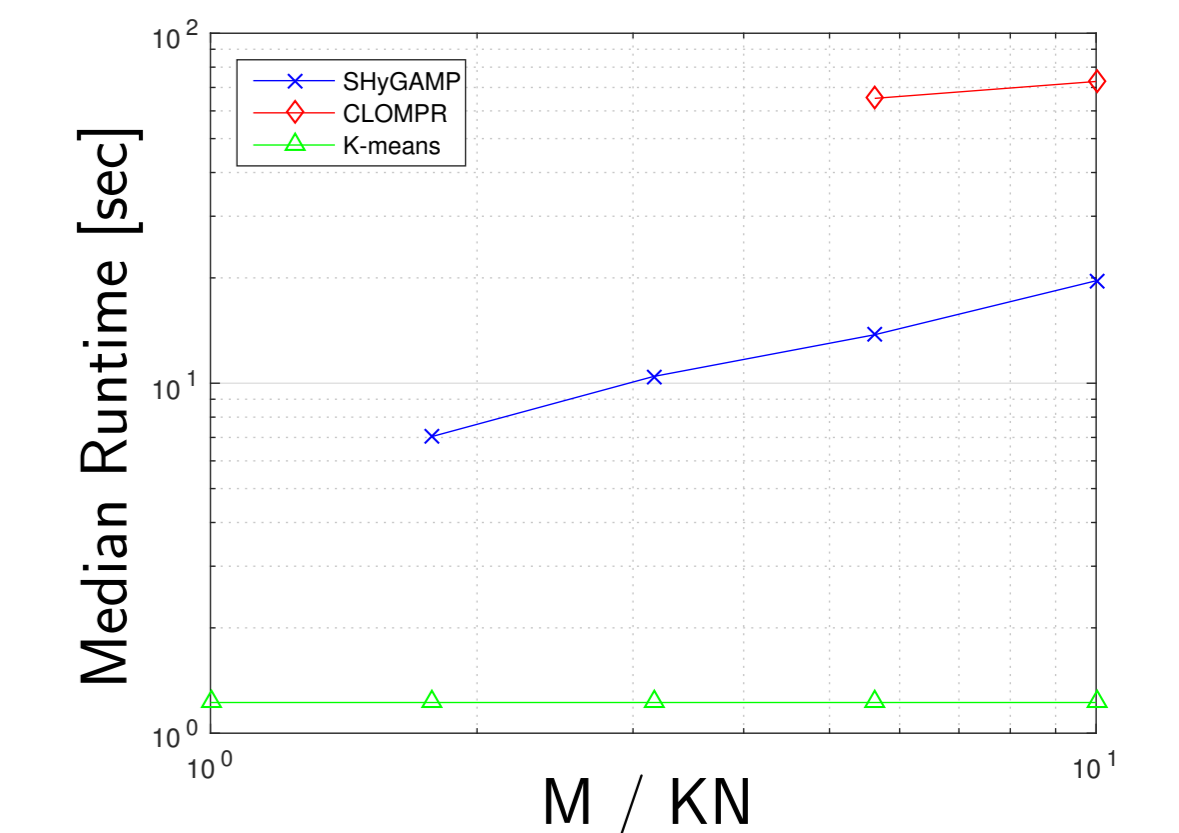
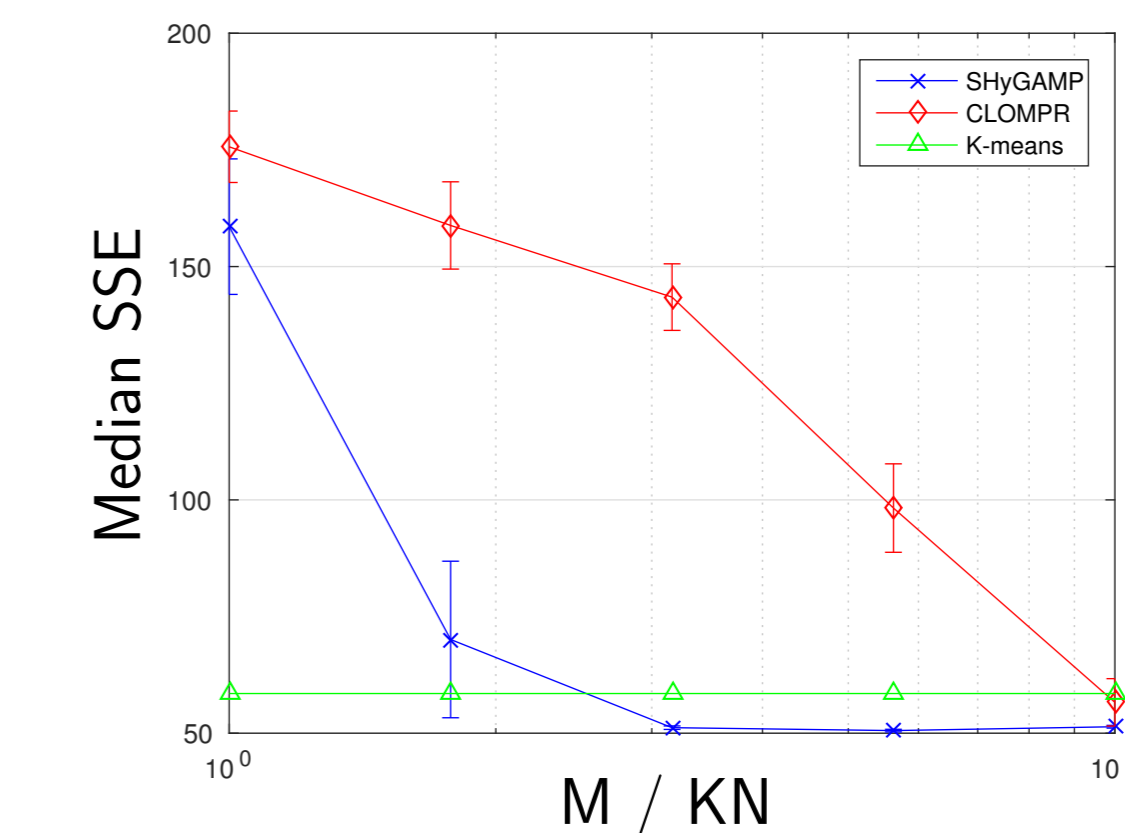
### Simulation: SSE vs $M$

- For each  $N \in \{50, 100\}$  and  $K \in \{5, 10\}$ , we tested several sketch lengths  $M \in [KN, 10KN]$ .
- We report the Median SSE and Median Runtime for SHyGAMP, CLOMPR++ and K-means++ over 10 trials. For SHyGAMP and CLOMPR++, we report runtime only when  $\text{SSE} < 2 \times \text{SSE}(\text{K-means++})$ .
- Compared to CLOMPR++, SHyGAMP has lower SSE and is faster at all tested  $M$ .

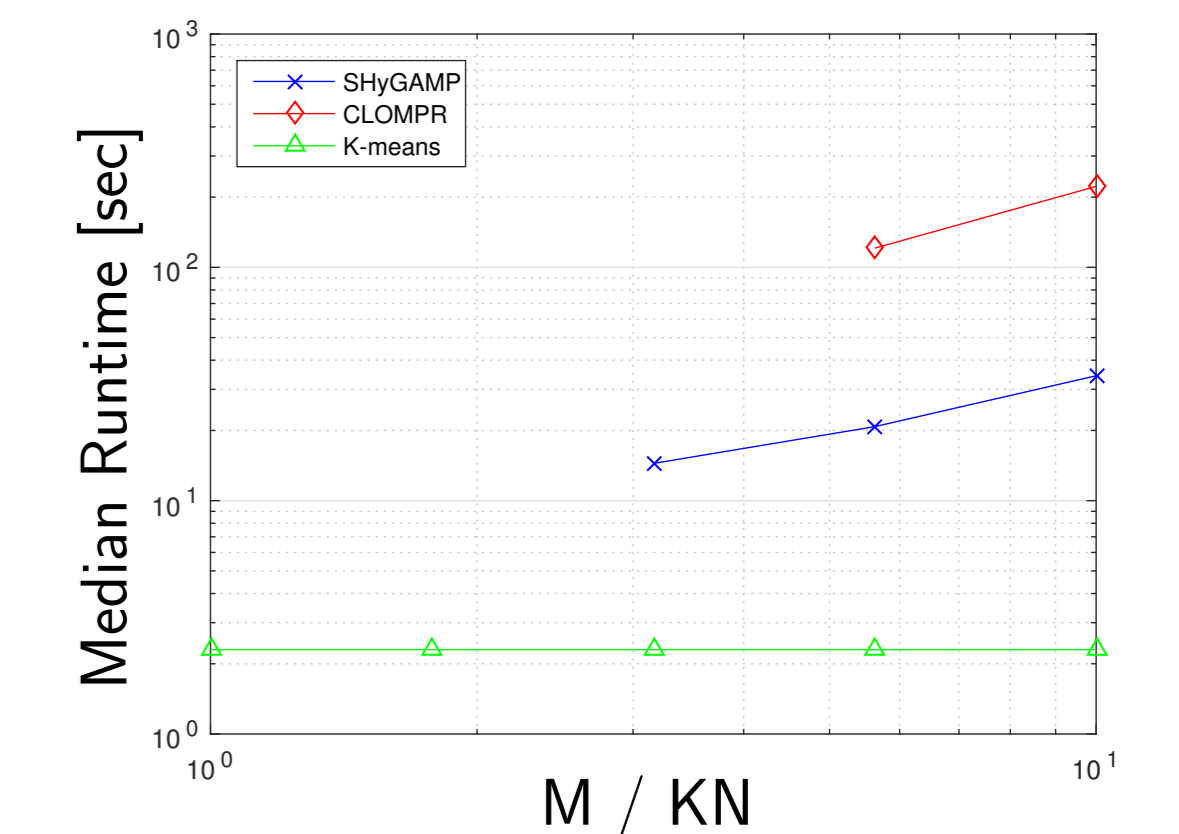
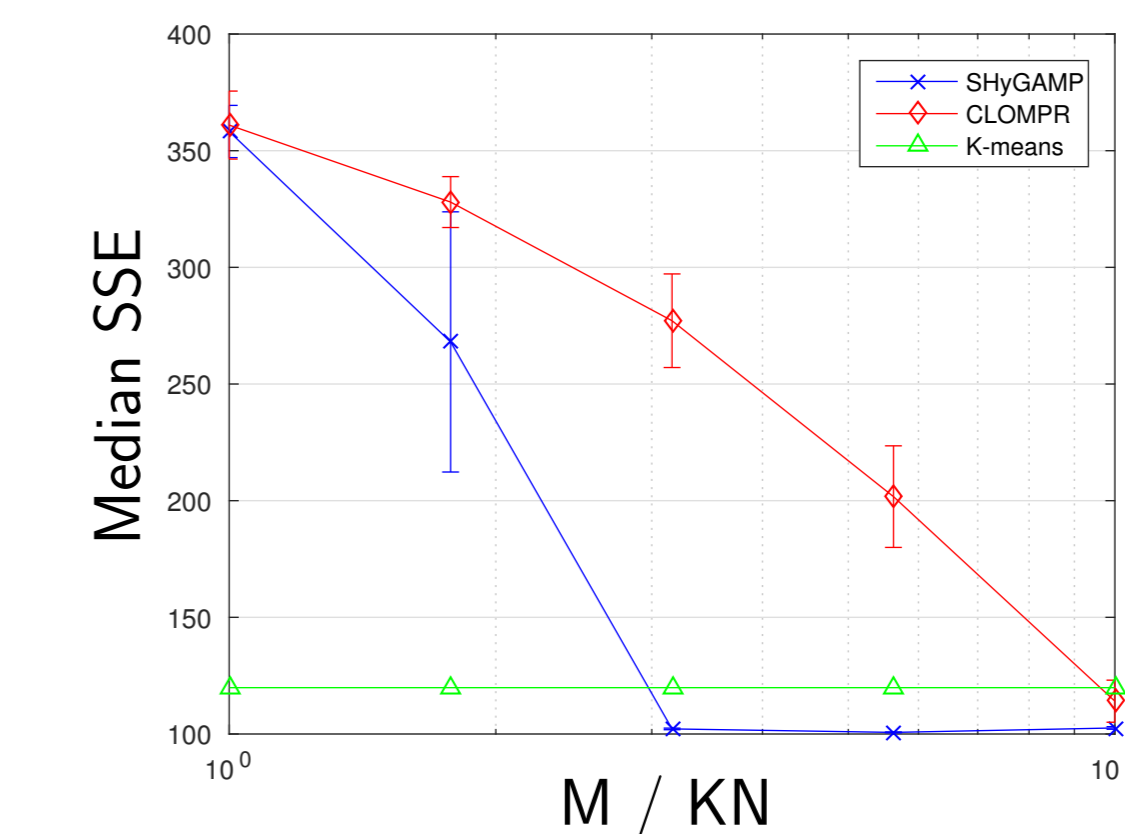
$K = 5, N = 100$ :



$K = 10, N = 50$ :



$K = 10, N = 100$ :



### Simulation: Classification Error vs Runtime

- Dimensions  $N = 20$  and  $K = 30$ . Training set with  $T = 10^4$  samples.
- Recovered cluster-centers used for classification on a test set with  $T = 5 \times 10^6$  samples.
- SHyGAMP and CLOMPR++ traces vary sketch size  $M$  logarithmically within  $[KN, 100KN]$ .
- K-means traces vary training subset size, in  $\{\frac{T}{20}, \frac{T}{25}, \dots, T\}$ , for a fixed # replicates in  $\{256, \dots, 4096\}$ .
- Results are the median of 5 trials (each trial used the same true centroids, but random train/test sets).
- SHyGAMP converged to the Bayes' Error Rate (BER) faster than K-means and CLOMPR++.

