Supplementary Materials for "SparseKmeans: Efficient K-means Clustering For Sparse Data"

Khoi Nguyen Pham Dang* MBZUAI khoinguyen.phamdang@mbzuai.ac.ae He-Zhe Lin* MBZUAI b07902028@csie.ntu.edu.tw Chih-Jen Lin National Taiwan Univ. / MBZUAI cjlin@csie.ntu.edu.tw

A STOPPING CONDITION OF K-MEANS ALGORITHM

SparseKmeans follows the implementation in scikit-learn to stop the iteration when the new centroids are close enough to their corresponding old ones. Specifically, the procedure is terminated when

$$\sum_{j=1}^{K} \|\boldsymbol{c}_{j} - \hat{\boldsymbol{c}}_{j}\|^{2} \le \tau \times \frac{1}{n} \sum_{k=1}^{n} \sigma_{k}^{2},$$
(A.1)

where c_j 's and \hat{c}_j 's are current and previous centroids, τ is a predefined parameter, and σ_k^2 is the variance of the *k*-th component in each x_i , respectively.

In SparseK means, we store the current centroids and previous centroids as two matrices C and \hat{C} :

$$C = \begin{bmatrix} - & c_1^T & - \\ - & c_2^T & - \\ & \vdots & \\ - & c_K^T & - \end{bmatrix} \text{ and } \hat{C} = \begin{bmatrix} - & \hat{c}_1^T & - \\ - & \hat{c}_2^T & - \\ & \vdots & \\ - & \hat{c}_K^T & - \end{bmatrix}.$$

To calculate $\|c_j - \hat{c}_j\|^2$ in (A.1), we first take the subtraction between *C* and \hat{C} :

$$C - \hat{C} = \begin{bmatrix} c_1^{I} - \hat{c}_1^{I} \\ c_2^{T} - \hat{c}_2^{T} \\ \vdots \\ c_K^{T} - \hat{c}_K^{T} \end{bmatrix},$$

and then calculate the squared norms of $C - \hat{C}$

$$\begin{bmatrix} \|\boldsymbol{c}_{1}^{T} - \hat{\boldsymbol{c}}_{1}^{T}\|^{2} \\ \|\boldsymbol{c}_{2}^{T} - \hat{\boldsymbol{c}}_{2}^{T}\|^{2} \\ \vdots \\ \|\boldsymbol{c}_{K}^{T} - \hat{\boldsymbol{c}}_{K}^{T}\|^{2} \end{bmatrix}$$
(A.2)

Summing all elements in (A.2), we have the left-hand side of (A.1).

B K-MEANS++ INITIALIZATION

The Initialization of c_j 's are often by choosing a subset of samples. In scikit-learn, they use K-means++ initialization [1] as its default setting. Compared to randomly choosing K points as the initial centroids, K-means++ reduces the number of iterations for convergence and minimizes the risk of poor clustering results.

K-means++ sequentially selects the initial points in a more optimal way by prioritizing points that are far from the previously chosen centroids. Therefore, the moment of choosing c_j , we already have c_1, \ldots, c_{j-1} . We then maintain a distance vector $d' \in \mathbb{R}^m$ with

$$d'_{i} = \min_{j'=1,...,j-1} ||\mathbf{x}_{i} - \mathbf{c}_{j'}||^{2},$$

which is the squared distance between x_i and the nearest point that has been chosen as a centroid. Then, based on d', we pick L candidates for c_j , where points further from the already chosen centroids are more likely to be selected. After selecting L candidates, we assign the one that yields the least total squared distance between data and centroids (called "potential") as c_j . The complete procedure of K-means++ is as follows.

- (1) Uniformly choose c_1 from $\{x_1, \ldots, x_m\}$ and initialize $d'_i = ||x_i c_1||^2$.
- (2) For j = 2, ..., K, do the following to select c_j .
 - (a) Sample $L = 2 + \log_2 K$ candidate points from $\{x_1, x_2, \dots, x_m\}$ using the distribution

$$\Pr(\mathbf{x}_i) = \frac{d'_i}{\sum_{i'=1}^m d'_{i'}}.$$

Denote $S = {s_1, \ldots, s_\ell}$ as the sampled points.

(b) Calculate the potential for assigning each $s \in S$ as c_j . Find the candidate point in *S* that minimizes the potential as c_j , i.e.,

$$c_j \in \arg\min_{s \in S} \sum_{i=1}^m \min\{d'_i, \|x_i - s\|^2\}.$$
 (B.1)

(c) Update

$$d'_i \leftarrow \min\{d'_i, \|\boldsymbol{x}_i - \boldsymbol{c}_j\|^2\}, i = 1, \dots, m.$$

C DETAILS OF ELKAN'S CLUSTER ASSIGNMENT

C.1 The Maintenance of $u(x_i)$ and $l(x_i, c_i)$

• For each x_i in cluster- ℓ_i , $u(x_i)$ is an upper bound of $||x_i - c_{\ell_i}||$. Since the centroids are updated after the cluster assignment, we need to update $u(x_i)$ correspondingly. Suppose \hat{c}_{ℓ_i} the centroid for cluster- ℓ_i before the centroid updates and $\hat{u}(x_i)$ is an upper bound of $||x_i - \hat{c}_{\ell_i}||$. To update $u(x_i)$ from $\hat{u}(x_i)$, by the triangle inequality, we have

$$\|x_i - \hat{c}_{\ell_i}\| + \|\hat{c}_{\ell_i} - c_{\ell_i}\| \ge \|x_i - c_{\ell_i}\|.$$

Therefore, after centroid updates, $u(x_i)$ can be obtained via

$$u(\mathbf{x}_i) \leftarrow \hat{u}(\mathbf{x}_i) + \|\mathbf{c}_{\ell_i} - \hat{\mathbf{c}}_{\ell_i}\|.$$
(C.1)

For each x_i and any j = 1,..., K, l(x_i, c_j) is a lower bound of ||x_i-c_j||. Similar to the situation in u(x_i), since c_j is changed afer the centroid updates, we need to update l(x_i, c_j) for all *i*. Suppose ĉ_j is the centroid of cluster-*j* before the centroid updates and l(x_i, ĉ_j) is a lower bound of ||x_i - ĉ_j||. By the triangle inequality and the definition of l(x_i, ĉ_j), we have

$$\|\mathbf{x}_i - \mathbf{c}_j\| \ge \|\mathbf{x}_i - \hat{\mathbf{c}}_j\| - \|\hat{\mathbf{c}}_j - \mathbf{c}_j\| \ge \hat{l}(\mathbf{x}_i, \hat{\mathbf{c}}_j) - \|\hat{\mathbf{c}}_j - \mathbf{c}_j\|.$$

^{*}Both authors contributed equally to this research.

Algorithm C.1: The original version of Elkan's cluster assignment in [2].

1 Compute $D \in \mathbb{R}^{K \times K}$ with $D_{i,i'} = \|\boldsymbol{c}_i - \boldsymbol{c}_{i'}\|/2$ 2 for i = 1, ..., m do if $u(\mathbf{x}) \leq \min_{j \neq \ell_i} D_{j,\ell_i}/2$ then 3 4 continue *bound tight* \leftarrow False 5 **for** j = 1, ..., K **do** 6 if $j \neq \ell_i$ and $u(\mathbf{x}_i) > D_{j,\ell_i}/2$ and $u(\mathbf{x}_i) > l(\mathbf{x}_i, \mathbf{c}_j)$ 7 then // filter clusters that x_i must not belong to if bound tight is False then 8 $d_i \leftarrow \|\mathbf{x}_i - \mathbf{c}_{\ell_i}\|$ 9 $u(\mathbf{x}_i) \leftarrow d_i, l(\mathbf{x}_i, \mathbf{c}_{\ell_i}) \leftarrow d_i$ 10 *bound tight* \leftarrow True 11 Go to Line 7 12 $l(\mathbf{x}_i, \mathbf{c}_j) \leftarrow \|\mathbf{x}_i - \mathbf{c}_j\|$ 13 if $d_i > ||x_i - c_j||$ then 14 $\ell_i \leftarrow j, u(x_i) \leftarrow ||x_i - c_j||$ 15

Therefore, after centroid updates, $l(x_i, c_j)$ can be updated by

$$l(\boldsymbol{x}_i, \boldsymbol{c}_j) \leftarrow \max \left\{ \hat{l}(\boldsymbol{x}_i, \hat{\boldsymbol{c}}_j) - \|\boldsymbol{c}_j - \hat{\boldsymbol{c}}_j\|, 0 \right\}.$$
(C.2)

C.2 Elkan's K-means Assignment in the Orignal Work [2]

In the main paper, we give the sketch of Elkan's K-means assignment in the original work of [2]. Here we provide a detailed version in Algorithm C.1. In particular, there are two additional settings:

(1) (Line 3 to Line 4) To avoid the for-loop over *j* in Algorithm 1, Elkan [2] uses considers the following condition

$$u(\mathbf{x}) \leq \frac{1}{2} \min_{j:j \neq \ell_i} \| \mathbf{c}_j - \mathbf{c}_{\ell_i} \|.$$
 (C.3)

If (C.3) holds, then condition (6) must hold for all $j \neq l_i$, which means x must remain in its original cluster.

(2) To further reduce the number of distance calculations, when the distance ||x_i - c_j|| may need to be calculated, we update u(x_i) to be the tighest upper bound ||x_i - c_{li}||. We illustrate the idea in the following figure:

$$\underbrace{\overbrace{1,\ldots,(j^*-1)}^{bound_loose}, j^*}_{\text{Line 8 invoked}}, (j^*+1),\ldots, K. \quad (C.4)$$

Suppose j^* is the first cluster satisfying neither (6) nor (7). For $j = 1, ..., j^* - 1$, we use $u(\mathbf{x}_i)$ to check condition (6) and (7). When $j = j^*$, we need both $||\mathbf{x}_i - \mathbf{c}_{\ell_i}||$ and $||\mathbf{x}_i - \mathbf{c}_{j^*}||$ to determine whether \mathbf{x}_i should stay in cluster- ℓ_i or go to cluster- j^* . Under this situation, we update $u(\mathbf{x}_i)$ to $||\mathbf{x}_i - \mathbf{c}_{\ell_i}||$, we go to Line 7 and recheck the conditions. Then, from $j = j^* + 1, ..., K$, we maintain $u(\mathbf{x}_i) = ||\mathbf{x}_i - \mathbf{c}_{\ell_i}||$ in Line 15 while ℓ_i may be changed. Since we use the smallest possible value of $u(\mathbf{x}_i)$, the number of distance calculations can be reduced.

D ADDITIONAL IMPLEMENTATION DETAILS FOR UPDATING CENTROIDS

In Section 3.1, we mentioned that when updating centroids, we need to maintain a vector d with $d_i = ||\mathbf{x}_i - \mathbf{c}_{\ell_i}||$, i.e., the distance to the nearest centroid. This allows us to replace the $u(\mathbf{x}_i)$ in (6) and (7) with d_i without additional costs. Here we explain why d is needed in centroid updates.

The main reason is, after cluster assignment, it is possible that none of x_i 's belongs to some cluster-j, so we cannot update c_j by (2). Suppose we have E empty clusters after the cluster assignment. Then scikit-learn would find E samples with the top-E largest d_i values, and reassign them to these empty clusters.

E DATASETS

This section describes the preprocessing method applied to the four data sets across all experiments. For binary classification data set Url^1 , we directly use the data instances for clustering. Other data sets including Amazon- $670K^2$, Wiki- $500K^3$ and Amazon- $3M^4$ are multi-label data sets, in which each instance can be associated with several labels. In the context of this work, to avoid confusion with clustering labels, we refer to the original labels of these mult-label data sets as meta-labels. Instead of clustering on data instances, we preprocessed the data to obtain a set of representations for meta-labels. These representations are then normalized and used as input for the K-means clustering algorithms. To clarify, let k be the total number of data instances, z_i be the TF-IDF representation of i^{th} meta-label. We call $Y \in \mathbb{R}^{k \times m}$ as a binary matrix, such that:

$$Y_{ij} = \begin{cases} 1 & \text{if } z_i \text{ is associated with the meta-label } j^{th} \\ 0 & \text{otherwise} \end{cases}$$
(E.1)

To obtain a meta-label representation, we compute the sum over all data instances corresponding to that meta-label.

$$\mathbf{x}_{j} = \frac{\sum_{i:Y_{ij}=1} z_{i}}{\|\sum_{i:Y_{ij}=1} z_{i}\|}.$$
 (E.2)

The process could be formalized as matrix multiplication:

$$X = Y^T Z \tag{E.3}$$

following by a normalization for each row x_i , in which $X \in \mathbb{R}^{m \times n}$ and $Z \in \mathbb{R}^{k \times n}$ are matrices of meta-label representations and instances representations, respectively.

REFERENCES

- David Arthur and Sergei Vassilvitskii. 2007. k-means++: the advantages of careful seeding. In Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms. 1027–1035.
- [2] Charles Elkan. 2003. Using the triangle inequality to accelerate k-means. In Proceedings of the Twentieth International Conference on International Conference on Machine Learning (ICML). 147–153.

¹https://www.csie.ntu.edu.tw/čjlin/libsvmtools/datasets/binary/url_combined_normalized.bz2
²https://www.csie.ntu.edu.tw/čjlin/libsvmtools/datasets/multilabel/Amazon-670K_tfdf_train_ver2.svm.bz2

³https://drive.google.com/open?id=1bGEcCagh8zaDV0ZNGsgF0QtwjcAm0Afk ⁴https://drive.google.com/open?id=187vt5vAkGI2mS2W0MZ2Qv48YKSJNbQv4