

Nonparametric Coreset for Clustering

Rachit Chhaya
IIT Gandhinagar

Jayesh Choudhari
Univ. of Warwick

Anirban Dasgupta
IIT Gandhinagar

Supratim Shit
Technion

k - Clustering

Given a set \mathbf{A} of n points in \mathbb{R}^d , the solution of the problem is \mathbf{X}^* with k set of centres such that

$$\mathbf{X}^* = \arg \min_{\mathbf{X}} f_{\mathbf{X}}(\mathbf{A})$$

where, $f_{\mathbf{X}}(\mathbf{A}) = \sum_{\mathbf{a} \in \mathbf{A}} \min_{\mathbf{x} \in \mathbf{X}} \|\mathbf{x} - \mathbf{a}\|^2$

Coreset

Given $\mathbf{A} \in \mathbb{R}^{n \times d}$, a coreset \mathbf{C} is a weighted subsample of \mathbf{A} , such that $|\mathbf{C}| \ll |\mathbf{A}|$ and with high probability, $\forall \mathbf{X}$ with k centres.

$$\|f_{\mathbf{X}}(\mathbf{C}) - f_{\mathbf{X}}(\mathbf{A})\| \leq \epsilon f_{\mathbf{X}}(\mathbf{A})$$

A coreset \mathbf{C} is **non-parametric** if its size is independent of both k (#cluster), and it still ensures the above guarantee $\forall \mathbf{X}$ with at most n centres in \mathbb{R}^d .

Such a coreset does not exist!!

Non-Parametric Coreset

Given \mathbf{A} , there exists a non parametric coreset with small additive error.

$$\|f_{\mathbf{X}}(\mathbf{C}) - f_{\mathbf{X}}(\mathbf{A})\| \leq \epsilon(f_{\mathbf{X}}(\mathbf{A}) + f_{\varphi}(\mathbf{A}))$$

The additive factor depends on the data. Here, $\varphi = \text{mean}(\mathbf{A})$.

To show its existence we rely on importance sampling which is based on sensitivity framework along with barrier potential functions.

For streaming inputs we define **online sensitivity scores**, that depends on the points the algorithm have seen so far.

Importance Score

Barrier Potential based Sensitivity Function

$$\sup_{\mathbf{X}} \frac{f_{\mathbf{X}}(\mathbf{a}_i)}{(1 + \epsilon)f_{\mathbf{X}}(\mathbf{A}_{i-1}) - f_{\mathbf{X}}(\mathbf{C}_{i-1}) + \epsilon f_{\varphi}(\mathbf{A}_i)}$$

$$\sup_{\mathbf{X}} \frac{f_{\mathbf{X}}(\mathbf{a}_i)}{f_{\mathbf{X}}(\mathbf{C}_{i-1}) - (1 - \epsilon)f_{\mathbf{X}}(\mathbf{A}_{i-1}) + \epsilon f_{\varphi}(\mathbf{A}_i)}$$

Expected Upper Bound

$$\frac{2f_{\varphi}^{\mathbf{M}_i}(\mathbf{a}_i)}{\mu_i \epsilon \sum_{j \leq i} f_{\varphi}^{\mathbf{M}_j}(\mathbf{a}_j)} + \frac{12}{\mu_i \epsilon (i-1)}$$

Getting a true upper bound is challenging!!

NonParametricFilter Result

Sampling points in \mathbf{C} based on the above sensitivity scores ensures the following $\forall \mathbf{X}$ with at most n centres,

$$\|f_{\mathbf{X}}(\mathbf{C}) - f_{\mathbf{X}}(\mathbf{A})\| \leq \epsilon(f_{\mathbf{X}}(\mathbf{A}) + f_{\varphi}(\mathbf{A}))$$

Coreset Size: $O\left(\frac{\log n}{\mu \epsilon^2} (\log n + \log(f_{\varphi}^{\mathbf{M}}(\mathbf{A})) - \log(f_{\varphi}^{\mathbf{M}_2}(\mathbf{a}_2)))\right)$

Online Coreset with Deterministic Guarantee!!

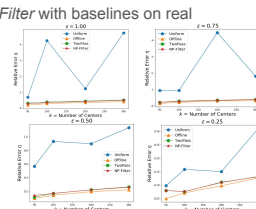
Experiments

Compare performance of *NonParametricFilter* with baselines on real world data.

After getting the coreset \mathbf{C}

- Run k-clustering on \mathbf{C} and \mathbf{A}
- Use these centres and report

$$\eta = \frac{|f_{\mathbf{X}_C}(\mathbf{A}) - f_{\mathbf{X}_A}(\mathbf{A})|}{f_{\mathbf{X}_A}(\mathbf{A})}$$



Algorithm 2 NonParametricFilter

Require: $\mathbf{a}_i, i = 1, \dots, n; \epsilon > 1; \epsilon \in (0, 1)$

Ensure: (\mathbf{C}, Ω)

$c^i = 2/\epsilon + 1; c^i = 2/\epsilon - 1; \varphi_0 = 0; S = 0; \mathbf{C}_0^i = \dots = \Omega_0^i = 0$

$\lambda = \|\mathbf{a}_1\|_{\min}; \nu = \|\mathbf{a}_1\|_{\max}$

while $i \leq n$ **do**

$\lambda = \min\{\lambda, \|\mathbf{a}_i\|_{\min}\}; \nu = \max\{\nu, \|\mathbf{a}_i\|_{\max}\}$

Update $\mathbf{M}_i; \mu_i = \lambda/\nu$

$\varphi_i = ((i-1)\varphi_{i-1} + \mathbf{a}_i)/i; S = S + f_{\varphi}^{\mathbf{M}_i}(\mathbf{a}_i)$

if $i = 1$ **then**

$p_i = 1$

else

$l_i^u = \frac{2f_{\varphi}^{\mathbf{M}_i}(\mathbf{a}_i)}{\epsilon \mu_i S} + \frac{12}{\epsilon \mu_i (i-1)}$

$l_i^l = \frac{2f_{\varphi}^{\mathbf{M}_i}(\mathbf{a}_i)}{\epsilon \mu_i S} + \frac{12}{\epsilon \mu_i (i-1)}$

$p_i = \min\{1, c^i l_i^u + c^i l_i^l\}$

end if

for $\forall j \in [i]$ **do**

$(c_j^i, \omega_j^i) = \begin{cases} (\mathbf{a}_i, 1/(p_i)) & \text{w. p. } p_i \\ (\mathbf{0}, 0) & \text{else} \end{cases}$

$(\mathbf{C}_i^j, \Omega_i^j) = (\mathbf{C}_{i-1}^j, \Omega_{i-1}^j) \cup (c_j^i, \omega_j^i)$

end for

$(\mathbf{C}, \Omega) = (\cup_{j \leq i} \mathbf{C}_{i-1}^j, \cup_{j \leq i} \Omega_{i-1}^j)$

end while

Return (\mathbf{C}, Ω)

Future Scope

- Computing the actual upper bound of the sensitivity scores.
- Coresets with deterministic guarantee for general loss functions.
- Showing lower bound of coresets for any loss functions.

References

- Lucic, Mario, Olivier Bachem, and Andreas Krause. "Strong coresets for hard and soft Bregman clustering with applications to exponential family mixtures." In Artificial Intelligence and statistics, pp. 1-9. PMLR, 2016.
- Bachem, Olivier, Mario Lucic, and Andreas Krause. "Scalable k-means clustering via lightweight coresets." Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, 2018.
- Batson, Joshua, Daniel A. Spielman, and Nikhil Srivastava. "Twice-ramanujan sparsifiers." SIAM Journal on Computing 41.6 (2012): 1704-1721.
- Feldman, Dan, and Michael Langberg. "A unified framework for approximating and clustering data." Proceedings of the forty-third annual ACM symposium on Theory of computing, 2011.
- Banerjee, Arindam, Srujana Merugu, Inderjit S. Dhillon, Joydeep Ghosh, and John Lafferty. "Clustering with Bregman divergences." Journal of machine learning research 6, no. 10 (2005).