

An Empirical Evaluation of k -Means Coresets*

Abstract

Coresets are among the most popular paradigms for summarizing data. In particular, there exist many high performance coresets for clustering problems such as k -means in both theory and practice. Curiously, there exists little work on comparing the quality of available k -means coresets.

In this paper we perform such an evaluation. First, we show that it is computationally hard to compare the quality of not only two different coreset algorithms, but also of two different output of a (randomized) coreset algorithm. In order to perform an empirical evaluation, we therefore have to work with heuristics. To this end, we propose and analyse a benchmark for coreset comparison. Using this benchmark and real-world data sets, we conduct an exhaustive evaluation of the most commonly used coreset implementations.

1 Introduction

The design and analysis of scalable algorithms has become an important research area over the past two decades. This is particularly important in data analysis, where even polynomial running time might not be enough to handle proverbial *big data* sets. One of the main approaches to deal with the scalability issue is to compress or sketch large data sets into smaller, more manageable ones. The aim of such compression methods is to preserve the properties of the original data, up to some small error, while significantly reducing the number of data points.

Among the most popular and successful paradigms in this line of research are *coresets*. Informally, given a data set A , a coreset $S \subset A$ with respect to a given set of queries Q and query function $f : A \times Q \rightarrow \mathbb{R}_{\geq 0}$ approximates the behaviour of A for all queries up to some multiplicative distortion D via

$$\sup_{q \in Q} \max \left(\frac{f(S, q)}{f(A, q)}, \frac{f(A, q)}{f(S, q)} \right) \leq D.$$

Coresets have been applied to a number of problems such as computational geometry [2, 6], linear algebra

[24, 28], and machine learning [30, 33]. But the by far most intensively studied and arguably most successful applications of the coreset framework is the k -clustering problem.

Here we are given n points A with (potential unit) weights $w : A \rightarrow \mathbb{R}_{\geq 0}$ in some metric space with distance function dist and aim to find k centers C such that

$$\text{cost}_A(C) := \frac{1}{n} \sum_{p \in A} \min_{c \in C} w(p) \cdot \text{dist}^z(p, c)$$

is minimized. The most popular variant of this problem is probably the k -means problem in d -dimensional Euclidean space where $z = 2$ and $\text{dist}(x, y) = \sqrt{\sum_{i=1}^d (x_i - y_i)^2}$.

A (k, ε) -coreset is now a subset $\Omega \subset A$ with weights $w : \Omega \rightarrow \mathbb{R}_{\geq 0}$ such that for any set of k centers C

$$(1.1) \quad \sup_C \max \left(\frac{\text{cost}_A(C)}{\text{cost}_\Omega(C)}, \frac{\text{cost}_\Omega(C)}{\text{cost}_A(C)} \right) \leq 1 + \varepsilon.$$

The coreset definition in Eq. (1.1) provides an upper bound for the distortion of all candidate solutions i.e., all possible k clusterings. A *weak coreset* is a relaxed guarantee that holds for optimal or nearly optimal clusterings of A instead of all clusterings.

In a long line of work spanning the last 20 years [4, 5, 7, 14, 17, 21, 20, 23, 5, 26, 36], the size of coresets has been steadily improved with the current state of the art yielding a coreset with $\tilde{O}(k\varepsilon^{-2} \cdot \min(d, \varepsilon^{-2}))$ points for a distortion $D \leq (1 + \varepsilon)$ due to Cohen-Addad, Saulpic, and Schwiegelshohn [11]¹.

While we have a good grasp of the theoretical guarantees of these algorithms, our understanding of the empirical performance is somewhat lacking. There exist a number of coreset implementations, but it is usually difficult to assess which implementation summarizes the data best. To accurately evaluate a given coreset, we would need to come up with a k clustering C which results in a maximal distortion. Solving this problem is likely difficult: related questions such as deciding whether a 3-dimensional point set A is an ε -net of a

*The full version of the paper can be accessed at <https://arxiv.org/abs/1902.09310>

¹We use $\tilde{O}(x)$ to hide $\log^c x$ terms for any constant c .

net of a set B with respect to convex ranges is co-NP hard [?] and it is similarly co-NP hard to decide whether a point set A is a *weak cores*et of a point set B (see ?? in the appendix).

Due to this difficulty, a common heuristic for evaluating coresets is as follows [1, 18]. First, compute a coreset Ω with the available algorithm(s) using some input data A . Then, run an optimization algorithm on Ω to compute a k clustering C . The *best* coreset algorithm is considered to be the one which yields a clustering with the smallest cost.

The drawback of this evaluation method is that it conflates the two separate tasks of coreset construction and optimization. An algorithm may yield a good clustering (with small cost) yet fail to produce a high quality coreset (with small distortion). Additionally, this method is more likely to measure the performance of the underlying optimization problem, rather than evaluating the coresets themselves.

The purpose of this study is to systematically evaluate the quality of various coreset algorithms for k -means. As such, we develop a new evaluation procedure which estimates the distortion of coreset algorithms. On real-world data sets, we observe that while the evaluated coreset algorithms are generally able to find solutions with comparable costs, there is a stark difference in their distortions. This shows that differences between optimization and compression are readily observable in practice.

As a complement to our evaluation procedure on real-world data sets, we propose a benchmark framework for generating synthetic data sets. We argue why this benchmark has properties that results in hard instances for all known coreset constructions. We also show how to efficiently estimate the distortion of a candidate coreset on the benchmark.

2 Coreset Algorithms

Though the algorithms vary in details, coreset constructions come in one of the following two flavours:

1. Movement-based constructions: Such algorithms compute a clustering with T points such that $\text{cost}_A(T) \ll \text{OPT}$, where OPT is the cost of an optimal k -means clustering. The coreset guarantee then follows as a consequence of the triangle inequality. These algorithms all have an exponential dependency on the dimension d , and therefore have been overtaken by sampling-based methods. Nevertheless, these constructions are more robust to various constrained clustering formulations [22, 35]

and continue to be popular. Examples from theory include [19, 21].

2. Importance sampling: Points are sampled proportionate to their sensitivity which for a point p is defined as $\text{sens}(p) := \sup_C \frac{\min_{c \in C} \text{dist}^2(p, c)}{\text{cost}_A(C)}$ and weighted by their inverse sampling probability. In terms of theoretical performance, sensitivity sampling has largely replaced movement based constructions, see for example [15, 27].

Of course, there exist algorithms that draw on techniques from both, see for example [11]. In what follows, we will survey implementations of various coreset constructions that we will evaluate later.

StreamKM++ [1] The popular k -means++ algorithm [3] computes a set of centers K by iteratively sampling a point p in A proportionate to $\min_{q \in K} \text{dist}^2(p, q)$ and adding it to K . The procedure terminates once the desired number of centers has been reached. The first center is typically picked uniformly at random. The StreamKM++ paper runs the k -means++ algorithms for T iterations, where T is the desired coreset size. At the end, every point q in K is weighted by the number of points in A closest to it. While the construction has elements of important sampling, the analysis is largely movement-based. The provable bound required for the algorithm to compute a coreset is $O\left(\frac{k \log n}{\delta^{d/2} \varepsilon^d} \cdot \log^{d/2} \frac{k \log n}{\delta^{d/2} \varepsilon^d}\right)$.

BICO [18] Combines the very fast, but poor quality clustering algorithm BIRCH [40] with the movement-based analysis from [19, 21]. The clustering is organized by way of a hierarchical decomposition: When adding a point p to one of the coreset points Ω at level i , it first finds the closest point q in Ω . If p is too far away from q , a new center is opened at p . Otherwise p is either added to q , or, if adding p to q increases the clustering cost of q beyond a certain threshold, the algorithm attempts to add p to the child-clusters of q . The procedure then continues recursively. The provable bound required for the algorithm to compute a coreset is $O(k \varepsilon^{-d-2} \log n)$.

Sensitivity Sampling [14] The simplest implementation of sensitivity sampling first computes an $(O(1), O(1))$ bicriteria K approximation², for example by running k -means++ for $2k$ iterations [39]. Let K be the $2k$ clustering thus computed and let

²An (α, β) bicriteria approximation computes an α approximation using $\beta \cdot k$ many centers.

K_i be an arbitrary cluster of K with center q_i . Subsequently, the algorithm picks points proportionate to $\frac{\text{dist}^2(p, q)}{\text{cost}_{K_i}(\{q_i\})} + \frac{1}{|K_i|}$ and weighs any point by its inverse sampling probability. Let $|\hat{K}_i|$ be the estimated number of points in the sample. Finally, the algorithm weighs each q_i by $(1+\varepsilon) \cdot |K_i| - |\hat{K}_i|$. The provable bound required for the algorithm to compute a coreset is $\tilde{O}(k d \varepsilon^{-4})$ ([14]), $\tilde{O}(k \varepsilon^{-6})$ ([23]), or $\tilde{O}(k^2 \varepsilon^{-4})$ ([5]).

Group Sampling [11] First, the algorithm computes an $O(1)$ approximation (or a bicriteria approximation) K . Subsequently, the algorithm preprocesses the input into groups such that (1) for any two points $p, p' \in K_i$, their cost is identical up to constant factors and (2) for any two clusters K_i, K_j , their cost is identical up to constant factors. In every group, Group-Sampling now samples points proportionate to their cost. The authors of [11] show that there always exist a partitioning into $\log^2 1/\varepsilon$ groups. Points not contained in a group are snapped to their closest center q in K . q is weighted by the number of points snapped to it. The provable bound required for the algorithm to compute a coreset is $\tilde{O}(k \varepsilon^{-4})$ ([11]).

Ray Maker [?] The algorithm computes an initial solution with k centers which is a constant factor approximation of the optimal clustering. Around each center, $O(1/\varepsilon^{d-1})$ random rays are created which span the hyperplane. Next, each point $p \in A$ is snapped to its closest ray resulting in a set of one-dimensional points associated with each ray. Afterwards, a coreset is created for each ray by computing an optimal 1D clustering with k^2/ε^2 centers and weighing each center by the number of points in each cluster. The final coreset is composed of the coresets computed for all the rays. The provable bound required for the algorithm to compute a coreset is $O(k^3 \cdot \varepsilon^{-d-1})$. The algorithm has recently received some attention due to its applicability to the fair clustering problem [22].

2.1 Dimension Reduction Finally, we also combine coreset constructions with a variety of dimension reduction techniques. Since the seminal paper by Feldman, Schmidt, and Sohler [16], most coreset algorithms have used some form of dimension reduction to eliminate the dependency on d , either by explicitly computing a low-dimensional embedding, see for example [16, 37], or by using the existence of a suitable embedding in the analysis [11, 23].

In particular, movement-based coresets often have an

exponential dependency on the dimension, which can be alleviated with some form of dimension reduction, both in theory [35] and in practice [25]. Here there are two main techniques.

Principal Component Analysis: Feldman,

Schmidt, and Sohler [16] showed that projecting an input A onto the first $O(k/\varepsilon^2)$ principal components is a coreset, albeit in low dimension. The analysis was subsequently tightened by [8] and extended to other center based cost functions by [36]. Although its target dimension is generally worse than those based on random projections and terminal embeddings, there is nevertheless reasons for using PCA regardless: It removes noise and thus may make it easier to compute a high quality coreset.

Terminal Embeddings: Given a set of points A in \mathbb{R}^d , a terminal embedding $f : \mathbb{R}^d \rightarrow \mathbb{R}^m$ preserves the pairwise distance between any point $p \in A$ and any point $q \in \mathbb{R}^d$ up to a $(1 \pm \varepsilon)$ factor. The statement is related to the famous Johnson-Lindenstrauss lemma but it is stronger as it does not apply to only the pairwise distances of A . Nevertheless, the same target dimension is sufficient. Terminal embeddings were studied by [13, 29, 34], with Narayanan and Nelson [34] achieving an optimal target dimension of $O(\varepsilon^{-2} \log n)$, where n is the number of points. For applications to coresets, we refer to [4, 11, 23].

For an overview on practical aspects of dimension reduction, we refer to Venkatsubramanian and Wang [38]. In our evaluation, we focus on using PCA. The performance of the algorithms when applying a Johnson-Lindenstrauss transformation does not affect the behaviour of an algorithm that only depends on pairwise distances. We note that terminal embeddings, combined with an iterative application of the coreset construction from [5], can reduce the target dimension to a factor $\tilde{O}(\varepsilon^{-2} \log k)$. This is mainly of theoretical interest, as in practice the deciding factor wrt the target dimension is the precision, rather than dependencies on $\log n$ and $\log k$.

3 Benchmark Construction

In this section, we describe our benchmark. The benchmark has a parameter α which controls the number of points and dimensions. For a given value of k the benchmark consists of $n = k^\alpha$ points and $d = \alpha \cdot k$ dimensions. It is recursively constructed as follows.

Denote by $\mathbf{1}_k$ the k -dimensional all 1 vector and by

v_i^1 the k dimensional vector with entries $(v_i^1)_j = \begin{cases} -\frac{1}{k} & \text{if } i \neq j \\ \frac{k-1}{k} & \text{if } i = j \end{cases}$. For $\ell \leq \alpha$, recursively define the k^ℓ

dimensional vector $v_i^\ell = \begin{bmatrix} (v_i^{\ell-1})_{1 \cdot k} \\ (v_i^{\ell-1})_{2 \cdot k} \\ \vdots \\ (v_i^{\ell-1})_{1 \cdot k} \end{bmatrix}$. Finally, set the

column $t = a \cdot k + b$, $a \in \{0, \dots, \alpha - 1\}$ and $b \in \{1, \dots, k\}$, of A to be $k^{\alpha-a+1}$ stacks of v_b^{a+1} .

To get a better feel for the construction, we have given two small example benchmarks for $k = 2$ and $k = 3$ in the appendix (see ??).

3.1 Properties of the Benchmark We now summarize the key properties of the benchmark. To this end, we require a few notions. Let A be the input matrix. We slightly abuse notation and refer to A_i as both the i th point as well as the i th row of the matrix A . For a clustering $\mathcal{C} = \{C_1, \dots, C_k\}$, we define that the $n \times k$ indicator matrix \tilde{X} induced by \mathcal{C} via

$$\tilde{X}_{i,j} = \begin{cases} 1 & \text{if } A_i \in C_j \\ 0 & \text{else.} \end{cases}$$

Furthermore, we will also use the $n \times k$ normalized clustering matrix X defined as

$$X_{i,j} = \begin{cases} \frac{1}{\sqrt{|C_i|}} & \text{if } A_i \in C_j \\ 0 & \text{else.} \end{cases}$$

We also recall the following lemma which will allow us to express the k -means cost of a clustering \mathcal{C} with optimally chosen centers in terms of the cost of X and A .

LEMMA 3.1. (FOLKLORE) *Let A be an arbitrary set of points and let $\mu(A) = \frac{1}{|A|} \sum_{p \in A} p$ be the mean. Then for any point c*

$$\sum_{p \in A} \|p - c\|^2 = \sum_{p \in A} \|p - \mu(A)\|^2 + |A| \cdot \|\mu(A) - c\|^2.$$

This lemma proves that for any given cluster C_j , the mean is the optimal choice of center. We also note that any two distinct columns of X are orthogonal. Furthermore $\frac{1}{n} \mathbf{1} \mathbf{1}^T A$ copies the mean into every entry of A . Combining these two observations, we see that the matrix $XX^T A$ maps the i th row of A onto the mean of the cluster it is assigned to. Finally, define the Frobenius norm of an $n \times d$ A by $\|A\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^d A_{i,j}^2}$. Then the k -means cost of the clustering \mathcal{C} is precisely $\|A - XX^T A\|_F^2$.

We also require the following distance measure on clusterings as proposed by Meila [31, 32]. Given two clusterings \mathcal{C} and \mathcal{C}' , the $k \times k$ confusion matrix M is defined as $M_{i,j} = |C_i \cap C'_j|$. Furthermore for the indicator matrices \tilde{X} and \tilde{X}' induced by \mathcal{C} and \mathcal{C}' we have the identity $M = \tilde{X}^T \tilde{X}'$. Denote by Π_k the set of all permutations over k elements. Then the distance between \mathcal{C} and \mathcal{C}' is defined as

$$d(\mathcal{C}, \mathcal{C}') = 1 - \frac{1}{n} \max_{\pi \in \Pi_k} \sum_{i=1}^k M_{i, \pi(i)}.$$

Observe that for clusters that are identical, their distance is 0. The maximum distance between any two k clusterings is always $\frac{k-1}{k}$.

We are now ready to state the desired properties of our benchmark. The benchmark was designed to generate many clusterings such that

1. The distance between these clustering is $1 - \frac{1}{k}$, i.e. it is maximized.
2. The clusterings have equal cost and the centers in each clustering have equal cost.
3. The clusterings are induced by a set of centers in \mathbb{R}^d .

A formal verification of these properties can be found in the appendix. The first and second property ensure that (equally good) solutions we use for evaluation are spread out through the solution space. In other words, we are less likely to only focus on a set of solutions \mathcal{S} for which a low distortion on one $S \in \mathcal{S}$ implies a low distortion for all elements of \mathcal{S} . The third property is important as these are the only clusterings the (standard) coresets guarantee has to apply to.

The solutions we consider are given as follows. For the columns $a \cdot k + 1, \dots, (a+1) \cdot k$, we define the clustering $\mathcal{C}^a = \{C_1^a, \dots, C_k^a\}$ with $A_i \in C_j^a$ if and only if $A_{i,j} > 0$. Let \tilde{X}^a and X^a denote the indicator matrix and clustering matrix, respectively, as induced by \mathcal{C}^a .

3.2 Benchmark Evaluation We now describe how we use the benchmark to measure the distortion of a coresets. Assume that the coresets only consists of input points³.

³It is not necessary for coresets constructions in general to consists of input points. One can adjust the evaluation in to account for this. But since all algorithms considered in this paper have the property and it makes the evaluation simpler, we will only consider coresets that are subsets of the original point set.

Consider the clustering $\mathcal{C}^a = \{C_1^a, \dots, C_k^a\}$ for some a and let Ω with weights $w : \Omega \rightarrow \mathbb{R}_{\geq 0}$ be the coreset. Note that there are α many such clusterings, for each value of a . We use $w(C_i^a \cap \Omega) := \sum_{p \in C_i^a \cap \Omega} w(p)$ to denote the mass of points of C_i^a in Ω . For every cluster C_i^a with $w(C_i^a \cap \Omega) \geq |C_i^a|(1 - \varepsilon)$, we place a center at $\mu(C_i^a)$. Conversely, if $w(C_i^a \cap \Omega) < |C_i^a|(1 - \varepsilon)$, we do not place a center at $\mu(C_i^a)$. We call such clusters *deficient*. Let \mathcal{S} be the total number of thus placed centers.

We now compare the cost as computed on the coreset and the true cost of \mathcal{S} . Due to Lemma 3.1 and the fact that all clusters have equal cost, we may write for any deficient cluster C_i^a $\text{cost}_{C_i^a}(\mathcal{S}) = \text{cost}_{C_j^a}(\{\mu(C_j^a)\}) + k^{\alpha-1} \|\mu(C_j^a) - \mu(C_h^a)\|_2^2$, where C_h^a is a non-deficient cluster. Thus, the cost is

$$\text{cost}_{C_i^a}(\mathcal{S}) \approx (1 + \frac{2}{\alpha}) \cdot \text{cost}_{C_j^a}(\{\mu(C_j^a)\}).$$

Conversely, the cost on the coreset is

$$\text{cost}_{\Omega \cap C_i^a}(\mathcal{S}) \approx \frac{w(C_i^a \cap \Omega)}{\text{cost}_{C_j^a}(\{\mu(C_j^a)\})} (1 + \frac{2}{\alpha}) \cdot \text{cost}_{C_j^a}(\{\mu(C_j^a)\}).$$

Thus for each deficient clustering individually, the distortion will be close to $\frac{k^{\alpha-1}}{w(C_i^a \cap \Omega)} \geq \frac{1}{1-\varepsilon}$. If there are many deficient clusters, then this will also be the overall distortion. For all possible (suitably discretized) thresholds for deficiency, we can now identify the clustering \mathcal{C}^a with a maximum number of deficient clusters and use the aforementioned construction to get a lower bound on the distortion.

3.3 Further Extensions On the benchmark we considered here, both Sensitivity Sampling, as well as Group Sampling are similar to uniform sampling and, indeed, uniform sampling could be used to construct a good coreset. We can eliminate uniform sampling as a viable algorithm for this instance by combining multiple benchmarks B_1, \dots, B_t with $\sum_{i=1}^t k_i = k$. Each benchmark then has size $\sum_{i=1}^t k_i^\alpha$. We then add an additive offset to the coordinates of each benchmark such that they do not interfere. In this case, uniform sampling does not work if the values of the k_i are different enough. Since it is well known that uniform sampling is not a viable coreset algorithm in both theory and practice, we only used the basic benchmark for our evaluations.

4 Experiments

In this section, we present how we evaluated different algorithms. First, we propose our evaluation procedure which gauges the quality of coresets. Then, we describe the data sets used for the empirical evaluation and our

experimental setup. Finally, we detail the outcome of the experiments and our interpretation of the results.

4.1 Evaluation Procedure Accurately evaluating a k -means coreset of a real-world data set requires constructing a solution \mathcal{S} (a set of k centers) which results in a maximal distortion. Finding such a solution, however, is difficult. Instead, we can estimate the quality of a given coreset by finding meaningful candidate solutions.

One approach is to use k -means++ as follows: compute a coreset Ω on a real-world data set A . Then, run k -means++ on Ω to find a set of k centers. Repeat the k -means++ algorithm 5 times and pick \mathcal{S} to be the set of centers with the largest distortion. The main advantage of this approach is that k -means++ can uncover natural cluster structures in the data.

Another approach is to randomly generate a candidate solution \mathcal{S} . For example, one could generate k random points inside the minimum enclosing ball (MEB) of a coreset Ω . This can be repeated 5 times. Then, a candidate solution is the set of points with the largest distortion. While it is very fast to generate candidate solutions in this manner, this method has its drawbacks. It is not readily apparent how to define a distribution of meaningful solutions from which to sample. Moreover, a randomly drawn solution, which does not exploit the behavior of a coreset construction, is less likely to yield a worst-case candidate solution. Nevertheless, we apply both k -means and random sampling inside the MEB to generate candidate solutions in our evaluations.

Granted the usefulness of evaluating coresets on real-world data sets, it can be tricky to gauge the general performance of coreset algorithms using only a small selection of data sets. For this reason, we used our benchmark to complement the evaluation on real-world data sets. The benchmark accomplishes two important tasks. First, the benchmark allows us to quickly find a bad solution because both good and bad clusterings are known a priori. It is unclear how to find bad clusterings for real-world data sets. Second, it is easier to make a fair comparison of different coreset constructions because the benchmark is known to generate hard instances for all known coreset algorithms. This cannot be said for real-world data sets. For the benchmark, we computed the distortion following the evaluation procedure described in Section 3.

Every randomized coreset construction was repeated 10 times. We aggregated the reported distortions by taking the average over all 10 evaluations. In addition, we preprocessed the data using PCA (compare Section 2.1).

4.2 Data sets We conducted experiments on five real-world data sets and four instances of our benchmark. Benchmark instances were generated to match approximately the sizes of the real-world data sets. The sizes of the data sets are summarized in Table 1. We now provide a brief description of each of the real-world data sets.

The *Census*⁴ dataset is a small subset of the Public Use Microdata Samples from 1990 US census. It consists of demographic information encoded as 68 categorical attributes of 2,458,285 individuals.

*Covertypes*⁵ is comprised of cartographic descriptions and forest cover type of four wilderness areas in the Roosevelt National Forest of Northern Colorado in the US. It consists of 581,012 records, 54 cartographic variables and one class variable. Although *Covertypes* was originally made for classification tasks, it is often used for clustering tasks by removing the class variable [1].

The data set with the fewest number of dimensions is *Tower*⁶. This data set consists of 4,915,200 rows and 3 features as it is a 2,560 by 1,920 picture of a tower on a hill where each pixel is represented by a RGB color value.

Inspired by [18], *Caltech* was created by computing SIFT features from the images in the Caltech101⁷ image database. This database contains pictures of objects partitioned into 101 categories. Disregarding the categories, we concatenated the 128-dimensional SIFT vectors from each image into one large data matrix with 3,680,458 rows and 128 columns.

*NYTimes*⁸ is a dataset composed of the bag-of-words (BOW) representations of 300,000 news articles from The New York Times. The vocabulary size of the text collection is 102,660. Due to the BOW encoding, *NYTimes* has a very large number of dimensions and is highly sparse. To make processing feasible, we reduced the number of dimensions to 100 using terminal embeddings.

4.3 Preprocessing & Experimental Setup To understand how denoising effects the quality of the outputted coresets, we applied Principal Component Analysis (PCA) on *Caltech*, *Census* and *Covertypes* by

⁴[https://archive.ics.uci.edu/ml/datasets/US+Census+Data+\(1990\)](https://archive.ics.uci.edu/ml/datasets/US+Census+Data+(1990))

⁵<https://archive.ics.uci.edu/ml/datasets/covertypes>

⁶<http://homepages.uni-paderborn.de/frahling/coremeans.html>

⁷http://www.vision.caltech.edu/Image_Datasets/Caltech101/

⁸<https://archive.ics.uci.edu/ml/datasets/bag+of+words>

	Data points	Dimensions
<i>Caltech</i>	3,680,458	128
<i>Census</i>	2,458,285	68
<i>Covertypes</i>	581,012	54
<i>NYTimes</i>	500,000	102,660
<i>Tower</i>	4,915,200	3

k	α	Data points	Dimensions
10	6	1,000,000	60
20	5	3,200,000	100
30	4	810,000	120
40	4	2,560,000	160

Table 1: Size of real world data sets (left) and evaluated benchmarks (right).

using the k singular vectors corresponding to the largest singular values. For these three data sets, we preserved the dimensions of the original data. The *NYTimes* dataset did not permit the preservation of dimensions as the number of dimensions is very large. In this case, we used PCA to reduce the dimensions to k . We did not perform any preprocessing on *Tower* due to its low dimensionality.

We followed the same experimental procedure with respect to the choice of parameter values for the algorithms as prior works [1, 18]. For the target coreset size, we used $200k$ for all our experiments. On *Caltech*, *Census*, *Covertypes* and *NYTimes*, we used k values in $\{10, 20, 30, 40, 50\}$, while for *Tower* we used larger cluster sizes $k \in \{20, 40, 60, 80, 100\}$. On the benchmark instances, we settled on $k \in \{10, 20, 30, 40\}$ as a reasonable trade-off between running time and data set size.

We implemented Sensitivity Sampling, Group Sampling, Ray Maker, and StreamKM++ in C++. The source code can be found on GitHub⁹. For BICO, we used the authors' reference implementation¹⁰. The source code was compiled with gcc 9.3.0. The experiments were performed on a machine with 14 cores (3.3 GHz) and 256 GB of memory.

4.4 Outcome of Experiments We summarized the distortions in Fig. 1, for numerical variance we refer to ?? in the appendix. All five algorithms are matched on the *Tower* dataset. The worst distortions across the algorithms are close to 1, and performance between

⁹Link to repository will be provided later.

¹⁰<https://ls2-www.cs.tu-dortmund.de/grav/en/bico>

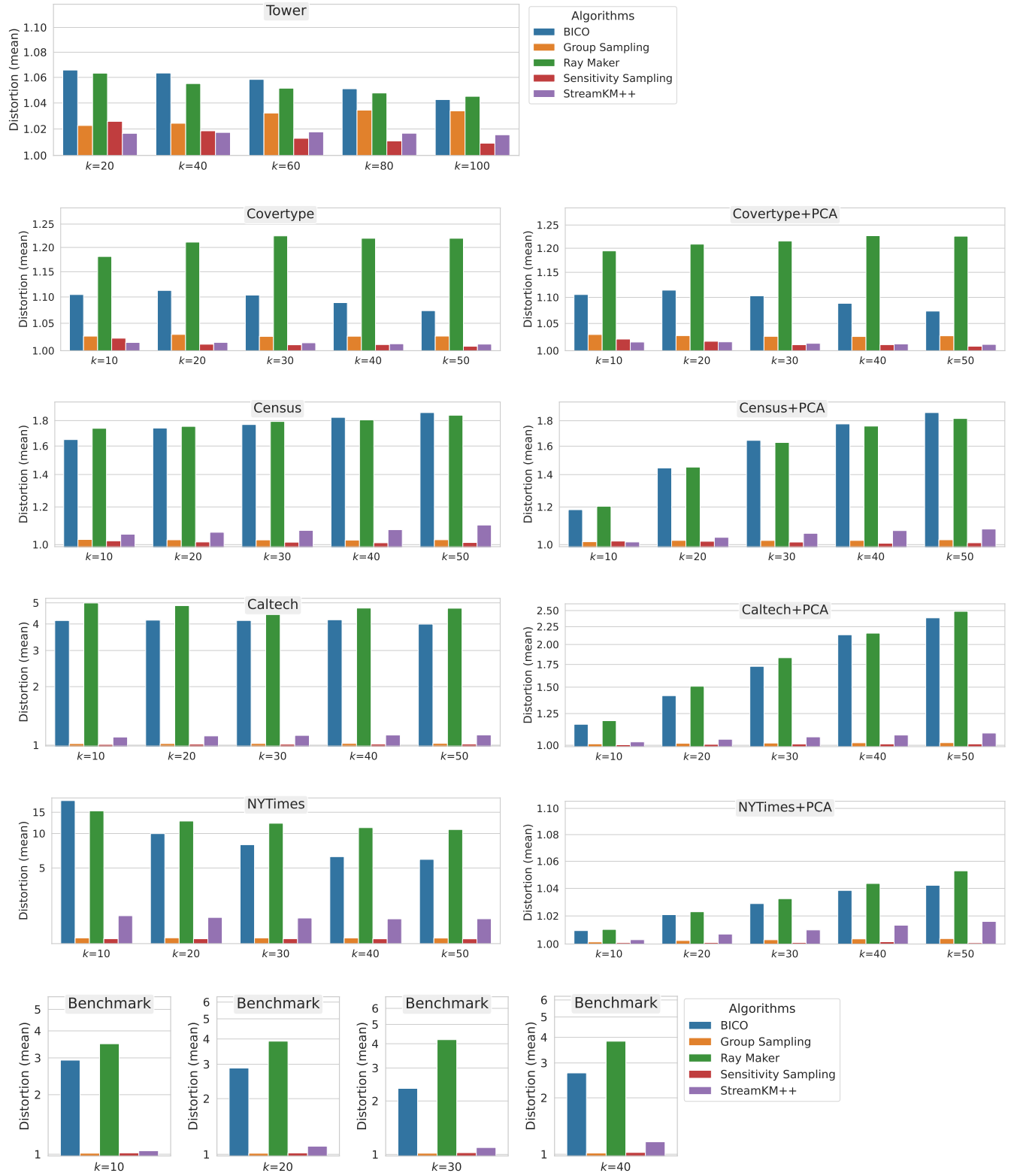


Figure 1: The distortions of the evaluated coreset algorithms on five real-world data sets and on four benchmark instances. The axis is non-linear as otherwise the bars for Sensitivity Sampling and Group Sampling would disappear on the plots. Their distortions are very close to 1.

the algorithms is negligible. The performance difference between sampling-based and movement-based methods become more pronounced as the number of dimensions increase. On *Coverttype* with its 54 features, Ray Maker performs the worst followed by BICO and Group Sampling while Sensitivity Sampling and StreamKM++ perform the best. Differences in performance are more noticeable on *Census*, *Caltech*, and *NYTimes* where methods based on importance sampling perform much better. Sensitivity Sampling and Group Sampling perform the best, StreamKM++ come in second while BICO and Ray Maker perform the worst across these data sets. On the *Benchmark*, Ray Maker is the worst while Sensitivity Sampling and Group Sampling are the best. StreamKM++ performs also very well compared to BICO.

Reducing noise with PCA helped boost the performance on real-world data sets with large number of dimensions. On *Coverttype*, PCA does not change the performance numbers by much. On *Census* ($d = 68$), the preprocessing step with PCA improves the performance of BICO and Ray Maker slightly for lower values of k . The distortions of BICO and Ray Maker are reduced markedly on *Caltech* ($d = 128$) after applying PCA.

4.5 Interpretation of Experimental Results

Optimization versus Compression While all five algorithms are equally matched when optimizing on the candidate coresets, coreset quality performance differ significantly (see Fig. 1). We omit tables and plots detailing the cost due to space restrictions. For all data sets, the obtained costs differed insignificantly for all values of k .

The quality of the coreset itself can be closely tied to the change in cost with increasing number of centers. It is not uncommon for the k -means cost of real-world data sets to drop significantly for larger values of k . Fig. 2 illustrates this behavior for several real-world data sets. The more the curve bends, the less of a difference there is between computing a coreset and a clustering with low cost. For data sets with a L-shaped cost curve, a coreset algorithm adding more centers to the coreset will seem to be performing well when evaluating it based on the outcome of the optimization. *Tower* is a good example of a data set where optimization is very close to compression. Its cost curve bends the most which indicates that adding more centers help reduce the cost. One of the strengths of the benchmark is that there is no way of reducing the cost without capturing the right subclusters within a benchmark instance. This means that the cost does not decrease markedly beyond

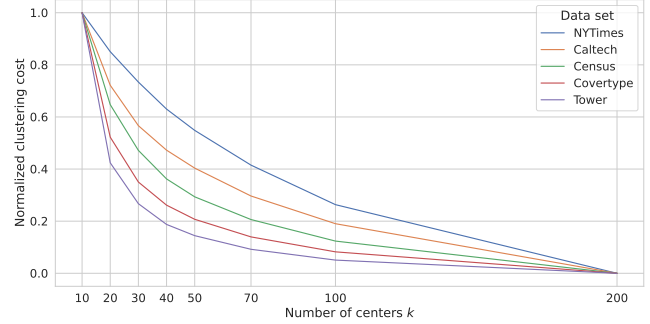


Figure 2: Depicts how clustering costs of five real-world data sets decrease as the number of centers increase. The most widely used data sets for evaluating coresets are *Tower*, *Coverttype*, and *Census*, while *Caltech* is rarely used and *NYTimes* has never been used before. Plotting the cost curve allows us to study whether we can observe a difference between coreset construction and optimization in a data set when evaluating a coreset based on cost.

a certain value of k even if more centers are added, see ?? in the appendix.

For BICO, Ray Maker, and StreamKM++, there is a correlation between the steepness of the cost curve for a data set and the distortion of the generated coreset. On data sets where the curve is less steep, we observed higher distortions. The effect is more pronounced for the movement-based algorithms (BICO and Ray Maker) than for StreamKM++. The other two sampling-based approaches (Group Sampling and Sensitivity Sampling) seem to be free from this behavior as they consistently generate high quality coresets irrespective of the shape of cost curve.

Movement-based versus Sampling-based Approaches In general, movement-based constructions perform the worst in terms of coreset quality. We observed that BICO and Ray Maker have the highest distortions across all data sets including on the benchmark instances. Among the sampling-based algorithms, Sensitivity Sampling performs well with Group Sampling generally being competitive. This runs contrary to theory where Group Sampling has the better (currently known) theoretical bounds. StreamKM++ is an interesting case. Like the movement-based methods, its distortion increases with the dimension. Nevertheless, it generally performs significantly better than BICO and Ray Maker. This can be attributed to the fact that the coreset produced by StreamKM++ consists entirely of k -means++ centers weighted by the number of points of a minimal cost assignment. This is similar to movement-

based algorithms such as BICO. Nevertheless, it also retains some of the performance from pure importance schemes.

In practice as well as in theory, the distortion of movement-based algorithms are affected by the dimension. By comparison, sampling-based algorithms are affected very little. Theoretically, there should not exist a difference, as the sampling bounds are independent of the dimension. What little effect can be observed is likely due to PCA making it easier to find low cost solutions that form the backbone of all coresets constructions. StreamKM++ is an interesting case, as it is still affected by the dimension, though less than the other movement based methods, due its performance without dimension reduction being significantly better than the worst-case bounds would suggest.

Impact of PCA On almost all our data sets, the performance improves when input data is preprocessed with PCA, especially for the movement-based algorithms. Empirically, the more noise is removed (i.e., small k value), the lower the distortion. Notice that k is the number of principal components that the input data is projected on to. The rest of the low variance components are treated as noise and removed. Method utilizing sampling (Group Sampling, Sensitivity Sampling and StreamKM++) are less effected by the preprocessing step. On *Covertime*, PCA does not change the distortions by much because almost all the variance in the data is explained by the first five principal components (see ??). On *Caltech* and *NYTimes*, the quality of the coresets by BICO and Ray Maker improves greatly because the noise removal is more aggressive. Even if the quality is much better for movement-based coresets constructions due to PCA, importance sampling methods are still superior when it comes to the quality of the compression. Summarizing, all methods benefit from PCA, and in case of movement based constructions, we consider PCA a necessary preprocessing step. For the sampling based methods, the computational expense of using PCA in preprocessing does not seem justify the comparatively meagre gains in coreset distortion.

5 Conclusion

In this work, we studied how to assess the quality of k -means coresets computed by state-of-the-art algorithms. Previous work generally measured the quality of optimization algorithms run on the coreset, which we empirically observed to be poor indicator of coreset quality. For real data sets, we sampled candidate clusterings and evaluated the worst case distortion on them. Complementing this, we also proposed a benchmark framework

which generates hard instances for known k -means coreset algorithms. Our experiments indicate a general advantage for algorithms based on importance sampling over movement-based methods aiming towards computing low-cost clusterings. Despite movement based methods running on very efficient code, it is necessary to compliment them with rather expensive dimension reduction methods, rendering what efficiency they might have over importance sampling somewhat moot.

Two results bear further investigation. First, the currently known provable coreset sizes for Sensitivity Sampling are worse than those provable via Group Sampling. Empirically, we observed the opposite: While Group Sampling is competitive, Sensitivity Sampling always outperforms it. Since Group Sampling requires somewhat cumbersome computational overhead, practical applications should prefer Sensitivity Sampling. In light of these results, a theoretical analysis for Sensitivity Sampling matching the performance of Group Sampling would be welcome.

The second point of interest focusses on the performance of StreamKM++. The distortion of this algorithm is significantly better than what one would expect from its theoretical analysis. Empirically, StreamKM++ is notably better than the other movement based constructions across all data sets, and especially on high dimensional data. While it is not competitive to the pure importance sampling algorithms, there are several reasons for investigating it further. It essentially only requires running k -means++ for additional iterations, which is already a nearly ubiquitous algorithm for the k -means problem. Although the other sampling-based coreset algorithms can also be readily implemented, doing so might be cumbersome. In particular, the theoretically (but not empirically) best algorithm Group Sampling requires extensive preprocessing steps. This begs the question whether there exist a better theoretical analysis for StreamKM++.

In addition, StreamKM++ currently weighs each point by the number of points assigned to it. It may also be possible to improve the performance of the algorithm in both theory and practise by using a different weighting scheme. We leave this as an open problem for future research.

References

- [1] Marcel R. Ackermann, Marcus Märtens, Christoph Raupach, Kamil Swierkot, Christiane Lammersen, and Christian Sohler. Streamkm++: A clustering algorithm for data streams. *ACM Journal of Experimental Algorithmics*, 17(1), 2012.

- [2] Pankaj K. Agarwal, Sariel Har-Peled, and Kasturi R. Varadarajan. Geometric approximation via coresets. In *Combinatorial and computational geometry, MSRI*, pages 1–30. University Press, 2005.
- [3] David Arthur and Sergei Vassilvitskii. k-means++: the advantages of careful seeding. In *Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2007, New Orleans, Louisiana, USA, January 7-9, 2007*, pages 1027–1035, 2007.
- [4] Luca Becchetti, Marc Bury, Vincent Cohen-Addad, Fabrizio Grandoni, and Chris Schwiegelshohn. Oblivious dimension reduction for k-means: beyond subspaces and the johnson-lindenstrauss lemma. In *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing, STOC 2019, Phoenix, AZ, USA, June 23-26, 2019*, pages 1039–1050, 2019.
- [5] Vladimir Braverman, Shaofeng H.-C. Jiang, Robert Krauthgamer, and Xuan Wu. Coresets for clustering in excluded-minor graphs and beyond. In Dániel Marx, editor, *Proceedings of the 2021 ACM-SIAM Symposium on Discrete Algorithms, SODA 2021, Virtual Conference, January 10 - 13, 2021*, pages 2679–2696. SIAM, 2021.
- [6] Timothy M. Chan. Dynamic coresets. *Discret. Comput. Geom.*, 42(3):469–488, 2009.
- [7] Ke Chen. On coresets for k-median and k-means clustering in metric and Euclidean spaces and their applications. *SIAM J. Comput.*, 39(3):923–947, 2009.
- [8] Michael B. Cohen, Sam Elder, Cameron Musco, Christopher Musco, and Madalina Persu. Dimensionality reduction for k-means clustering and low rank approximation. In *Proceedings of the Forty-Seventh Annual ACM on Symposium on Theory of Computing, STOC 2015, Portland, OR, USA, June 14-17, 2015*, pages 163–172, 2015.
- [9] Vincent Cohen-Addad and Karthik C. S. Inapproximability of clustering in lp metrics. In David Zuckerman, editor, *60th IEEE Annual Symposium on Foundations of Computer Science, FOCS 2019, Baltimore, Maryland, USA, November 9-12, 2019*, pages 519–539. IEEE Computer Society, 2019.
- [10] Vincent Cohen-Addad, Karthik C. S., and Euiwoong Lee. On approximability of clustering problems without candidate centers. In Dániel Marx, editor, *Proceedings of the 2021 ACM-SIAM Symposium on Discrete Algorithms, SODA 2021, Virtual Conference, January 10 - 13, 2021*, pages 2635–2648. SIAM, 2021.
- [11] Vincent Cohen-Addad, David Saulpic, and Chris Schwiegelshohn. A new coreset framework for clustering. In Samir Khuller and Virginia Vassilevska Williams, editors, *STOC '21: 53rd Annual ACM SIGACT Symposium on Theory of Computing, Virtual Event, Italy, June 21-25, 2021*, pages 169–182. ACM, 2021.
- [12] Vincent Cohen-Addad, David Saulpic, and Chris Schwiegelshohn. A new coreset framework for clustering. In *Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing*, pages 169–182, 2021.
- [13] Michael Elkin, Arnold Filtser, and Ofer Neiman. Terminal embeddings. *Theor. Comput. Sci.*, 697:1–36, 2017.
- [14] Dan Feldman and Michael Langberg. A unified framework for approximating and clustering data. In *Proceedings of the 43rd ACM Symposium on Theory of Computing, STOC 2011, San Jose, CA, USA, 6-8 June 2011*, pages 569–578, 2011.
- [15] Dan Feldman and Michael Langberg. A unified framework for approximating and clustering data. In *Proceedings of the 43rd ACM Symposium on Theory of Computing, STOC 2011, San Jose, CA, USA, 6-8 June 2011*, pages 569–578, 2011.
- [16] Dan Feldman, Melanie Schmidt, and Christian Sohler. Turning big data into tiny data: Constant-size coresets for k-means, PCA and projective clustering. In *Proceedings of the Twenty-Fourth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2013, New Orleans, Louisiana, USA, January 6-8, 2013*, pages 1434–1453, 2013.
- [17] Dan Feldman, Melanie Schmidt, and Christian Sohler. Turning big data into tiny data: Constant-size coresets for k-means, pca, and projective clustering. *SIAM J. Comput.*, 49(3):601–657, 2020.
- [18] Hendrik Fichtenberger, Marc Gillé, Melanie Schmidt, Chris Schwiegelshohn, and Christian Sohler. BICO: BIRCH meets coresets for k-means clustering. In *Algorithms - ESA 2013 - 21st Annual European Symposium, Sophia Antipolis, France, September 2-4, 2013. Proceedings*, pages 481–492, 2013.
- [19] Gereon Frahling and Christian Sohler. Coresets in dynamic geometric data streams. In *Proceedings of the 37th Annual ACM Symposium on Theory of Computing (STOC)*, pages 209–217, 2005.
- [20] Sariel Har-Peled and Akash Kushal. Smaller coresets for k-median and k-means clustering. *Discrete & Computational Geometry*, 37(1):3–19, 2007.
- [21] Sariel Har-Peled and Soham Mazumdar. On coresets for k-means and k-median clustering. In *Proceedings of the 36th Annual ACM Symposium on Theory of Computing, Chicago, IL, USA, June 13-16, 2004*, pages 291–300, 2004.
- [22] Lingxiao Huang, Shaofeng H.-C. Jiang, and Nisheeth K. Vishnoi. Coresets for clustering with fairness constraints. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 7587–7598, 2019.
- [23] Lingxiao Huang and Nisheeth K. Vishnoi. Coresets for clustering in euclidean spaces: importance sampling is nearly optimal. In Konstantin Makarychev, Yury Makarychev, Madhur Tulsiani, Gautam Kamath, and Julia Chuzhoy, editors, *Proceedings of the 52nd Annual*

- ACM SIGACT Symposium on Theory of Computing, STOC 2020, Chicago, IL, USA, June 22-26, 2020*, pages 1416–1429. ACM, 2020.
- [24] Piotr Indyk, Sepideh Mahabadi, Shayan Oveis Gharan, and Alireza Rezaei. Composable core-sets for determinant maximization problems via spectral spanners. In Shuchi Chawla, editor, *Proceedings of the 2020 ACM-SIAM Symposium on Discrete Algorithms, SODA 2020, Salt Lake City, UT, USA, January 5-8, 2020*, pages 1675–1694. SIAM, 2020.
 - [25] Jan-Philipp W. Kappmeier, Daniel R. Schmidt, and Melanie Schmidt. Solving k-means on high-dimensional big data. In Evripidis Bampis, editor, *Experimental Algorithms - 14th International Symposium, SEA 2015, Paris, France, June 29 - July 1, 2015, Proceedings*, volume 9125 of *Lecture Notes in Computer Science*, pages 259–270. Springer, 2015.
 - [26] Michael Langberg and Leonard J. Schulman. Universal ϵ -approximators for integrals. In *Proceedings of the Twenty-First Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2010, Austin, Texas, USA, January 17-19, 2010*, pages 598–607, 2010.
 - [27] Michael Langberg and Leonard J. Schulman. Universal ϵ -approximators for integrals. In *Proceedings of the Twenty-First Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2010, Austin, Texas, USA, January 17-19, 2010*, pages 598–607, 2010.
 - [28] Alaa Maalouf, Ibrahim Jubran, and Dan Feldman. Fast and accurate least-mean-squares solvers. In *Advances in Neural Information Processing Systems*, pages 8307–8318, 2019.
 - [29] Sepideh Mahabadi, Konstantin Makarychev, Yury Makarychev, and Ilya P. Razenshteyn. Nonlinear dimension reduction via outer bi-lipschitz extensions. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2018, Los Angeles, CA, USA, June 25-29, 2018*, pages 1088–1101, 2018.
 - [30] Tung Mai, Anup B. Rao, and Cameron Musco. Coresets for classification - simplified and strengthened. *CoRR*, abs/2106.04254, 2021.
 - [31] Marina Meila. Comparing clusterings: an axiomatic view. In Luc De Raedt and Stefan Wrobel, editors, *Machine Learning, Proceedings of the Twenty-Second International Conference (ICML 2005), Bonn, Germany, August 7-11, 2005*, volume 119 of *ACM International Conference Proceeding Series*, pages 577–584. ACM, 2005.
 - [32] Marina Meila. The uniqueness of a good optimum for k-means. In William W. Cohen and Andrew W. Moore, editors, *Machine Learning, Proceedings of the Twenty-Third International Conference (ICML 2006), Pittsburgh, Pennsylvania, USA, June 25-29, 2006*, volume 148 of *ACM International Conference Proceeding Series*, pages 625–632. ACM, 2006.
 - [33] Alexander Munteanu, Chris Schwiegelshohn, Christian Sohler, and David P. Woodruff. On coresets for logistic regression. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 6562–6571, 2018.
 - [34] Shyam Narayanan and Jelani Nelson. Optimal terminal dimensionality reduction in euclidean space. In Moses Charikar and Edith Cohen, editors, *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing, STOC 2019, Phoenix, AZ, USA, June 23-26, 2019*, pages 1064–1069. ACM, 2019.
 - [35] Melanie Schmidt, Chris Schwiegelshohn, and Christian Sohler. Fair coresets and streaming algorithms for fair k-means. In *Approximation and Online Algorithms - 17th International Workshop, WAOA 2019, Munich, Germany, September 12-13, 2019, Revised Selected Papers*, pages 232–251, 2019.
 - [36] Christian Sohler and David P. Woodruff. Strong coresets for k-median and subspace approximation: Goodbye dimension. In *59th IEEE Annual Symposium on Foundations of Computer Science, FOCS 2018, Paris, France, October 7-9, 2018*, pages 802–813, 2018.
 - [37] Christian Sohler and David P. Woodruff. Strong coresets for k-median and subspace approximation: Goodbye dimension. *CoRR*, abs/1809.02961, 2018.
 - [38] Suresh Venkatasubramanian and Qiusi Wang. The johnson-lindenstrauss transform: An empirical study. In Matthias Müller-Hannemann and Renato Fonseca F. Werneck, editors, *Proceedings of the Thirteenth Workshop on Algorithm Engineering and Experiments, ALENEX 2011, Holiday Inn San Francisco Golden Gateway, San Francisco, California, USA, January 22, 2011*, pages 164–173. SIAM, 2011.
 - [39] Dennis Wei. A constant-factor bi-criteria approximation guarantee for k-means++. In Daniel D. Lee, Masashi Sugiyama, Ulrike von Luxburg, Isabelle Guyon, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 604–612, 2016.
 - [40] Tian Zhang, Raghu Ramakrishnan, and Miron Livny. BIRCH: A new data clustering algorithm and its applications. *Data Min. Knowl. Discov.*, 1(2):141–182, 1997.