

# An Empirical Evaluation of $k$ -Means Coresets

Anonymous author

Anonymous affiliation

Anonymous author

Anonymous affiliation

## 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 no work on comparing the quality of available  $k$ -means coresets.

In this paper we perform such an evaluation. We describe the difficulty in comparing the quality of different coreset algorithms. In order to perform an empirical evaluation, we therefore have to work with heuristics. To this end, we propose a new procedure for evaluating coresets on real-world datasets. To complement this, we propose a benchmark framework for generating provably hard synthetic data sets to evaluate coresets. Using this benchmark and real-world data sets, we conduct an exhaustive evaluation of the most commonly used coreset implementations.

**2012 ACM Subject Classification** Replace ccsdesc macro with valid one

**Keywords and phrases** Dummy keyword

**Digital Object Identifier** 10.4230/LIPIcs.SEA.2022.23

**Acknowledgements** Anonymous acknowledgements

## 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  $\text{dist}$  and aim to find  $k$  centers  $C$  such that

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

Omar: Shouldn't the cost be:

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

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  $(k, \varepsilon)$ -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  $C$

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

The coreset definition in Equation (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, 13, 16, 21, 19, 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]<sup>1</sup>.

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 access 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  $\varepsilon$ -net of a 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 coreset* of a point set  $B$ .

Due to this difficulty, a common heuristic for evaluating coresets is as follows [1, 17]. First, compute a coreset  $\Omega$  with the available algorithm(s) using some input data  $A$ . Then, run an optimization algorithm on  $\Omega$  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 mixes up 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.

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

## 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  $\text{OPT}$  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 [18, 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 [14, 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  $O\left(\frac{k \log n}{\delta^{d/2} \epsilon^d} \cdot \log^{d/2} \frac{k \log n}{\delta^{d/2} \epsilon^d}\right)$ .

**BICO [17]** Combines the very fast, but poor quality clustering algorithm BIRCH [40] with the movement-based analysis from [18, 21]. The clustering is organized by way of a hierarchical decomposition: When adding a point  $p$  to one of the coreset points  $\Omega$  at level  $i$ , it first finds the closest point  $q$  in  $\Omega$ . 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  $O(k \log n \epsilon^{-d-2})$ .

**Sensitivity Sampling [13]** The simplest implementation of sensitivity sampling first computes an  $(O(1), O(1))$  bicriteria  $K$  approximation<sup>2</sup>, for example by running  $k$ -means++ for  $2k$  iterations [39]. Let  $K$  be the  $2k$  clustering thus computed and let  $K_i$  be an arbitrary cluster of  $K$  with center  $q_i$ . Subsequently, the algorithm picks  $T - 2k$  points proportionate to  $\frac{\text{dist}^2(p, q_i)}{\text{cost}_{K_i}(\{q_i\})} + \frac{1}{|K_i|}$ . Let  $|\hat{K}_i|$  be the estimated number of points in the sample. Finally, the algorithm weighs each  $q_i$  by  $(1 + \epsilon) \cdot |K_i| - |\hat{K}_i|$ . The provable bound required for the algorithm to compute a coreset is  $\tilde{O}(k d \epsilon^{-4})$  ([13]),  $\tilde{O}(k \epsilon^{-6})$  ([23]), or  $\tilde{O}(k^2 \epsilon^{-4})$  ([5]).

<sup>2</sup> An  $(\alpha, \beta)$  bicriteria approximation computes an  $\alpha$  approximation using  $\beta \cdot k$  many centers.

**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  $\log^2 1/\epsilon$  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  $\tilde{O}(k\epsilon^{-4})$  ([11]).

**Ray Maker [20]** The algorithm computes an initial solution with  $k$  centers which is a constant factor approximation of the optimal clustering. Around each center,  $O(1/\epsilon^{d-1})$  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  $k^2/\epsilon^2$  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  $O(n + \text{poly}(k, \log n, 1/\epsilon) + \text{func}(k, \epsilon))$  where  $\text{poly}$  is a polynomial and  $\text{func}(k, \epsilon)$  is a function which depends on  $k$  and  $\epsilon$ .

## 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 [15], 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 [15, 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 [15] showed that projecting an input  $A$  onto the first  $O(k/\epsilon^2)$  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  $(1 \pm \epsilon)$  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 [12, 29, 34], with Narayanan and Nelson [34] achieving an optimal target dimension of  $O(\epsilon^{-2} \log n)$ , 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].

### 3 Hardness of Coreset Evaluation and a Benchmark

In this section, we first show that it is in general co-NP hard to evaluate the coreset distortion, given two point sets  $A$  and  $B$ . Thereafter we describe the benchmark and its properties.

► **Proposition 1.** *Given two point sets  $A$  and  $B$  in  $\mathbb{R}^d$  and a sufficiently small (constant)  $\varepsilon > 0$ , it is co-NP hard to decide whether  $A$  is a  $(k, \varepsilon)$ -coreset of  $B$ .*

**Proof.** First, we recall that for some  $\varepsilon_0$  and candidate clustering cost  $V$ , it is NP-hard to decide whether there exists a clustering  $C$  with cost in  $\text{cost}_A(C) \leq V$  and  $\text{cost}_B(C) \geq (1 + \varepsilon_0) \cdot V$ . Conversely, it is co-NP-hard to decide whether there exists no set of centers  $C$  such that  $\text{cost}_A(C) \leq V$  and  $\text{cost}_B(C) \geq (1 + \varepsilon_0) \cdot V$ . ◀

We remark that the possible values for  $\varepsilon_0$  are determined by the current APX-hardness results. Assuming  $\text{NP} \neq \text{P}$ ,  $\varepsilon_0 \approx 1.07$  and assuming UCG,  $\varepsilon_0 \approx 1.17$  [10, 9] for  $k$ -means in Euclidean spaces.

#### 3.1 Benchmark Construction

In this section, we describe our benchmark. The benchmark has a parameter  $\alpha$  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  $\mathbb{1}_k$  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  $k^\ell$  dimensional vector

$$v_i^\ell = \begin{bmatrix} (v_i^{\ell-1})_1 \cdot \mathbb{1}_k \\ (v_i^{\ell-1})_2 \cdot \mathbb{1}_k \\ \vdots \\ (v_i^{\ell-1})_{k^{\ell-1}} \cdot \mathbb{1}_k \end{bmatrix}. \text{ Finally, set the column } t = a \cdot k + b, a \in \{0, \dots, \alpha - 1\} \text{ and } b \in \{1, \dots, k\},$$

of  $A$  to be  $k^{\alpha-a}$  stacks of  $v_b^{a+1}$ .

To get a better feel for the construction, we have given two example inputs in Figure 1.

$$\begin{bmatrix} \frac{1}{2} & -\frac{1}{2} & \frac{1}{2} & -\frac{1}{2} & \frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} & \frac{1}{2} & -\frac{1}{2} & \frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} & \frac{1}{2} & \frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} & -\frac{1}{2} & \frac{1}{2} & \frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} & \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} & \frac{1}{2} \\ -\frac{1}{2} & \frac{1}{2} & \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} & \frac{1}{2} \end{bmatrix} \quad \begin{bmatrix} \frac{2}{3} & -\frac{1}{3} & -\frac{1}{3} & \frac{2}{3} & -\frac{1}{3} & -\frac{1}{3} \\ -\frac{1}{3} & \frac{2}{3} & -\frac{1}{3} & -\frac{1}{3} & \frac{2}{3} & -\frac{1}{3} \\ -\frac{1}{3} & -\frac{1}{3} & \frac{2}{3} & \frac{2}{3} & -\frac{1}{3} & -\frac{1}{3} \\ \frac{2}{3} & -\frac{1}{3} & -\frac{1}{3} & -\frac{1}{3} & \frac{2}{3} & -\frac{1}{3} \\ -\frac{1}{3} & \frac{2}{3} & -\frac{1}{3} & -\frac{1}{3} & -\frac{1}{3} & \frac{2}{3} \\ -\frac{1}{3} & -\frac{1}{3} & \frac{2}{3} & \frac{2}{3} & -\frac{1}{3} & -\frac{1}{3} \end{bmatrix}$$

■ **Figure 1** Benchmark construction for  $k = 2$  and  $\alpha = 3$  (left) and  $k = 3$  and  $\alpha = 2$  (right).

#### 3.2 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

## 23:6 An Empirical Evaluation of $k$ -Means Coresets

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

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

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

► **Lemma 2** (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  $\Pi_k$  the set of all permutations over  $k$  elements. Then the distance between  $\mathcal{C}$  and  $\mathcal{C}'$  is defined as

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

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

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

- 185 1. The distance between these clustering is maximized.
- 186 2. The clusterings have equal cost.
- 187 3. The clusterings are induced by a set of centers in  $\mathbb{R}^d$ .

188 The first and second property ensure that (equally good) solutions we use for evaluation  
189 are spread out through the solution space. In other words, we are less likely to only focus  
190 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) coreset 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  $\mathcal{C}^a$ .

► **Fact 3.** For  $a \neq a'$ , we have  $d(\mathcal{C}^a, \mathcal{C}^{a'}) = 1 - 1/k$ .

**Proof.** Consider an arbitrary vector  $v_i^\ell$ . By construction, the positive entries of  $v_i^\ell$  range from  $k^{\ell-1} \cdot i + 1$  to  $k^{\ell-1} \cdot (i + 1)$ . Similarly, the positive entries for the vector  $v_j^{\ell-1}$  range from range from  $k^{\ell-2} \cdot j + 1$  to  $k^{\ell-2} \cdot (j + 1)$ . Therefore, concatenating  $v_j^{\ell-1}$   $k$  times into a vector  $v'$ ,  $v'$  and  $v_i^\ell$  can share at most one positive coordinate. Inductively, the same holds true for any concatenation of vectors  $v_j^{\ell-h}$ . Thus, the two clusters induced by the columns formed by concatenating the vectors  $v$  can share only a  $1/k$  fraction of the points. Since each cluster consists of exactly  $k^\alpha/k = k^{\alpha-1}$  points, the confusion matrix  $M$  only has entries  $\frac{n}{k^2}$  and for any permutation  $\pi$ , we have  $d(\mathcal{C}^a, \mathcal{C}^{a'}) = 1 - 1/k$ . ◀

► **Fact 4.** For any  $\mathcal{C}_j^a$ , we have  $\text{cost}_{\mathcal{C}_j^a}(\{\mu(C_j^a)\}) = (\alpha - 1) \cdot k^{\alpha-2} \cdot (k - 1)$ .

**Proof.** Without loss of generality, we consider  $\mathcal{C}_1^0$ ; the proof is analogous for the other choices of  $j$  and  $a$ . We first note that for any point  $A_i \in C_1^0$ , the coordinates  $A_{i,\ell}$  are identical for  $\ell < k$ . Furthermore for the column  $\ell \geq k$ , we have by construction  $\sum_{A_i \in C_j^a} A_{i,\ell} = k^{\alpha-1} \cdot \frac{k-1}{k} + (k^\alpha - k^{\alpha-1}) \frac{1}{k} = k^{\alpha-1} \cdot (\frac{k-1}{k} - (k-1) \frac{1}{k}) = 0$ . Therefore, the mean of  $\mathcal{C}_1^0$  satisfies  $\mu(C_1^0)_\ell = \begin{cases} A_{i,\ell} & \text{if } \ell < k \\ 0 & \text{else.} \end{cases}$ . Thus, the cost is precisely  $(\alpha - 1) \cdot k^{\alpha-1} \cdot \left( \left( \frac{k-1}{k} \right)^2 + \left( \frac{1}{k} \right)^2 \right) = (\alpha - 1) \cdot k^{\alpha-2} \cdot (k - 1)$ . ◀

Finally, we show that the means for the clustering  $\mathcal{C}^a$  also induce  $\mathcal{C}^a$ .

► **Fact 5.** For a clustering  $\mathcal{C}^a$ , let  $\mu(C_j^a)$  denote the mean of cluster  $C_j^a$ . Then every point of *Omar*:  $A_i$  of  $\mathcal{C}_j^a$  is assigned to its closest center. Moreover, every point  $A_i$  of  $\mathcal{C}_j^a$  has equal distance to any center  $\mu(C_h^a)$  with  $h \neq j$ .

**Proof.** Again, we assume without loss of generality  $a = 0$ . Let  $A_i$  be an arbitrary point of cluster  $C_h^a$  and consider the mean  $\mu(C_j^a)_\ell = \begin{cases} A_{i,\ell} & \text{if } \ell < k \\ 0 & \text{else.} \end{cases}$  of cluster  $C_j^a$ . By definition, the positive coordinates of  $A_i$  are not equal to the positive coordinates of  $\mu(C_j^a)$ . The only difference in coordinates between the means of  $\mu(C_j^a)$  and  $\mu(C_h^a)$  are the first  $k$  coordinates, as the rest are 0. But here the coordinates of  $\mu(C_h^a)$  and  $A_i$  are identical, hence  $\mu(C_j^a)$  cannot be closer to  $A_i$ .

To prove that the distances between  $A_i$  and any  $\mu(C_h^a)$  with  $h \neq j$  are equal, again consider that any difference can only exist among the first  $k$  coordinates. Here, we have  $\mu(C_h^a)_h = \frac{k-1}{k}$ , and the remaining columns are  $-\frac{1}{k}$ . Since  $A_{i,h} = -\frac{1}{k}$  for any  $h \neq j$ , the claim follows. ◀

### 3.3 Benchmark Evaluation

We now describe how we use the benchmark to measure the distortion of a coreset. Assume that the coresets only consists of input point<sup>3</sup>

<sup>3</sup> It is not necessary for coreset 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



Consider the clustering  $\mathcal{C}^a = \{C_1^a, \dots, C_k^a\}$  for some  $a$  and let  $\Omega$  with weights  $w : \Omega \rightarrow \mathbb{R}_{\geq 0}$  be the coreset. 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  $\Omega$ . 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  $w(C_i^a \cap \Omega) \leq |C_i^a|(1 - \varepsilon)$  Omar: Replace  $\leq$  with  $<$ ?, 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. Omar: Replace all  $C_i$  with  $C_i^a$ ?

We now compare the cost as computed on the coreset and the true cost of  $\mathcal{S}$ . Due to Lemma 2 and Fact 5, 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 due to Fact 4

$$\begin{aligned} \text{cost}_{C_i^a}(\mathcal{S}) &= (\alpha - 1) \cdot k^{\alpha-2} \cdot (k - 1) + k^{\alpha-1} \cdot 2 \\ &\approx \left(1 + \frac{2}{\alpha}\right) \cdot \text{cost}_{C_j^a}(\{\mu(C_j^a)\}). \end{aligned}$$

Conversely, the cost on the coreset is

$$\begin{aligned} &\text{cost}_{\Omega \cap C_i^a}(\mathcal{S}) \\ &= w(C_i^a \cap \Omega) \left( (\alpha - 1) \cdot \left( \left( \frac{k-1}{k} \right)^2 + (k-1) \left( \frac{1}{k} \right)^2 \right) + 2 \right) \\ &\approx \frac{w(C_i^a \cap \Omega)}{\text{cost}_{C_j^a}(\{\mu(C_j^a)\})} \left(1 + \frac{2}{\alpha}\right) \cdot \text{cost}_{C_j^a}(\{\mu(C_j^a)\}). \end{aligned}$$

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.

## 3.4 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 a new 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 observations on the results.

### 4.1 Evaluation Procedure

Omar: Removed subsection “Motivation for Introducing New Evaluation Procedure” and wrote Evaluation Procedure in a slightly different way. What do you think? Finding a

---

and it makes the evaluation simpler, we will only consider coresets that are subsets of the original point set.



$k$ -clustering which results in a maximal distortion is difficult. Instead, we can estimate the quality of coresets on real-world data sets as follows: we first compute a coreset  $\Omega$ . Then, run  $k$ -means++ on  $\Omega$  to find a set of  $k$  centers  $C$ . We repeat the  $k$ -means++ algorithm 5 times to pick the set of centers with the largest distortion.

While it is useful to evaluate 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. Omar: The paragraph above can probably be reduced to two sentence about how we use the benchark, but I wanted to emphasize why the benchmark is so useful. I am okay if it gets removed due to lack of space.

Every randomized coreset construction was repeated 10 times. We aggregated the reported distortions by taking the maximum over all 10 evaluations. Chris: I don't get this. Did you take  $D = \frac{1}{10} \sum D_i$ , where  $D_i$  is the maximum distortion in the  $i$ th evaluation or  $D = \max D_i$ ? Because I would understand using the mean (or median), but not so much using the max both times. We repeat the experiment to filter out outliers and bad runs; using the max means we always take the biggest outlier. Omar: I did the later, meaning I took the maximum distortion over all 10 repetitions of an experiment. This is what I do exact.

For each iteration  $i \in \{1, 2, \dots, 10\}$  of an experiment:

1. Compute a coreset  $\Omega_i$  on  $A$
2. Run  $k$ -means++ algorithm 5 times.  
Each run of  $k$ -means++ on  $\Omega_i$  outputs a clustering  $C_{ij}$
3. Then, the best lower bound on the distortion fir  $\Omega_i$  is the solution with the largest distortion:

$$D_i = \max_j \left( \max \left( \frac{\text{cost}_A(C_{ij})}{\text{cost}_{\Omega_i}(C_{ij})}, \frac{\text{cost}_{\Omega_i}(C_{ij})}{\text{cost}_A(C_{ij})} \right) \right)$$

Having computed 10 distortions, I pick  $D = \max D_i$  because it captures the worst possible distortion that an algorithm can produce. Since the coreset definition does not say anything about an algorithm generating a coreset with a certain distortion with high probability, I assumed that a “bad” run with extremely high distortion should give an estimate of the general quality of the evaluated coreset construction. The idea is that an outlier is not a failure mode of the algorithm but an indication of how good it actually is. Not sure if this makes sense. What do you think?

In any case, I computed the the average distortions with variances over the 10 runs. I created an appendix with these results, see ?? and ?. I also regenerated the plots so they now show the average distortion  $D = \frac{1}{10} \sum D_i$ . In addition, we preprocessed the data using the dimension reduction techniques described in Section 2.

## 4.2 Data sets

We conducted experiments on five real-world data sets and four instances of our benchmark. The sizes of the real-world data sets are summarized in Table 1. Benchmark instances were

generated to match approximately the sizes of the real-world data sets. The chosen parameter values and the corresponding instance sizes are shown in Table 2. We now provide a brief description of each of the real-world data sets.

The *Census*<sup>4</sup> 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*<sup>5</sup> 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*<sup>6</sup>. 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 [17], *Caltech* was created by computing SIFT features from the images in the Caltech101<sup>7</sup> 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*<sup>8</sup> 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 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, 17]. 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 GitHub<sup>9</sup>. For BICO, we used the authors'

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

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

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

<sup>7</sup> [http://www.vision.caltech.edu/Image\\_Datasets/Caltech101/](http://www.vision.caltech.edu/Image_Datasets/Caltech101/)

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

<sup>9</sup> Link to repository will be provided later.

■ **Table 1** The sizes of the real-world data sets used for the experimental evaluation

	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

■ **Table 2** The parameter values and the sizes of the benchmark instances used for the experimental evaluation.

$k$	$\alpha$	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

reference implementation<sup>10</sup>. 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 Figure 4. All five algorithms are matched on the *Tower* dataset. The worst distortions across the algorithms are close to 1, and performance between the algorithms is negligible. The performance difference between sampling-based and movement-based methods become more pronounced as the number of dimensions increase. On *Covertypes* 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.

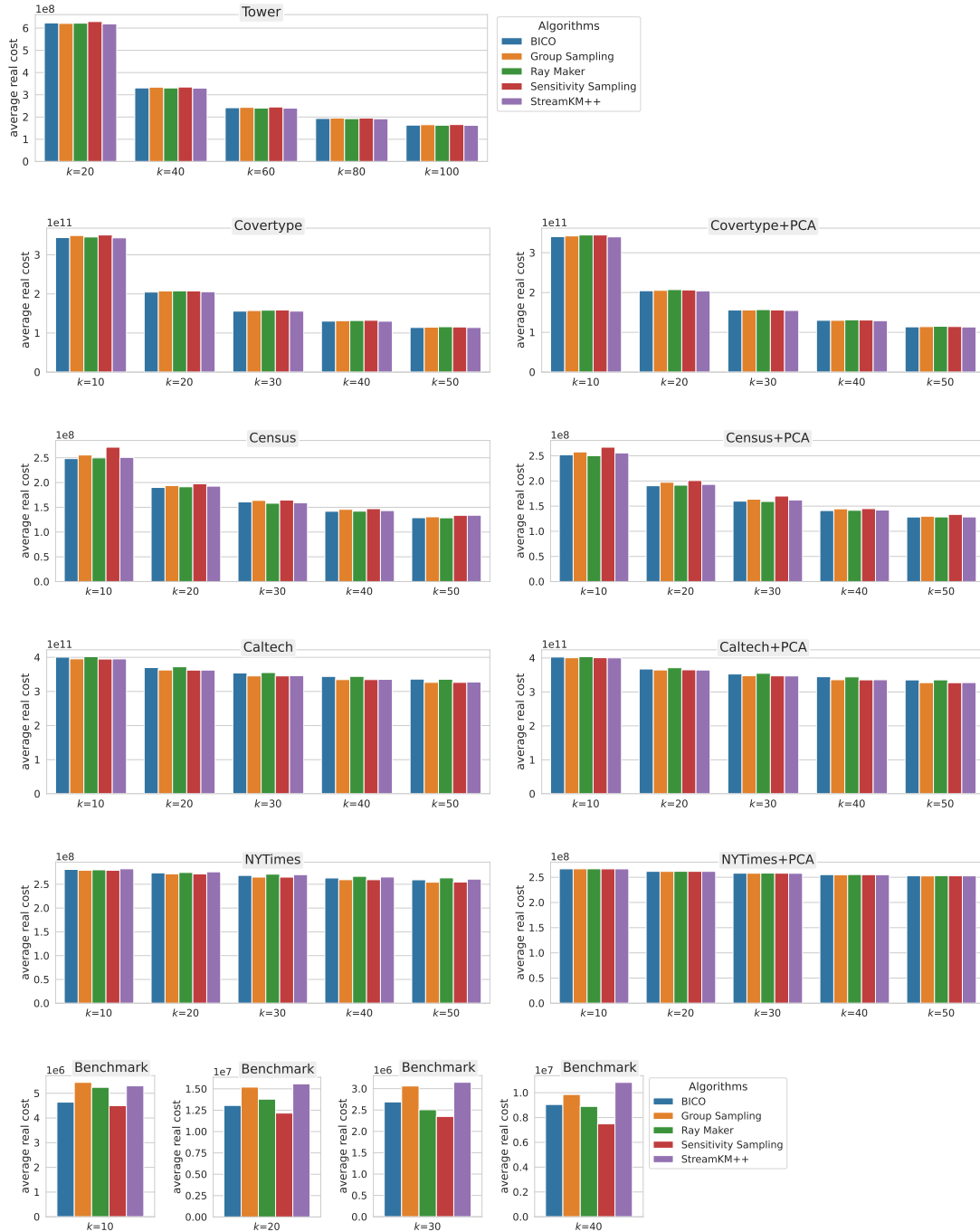
Reducing noise with PCA helped boost the performance on real-world data sets with large number of dimensions. On *Covertypes*, PCA does not change the performance numbers by much. On *Census* ( $d = 68$ ), the preprocessing step with PCA improves the performance of BICO and Ray Maker slightly for lower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

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

■ **Figure 2** 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.



**Figure 3** The average costs of running the evaluated coreset algorithms multiple times on different data sets. In general, the five coreset algorithms are able to compute coresets which result in solutions with comparable costs on the different real-world data sets. The differences in cost is more noticeable on the benchmark instances. Here, Sensitivity Sampling is the winner because it seems to be better at capturing the correct “clusters” inherent in the benchmark instances.



## 365 4.5 Interpretation of Experimental Results

366 Chris: Can we have one central take-away for each? I wrote a suggestion, let me know how  
367 you feel about them Omar: Yes! See my notes below.

### 368 Optimization versus Compression

369 For some data sets, optimization is very close to compression. It is not uncommon for the  
370  $k$ -means cost of real-world data sets to drop significantly for larger values of  $k$ . Figure 2  
371 illustrates this behavior for several real-world data sets. The more the curve bends, the less  
372 of a difference there is between computing a coreset and a clustering with low cost. For data  
373 sets with a L-shaped cost curve, a coreset algorithm adding more centers to the coreset will  
374 seem to be performing well when evaluating it based on the outcome of the optimization.  
375 *Tower* is a good example of a data set where optimization is very close to compression. Its  
376 cost curve bends the most which indicates that adding more centers help reduce the cost.  
377 One of the strengths of the benchmark is that there is no way of reducing the cost without  
378 capturing the right subclusters within a benchmark instance. This means that the cost does  
379 not decrease markedly beyond a certain value of  $k$  even if more centers are added as depicted  
380 in Figure 3.

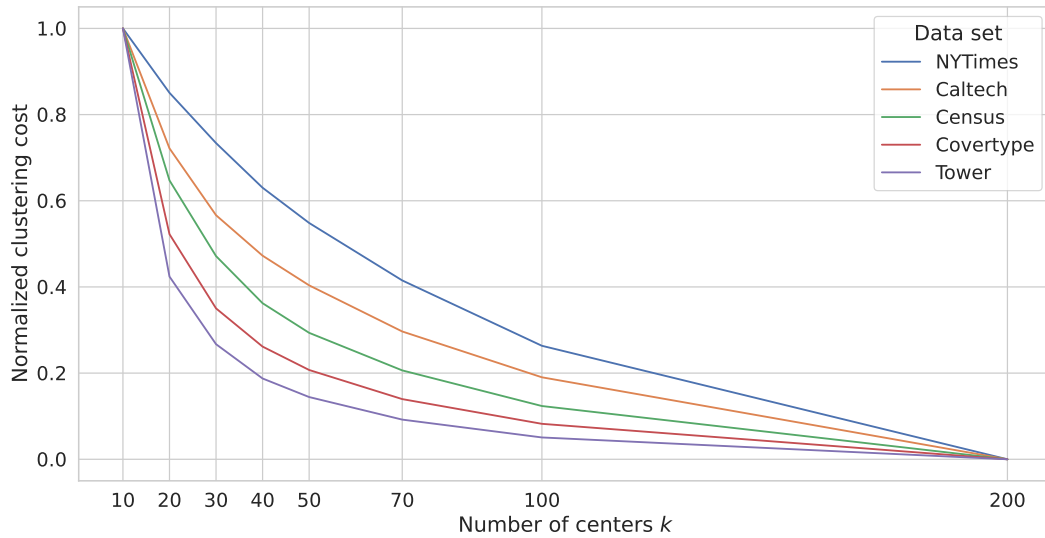
381 For BICO, Ray Maker, and StreamKM++, there is a correlation between the steepness  
382 of the cost curve for a data set and the distortion of the generated coreset. Chris: excellent.  
383 please add the aforementioned stuff here (plots and the parts from the paragraph that  
384 we will probably erase). Omar: Added at the beginning of the subsection for better text  
385 flow. For data sets where the curve is less steep, we observed higher distortions. The effect  
386 is more pronounced for the movement-based algorithms (BICO and Ray Maker) than for  
387 StreamKM++ which performs sampling. The other two sampling-based approaches (Group  
388 Sampling and Sensitivity Sampling) seem to be free from this behavior as they consistently  
389 generate high quality coresets irrespective of the shape of cost curve.

390 Chris: The aforementioned correlation between cost and distortion for movement and the  
391 fact that sampling seems to be rather oblivious to this. Omar: I agree, this take-away seems  
392 to be clear from the text.

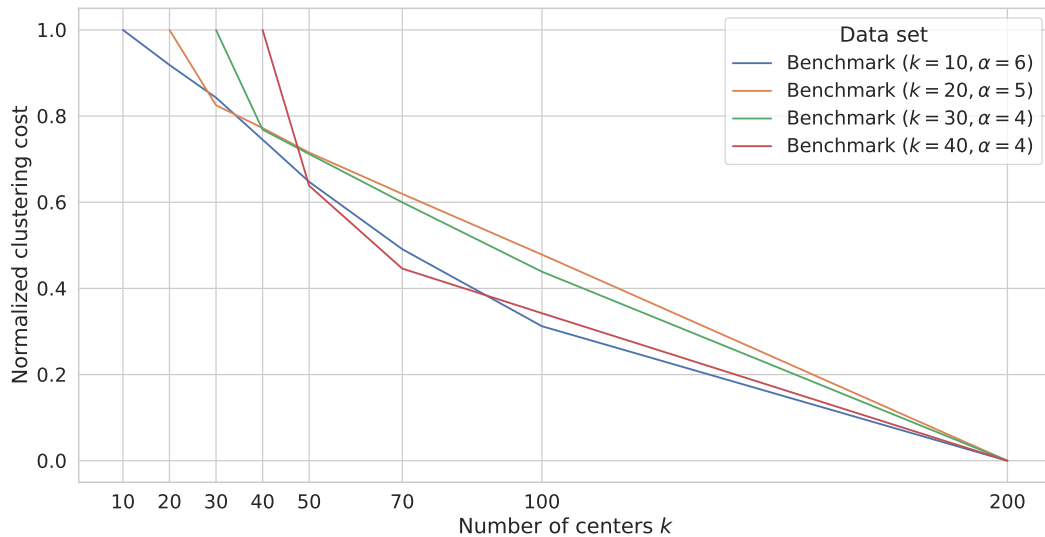
### 393 Movement-based versus Sampling-based Approaches

394 While all five algorithms are equally matched when comparing clustering costs (see Figure 5),  
395 coreset quality performance differ significantly (see Figure 4). In general, movement-based  
396 constructions perform the worst in terms of coreset quality. We observed that BICO and  
397 Ray Maker have the highest distortions across all data sets including on the benchmark  
398 instances. Among the sampling-based algorithms, Sensitive Sampling performs well with  
399 Group Sampling generally being competitive. StreamKM++ is an interesting case. Like  
400 the movement-based methods, its distortion increases with the dimension. Nevertheless, it  
401 generally performs significantly better than BICO and Ray Maker. This can be attributed to  
402 the fact that the coreset produced by StreamKM++ consists entirely of  $k$ -means++ centers  
403 weighted by the number of points of a minimal cost assignment. This is similar to movement-  
404 based algorithms such as BICO. Nevertheless, it also retains some of the performance from  
405 pure importance schemes. Chris: Can we say something about the impact of dimension?  
406 I am not sure I read the plot right (the non-linearity of the axis makes it a bit difficult),  
407 but is the distortion of the sampling based stuff affected by the dimension, or more or less  
408 independent? Also, I think we want to highlight more one key take away message: Dimension  
409 affects distortion for movement based in practise, not just in theory. Omar: Yes, we can say

**Figure 4** 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*, *Covertypes*, 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.

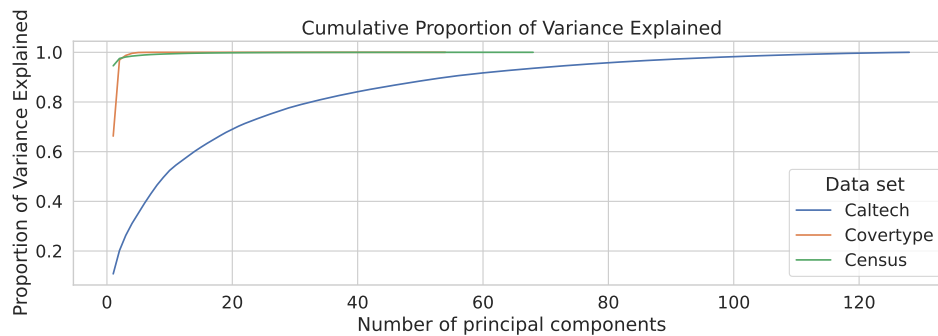


**Figure 5** Shows the clustering costs of four instances of the benchmark framework as a function of the number of centers. In contrast to real-world data sets, the costs do not decrease rapidly as more cluster centers are added.





■ **Figure 6** The cumulative proportion of explained variance by principal components on *Caltech*, *Covertypes*, and *Census*.



410 something about the impact of dimension. Please see ?? for the distortion numbers. What  
 411 do you think of adding the following: In practice as well as in theory, the distortion of  
 412 movement-based algorithms are affected by the dimension. By comparison, sampling-based  
 413 algorithms are affected very little by the dimension. The distortion of StreamKM++ increase  
 414 only slightly with the dimension while Group Sampling and Sensitivity Sampling are largely  
 415 unaffected.

## 416 Impact of PCA

417 On almost all our data sets, the performance improves when input data is preprocessed with  
 418 PCA, especially for the movement-based algorithms. Empirically, the more noise is removed  
 419 (i.e., small  $k$  value), the lower the distortion. Notice that  $k$  is the number of principal  
 420 components that the input data is projected on to. The rest of the low variance components  
 421 are treated as noise and removed. Method utilizing sampling (Group Sampling, Sensitivity  
 422 Sampling and StreamKM++) are less effected by the preprocessing step. On *Covertypes*,  
 423 PCA does not change the distortions by much because almost all the variance in the data is  
 424 explained by the first five principal components (see Figure 6). On *Caltech* and *NYTimes*,  
 425 the quality of the coresets by BICO and Ray Maker improves greatly because the noise  
 426 removal is more aggressive. Even if the quality is much better for movement-based coreset  
 427 constructions due to PCA transformation, importance sampling methods are still superior  
 428 when it comes to the quality of the compression. Chris: Take away message at the end: All  
 429 coreset constructions are affected by the dimension (even sampling, though markedly less  
 430 so). We can say, PCA is computationally very expensive, but may be nevertheless be worth  
 431 it. Omar: Agreed, that is a concluding message!

## 432 5 Conclusion

433 In this work, we studied how to assess the quality of  $k$ -means coresets computed by state-of-  
 434 the-art algorithms. It is generally hard to measure the quality of a coreset because it requires  
 435 computing the worst-case distortion. Due to this difficulty, earlier works evaluated coresets  
 436 by the outcome of an optimization algorithm. This method of comparison has the drawback  
 437 that it is more likely to measure the performance of the underlying optimization problem,  
 438 rather than evaluating coresets. As a alternative, we proposed a new evaluation procedure  
 439 which can estimate the quality of coresets on real-world data sets. To complement this, we  
 440 also proposed a benchmark framework which provably generates hard instances for all known

$k$ -means coreset algorithms. We evaluated the quality of five  $k$ -means coreset algorithms on five real-world datasets and four instances of the proposed benchmark. We experimented with both movement-based and sampling-based coreset algorithms. We found that while all algorithms produce coresets which yield similar low cost clusterings, sampling-based methods are superior to movement-based algorithms. Chris: Add the  $k$ -means++ question here. I'm not sure if I ever mentioned that before. If you came up with it by yourself, I am very impressed. It also flowed naturally in what you were writing. Omar: I wish I came up with it myself! I do not yet understand the proofs. Anyway, I moved the section "BICO versus StreamKM++" to here and tried to make it shorter. What do you think?

Surprisingly, the quality of the coresets produced by StreamKM++ is significantly better than what one would expect from its theoretical analysis. StreamKM++ is derived from a movement-based construction similar to BICO. Empirically, StreamKM++ is notably better than BICO across all data sets, and especially on high dimensional data. This begs the question whether there exist a better theoretical analysis for StreamKM++. We leave this as an open problem for future research.

## References

- 1 Marcel R. Ackermann, Marcus Mörtens, Christoph Raupach, Kamil Swierkot, Christiane Lammersen, and Christian Sohler. Streamkm++: A clustering algorithm for data streams. *ACM Journal of Experimental Algorithmics*, 17(1), 2012. URL: <http://doi.acm.org/10.1145/2133803.2184450>, doi:10.1145/2133803.2184450.
- 2 Pankaj K. Agarwal, Sarel 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. URL: <http://dl.acm.org/citation.cfm?id=1283383.1283494>.
- 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. doi:10.1145/3313276.3316318.
- 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. doi:10.1137/1.9781611976465.159.
- 6 Timothy M. Chan. Dynamic coresets. *Discret. Comput. Geom.*, 42(3):469–488, 2009. doi:10.1007/s00454-009-9165-3.
- 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  $l_p$  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. doi:10.1109/FOCS.2019.00040.

- 489 10 Vincent Cohen-Addad, Karthik C. S., and Euiwoong Lee. On approximability of clustering  
490 problems without candidate centers. In Dániel Marx, editor, *Proceedings of the 2021 ACM-*  
491 *SIAM Symposium on Discrete Algorithms, SODA 2021, Virtual Conference, January 10 - 13,*  
492 *2021*, pages 2635–2648. SIAM, 2021. doi:10.1137/1.9781611976465.156.
- 493 11 Vincent Cohen-Addad, David Saulpic, and Chris Schwiegelshohn. A new coreset framework  
494 for clustering. In Samir Khuller and Virginia Vassilevska Williams, editors, *STOC '21: 53rd*  
495 *Annual ACM SIGACT Symposium on Theory of Computing, Virtual Event, Italy, June 21-25,*  
496 *2021*, pages 169–182. ACM, 2021.
- 497 12 Michael Elkin, Arnold Filtser, and Ofer Neiman. Terminal embeddings. *Theor. Comput. Sci.*,  
498 697:1–36, 2017. doi:10.1016/j.tcs.2017.06.021.
- 499 13 Dan Feldman and Michael Langberg. A unified framework for approximating and clustering  
500 data. In *Proceedings of the 43rd ACM Symposium on Theory of Computing, STOC 2011, San*  
501 *Jose, CA, USA, 6-8 June 2011*, pages 569–578, 2011.
- 502 14 Dan Feldman and Michael Langberg. A unified framework for approximating and clustering  
503 data. In *Proceedings of the 43rd ACM Symposium on Theory of Computing, STOC 2011, San*  
504 *Jose, CA, USA, 6-8 June 2011*, pages 569–578, 2011.
- 505 15 Dan Feldman, Melanie Schmidt, and Christian Sohler. Turning big data into tiny data:  
506 Constant-size coresets for  $k$ -means, PCA and projective clustering. In *Proceedings of the*  
507 *Twenty-Fourth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2013, New*  
508 *Orleans, Louisiana, USA, January 6-8, 2013*, pages 1434–1453, 2013.
- 509 16 Dan Feldman, Melanie Schmidt, and Christian Sohler. Turning big data into tiny data:  
510 Constant-size coresets for  $k$ -means, pca, and projective clustering. *SIAM J. Comput.*, 49(3):601–  
511 657, 2020. doi:10.1137/18M1209854.
- 512 17 Hendrik Fichtenberger, Marc Gillé, Melanie Schmidt, Chris Schwiegelshohn, and Christian  
513 Sohler. BICO: BIRCH meets coresets for  $k$ -means clustering. In *Algorithms - ESA 2013 -*  
514 *21st Annual European Symposium, Sophia Antipolis, France, September 2-4, 2013. Proceedings,*  
515 pages 481–492, 2013.
- 516 18 Gereon Frahling and Christian Sohler. Coresets in dynamic geometric data streams. In  
517 *Proceedings of the 37th Annual ACM Symposium on Theory of Computing (STOC)*, pages  
518 209–217, 2005.
- 519 19 Sariel Har-Peled and Akash Kushal. Smaller coresets for  $k$ -median and  $k$ -means clustering.  
520 *Discrete & Computational Geometry*, 37(1):3–19, 2007.
- 521 20 Sariel Har-Peled and Akash Kushal. Smaller coresets for  $k$ -median and  $k$ -means clustering.  
522 *Discrete & Computational Geometry*, 37(1):3–19, 2007.
- 523 21 Sariel Har-Peled and Soham Mazumdar. On coresets for  $k$ -means and  $k$ -median clustering. In  
524 *Proceedings of the 36th Annual ACM Symposium on Theory of Computing, Chicago, IL, USA,*  
525 *June 13-16, 2004*, pages 291–300, 2004.
- 526 22 Lingxiao Huang, Shaofeng H.-C. Jiang, and Nisheeth K. Vishnoi. Coresets for clustering  
527 with fairness constraints. In *Advances in Neural Information Processing Systems 32: Annual*  
528 *Conference on Neural Information Processing Systems 2019, NeurIPS 2019, 8-14 December*  
529 *2019, Vancouver, BC, Canada*, pages 7587–7598, 2019. URL: [http://papers.nips.cc/paper/](http://papers.nips.cc/paper/8976-coresets-for-clustering-with-fairness-constraints)  
530 [8976-coresets-for-clustering-with-fairness-constraints](http://papers.nips.cc/paper/8976-coresets-for-clustering-with-fairness-constraints).
- 531 23 Lingxiao Huang and Nisheeth K. Vishnoi. Coresets for clustering in euclidean spaces: im-  
532 portance sampling is nearly optimal. In Konstantin Makarychev, Yury Makarychev, Madhur  
533 Tulsiani, Gautam Kamath, and Julia Chuzhoy, editors, *Proceedings of the 52nd Annual ACM*  
534 *SIGACT Symposium on Theory of Computing, STOC 2020, Chicago, IL, USA, June 22-26,*  
535 *2020*, pages 1416–1429. ACM, 2020. doi:10.1145/3357713.3384296.
- 536 24 Piotr Indyk, Sepideh Mahabadi, Shayan Oveis Gharan, and Alireza Rezaei. Composable  
537 core-sets for determinant maximization problems via spectral spanners. In Shuchi Chawla,  
538 editor, *Proceedings of the 2020 ACM-SIAM Symposium on Discrete Algorithms, SODA 2020,*  
539 *Salt Lake City, UT, USA, January 5-8, 2020*, pages 1675–1694. SIAM, 2020. doi:10.1137/1.  
540 [9781611975994.103](https://doi.org/10.1137/1.9781611975994.103).

- 541 25 Jan-Philipp W. Kappmeier, Daniel R. Schmidt, and Melanie Schmidt. Solving k-means on  
542 high-dimensional big data. In *Experimental Algorithms - 14th International Symposium,*  
543 *SEA 2015, Paris, France, June 29 - July 1, 2015, Proceedings*, pages 259–270, 2015. doi:  
544 10.1007/978-3-319-20086-6\_20.
- 545 26 Michael Langberg and Leonard J. Schulman. Universal  $\varepsilon$ -approximators for integrals. In  
546 *Proceedings of the Twenty-First Annual ACM-SIAM Symposium on Discrete Algorithms,*  
547 *SODA 2010, Austin, Texas, USA, January 17-19, 2010*, pages 598–607, 2010.
- 548 27 Michael Langberg and Leonard J. Schulman. Universal  $\varepsilon$ -approximators for integrals. In  
549 *Proceedings of the Twenty-First Annual ACM-SIAM Symposium on Discrete Algorithms,*  
550 *SODA 2010, Austin, Texas, USA, January 17-19, 2010*, pages 598–607, 2010. URL: <http://dx.doi.org/10.1137/1.9781611973075.50>, doi:10.1137/1.9781611973075.50.
- 551 28 Alaa Maalouf, Ibrahim Jubran, and Dan Feldman. Fast and accurate least-mean-squares  
552 solvers. In *Advances in Neural Information Processing Systems*, pages 8307–8318, 2019.
- 553 29 Sepideh Mahabadi, Konstantin Makarychev, Yury Makarychev, and Ilya P. Razenshteyn.  
554 Nonlinear dimension reduction via outer bi-lipschitz extensions. In *Proceedings of the 50th*  
555 *Annual ACM SIGACT Symposium on Theory of Computing, STOC 2018, Los Angeles, CA,*  
556 *USA, June 25-29, 2018*, pages 1088–1101, 2018. URL: <http://doi.acm.org/10.1145/3188745.3188828>, doi:10.1145/3188745.3188828.
- 557 30 Tung Mai, Anup B. Rao, and Cameron Musco. Coresets for classification - simplified and  
558 strengthened. *CoRR*, abs/2106.04254, 2021. URL: <https://arxiv.org/abs/2106.04254>,  
559 [arXiv:2106.04254](https://arxiv.org/abs/2106.04254).
- 560 31 Marina Meila. Comparing clusterings: an axiomatic view. In Luc De Raedt and Stefan Wrobel,  
561 editors, *Machine Learning, Proceedings of the Twenty-Second International Conference (ICML*  
562 *2005), Bonn, Germany, August 7-11, 2005*, volume 119 of *ACM International Conference*  
563 *Proceeding Series*, pages 577–584. ACM, 2005. doi:10.1145/1102351.1102424.
- 564 32 Marina Meila. The uniqueness of a good optimum for k-means. In William W. Cohen and  
565 Andrew W. Moore, editors, *Machine Learning, Proceedings of the Twenty-Third International*  
566 *Conference (ICML 2006), Pittsburgh, Pennsylvania, USA, June 25-29, 2006*, volume 148 of  
567 *ACM International Conference Proceeding Series*, pages 625–632. ACM, 2006. doi:10.1145/  
568 1143844.1143923.
- 569 33 Alexander Munteanu, Chris Schwiegelshohn, Christian Sohler, and David P. Wood-  
570 ruff. On coresets for logistic regression. In Samy Bengio, Hanna M. Wallach, Hugo  
571 Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Ad-*  
572 *vances in Neural Information Processing Systems 31: Annual Conference on Neural*  
573 *Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal,*  
574 *Canada*, pages 6562–6571, 2018. URL: [https://proceedings.neurips.cc/paper/2018/hash/](https://proceedings.neurips.cc/paper/2018/hash/63bfd6e8f26d1d3537f4c5038264ef36-Abstract.html)  
575 [63bfd6e8f26d1d3537f4c5038264ef36-Abstract.html](https://proceedings.neurips.cc/paper/2018/hash/63bfd6e8f26d1d3537f4c5038264ef36-Abstract.html).
- 576 34 Shyam Narayanan and Jelani Nelson. Optimal terminal dimensionality reduction in euclidean  
577 space. In Moses Charikar and Edith Cohen, editors, *Proceedings of the 51st Annual ACM*  
578 *SIGACT Symposium on Theory of Computing, STOC 2019, Phoenix, AZ, USA, June 23-26,*  
579 *2019*, pages 1064–1069. ACM, 2019. doi:10.1145/3313276.3316307.
- 580 35 Melanie Schmidt, Chris Schwiegelshohn, and Christian Sohler. Fair coresets and streaming  
581 algorithms for fair k-means. In *Approximation and Online Algorithms - 17th International*  
582 *Workshop, WAOA 2019, Munich, Germany, September 12-13, 2019, Revised Selected Papers,*  
583 *pages 232–251, 2019. doi:10.1007/978-3-030-39479-0\_16.*
- 584 36 Christian Sohler and David P. Woodruff. Strong coresets for k-median and subspace ap-  
585 proximation: Goodbye dimension. In *59th IEEE Annual Symposium on Foundations of*  
586 *Computer Science, FOCS 2018, Paris, France, October 7-9, 2018*, pages 802–813, 2018.  
587 doi:10.1109/FOCS.2018.00081.
- 588 37 Christian Sohler and David P. Woodruff. Strong coresets for k-median and subspace approx-  
589 imation: Goodbye dimension. *CoRR*, abs/1809.02961, 2018. URL: [http://arxiv.org/abs/](http://arxiv.org/abs/1809.02961)  
590 [1809.02961](http://arxiv.org/abs/1809.02961), [arXiv:1809.02961](http://arxiv.org/abs/1809.02961).
- 591 592

- 593    **38**    Suresh Venkatasubramanian and Qiushi Wang. The johnson-lindenstrauss transform: An  
594    empirical study. In Matthias Müller-Hannemann and Renato Fonseca F. Werneck, editors,  
595    *Proceedings of the Thirteenth Workshop on Algorithm Engineering and Experiments, ALENEX*  
596    *2011, Holiday Inn San Francisco Golden Gateway, San Francisco, California, USA, January*  
597    *22, 2011*, pages 164–173. SIAM, 2011.
- 598    **39**    Dennis Wei. A constant-factor bi-criteria approximation guarantee for k-means++. In  
599    Daniel D. Lee, Masashi Sugiyama, Ulrike von Luxburg, Isabelle Guyon, and Roman Garnett,  
600    editors, *Advances in Neural Information Processing Systems 29: Annual Conference on Neural*  
601    *Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 604–612,  
602    2016.
- 603    **40**    Tian Zhang, Raghu Ramakrishnan, and Miron Livny. BIRCH: A new data clustering algorithm  
604    and its applications. *Data Min. Knowl. Discov.*, 1(2):141–182, 1997. URL: [http://dx.doi.](http://dx.doi.org/10.1023/A:1009783824328)  
605    [org/10.1023/A:1009783824328](http://dx.doi.org/10.1023/A:1009783824328), doi:10.1023/A:1009783824328.

606

**A** Distortion results

■ **Table 3** Distortions of the evaluation algorithms on all data sets without PCA preprocessing. Each cell specify the mean distortion along with the standard deviation in brackets of 10 repetition of the experiment.

Data set	$k$	BICO	Group Sampling	Ray Maker	Sensitivity Sampling	StreamKM++
Benchmark	10	2.93 (0.144)	1.01 (0.002)	3.49 (0.055)	1.02 (0.003)	1.04 (0.002)
	20	2.87 (0.137)	1.01 (0.002)	3.90 (0.112)	1.02 (0.003)	1.11 (0.002)
	30	2.34 (0.070)	1.02 (0.001)	4.19 (0.137)	1.02 (0.002)	1.10 (0.002)
	40	2.68 (0.230)	1.02 (0.001)	3.83 (0.172)	1.02 (0.002)	1.17 (0.004)
Caltech	10	4.15 (0.485)	1.02 (0.003)	4.99 (0.124)	1.01 (0.004)	1.10 (0.004)
	20	4.16 (0.708)	1.02 (0.003)	4.86 (0.134)	1.01 (0.003)	1.12 (0.004)
	30	4.15 (0.416)	1.02 (0.001)	4.78 (0.117)	1.01 (0.002)	1.13 (0.003)
	40	4.18 (0.486)	1.02 (0.002)	4.73 (0.068)	1.01 (0.002)	1.13 (0.001)
	50	3.99 (0.560)	1.02 (0.002)	4.72 (0.090)	1.01 (0.002)	1.13 (0.002)
Census	10	1.65 (0.082)	1.03 (0.011)	1.74 (0.044)	1.02 (0.006)	1.05 (0.005)
	20	1.74 (0.054)	1.02 (0.004)	1.75 (0.038)	1.01 (0.005)	1.06 (0.003)
	30	1.77 (0.063)	1.02 (0.003)	1.79 (0.032)	1.01 (0.005)	1.07 (0.002)
	40	1.83 (0.038)	1.02 (0.002)	1.81 (0.021)	1.01 (0.005)	1.08 (0.003)
	50	1.87 (0.088)	1.02 (0.002)	1.85 (0.044)	1.01 (0.003)	1.10 (0.059)
Coverttype	10	1.10 (0.002)	1.03 (0.011)	1.18 (0.016)	1.02 (0.007)	1.01 (0.002)
	20	1.11 (0.004)	1.03 (0.010)	1.21 (0.015)	1.01 (0.007)	1.01 (0.002)
	30	1.10 (0.003)	1.03 (0.005)	1.22 (0.008)	1.01 (0.005)	1.01 (0.002)
	40	1.09 (0.001)	1.03 (0.005)	1.22 (0.012)	1.01 (0.005)	1.01 (0.001)
	50	1.07 (0.001)	1.03 (0.003)	1.22 (0.008)	1.01 (0.002)	1.01 (0.001)
NYTimes	10	18.59 (5.356)	1.04 (0.003)	15.37 (0.740)	1.02 (0.004)	1.75 (0.105)
	20	9.98 (1.559)	1.04 (0.003)	12.71 (0.448)	1.01 (0.004)	1.69 (0.024)
	30	8.03 (2.156)	1.04 (0.002)	12.18 (0.974)	1.02 (0.003)	1.66 (0.030)
	40	6.31 (0.436)	1.03 (0.003)	11.20 (0.742)	1.01 (0.002)	1.63 (0.024)
	50	5.97 (0.381)	1.03 (0.002)	10.81 (0.571)	1.01 (0.002)	1.63 (0.026)
Tower	20	1.07 (0.003)	1.02 (0.007)	1.06 (0.004)	1.03 (0.009)	1.02 (0.001)
	40	1.06 (0.004)	1.02 (0.006)	1.06 (0.002)	1.02 (0.007)	1.02 (0.001)
	60	1.06 (0.001)	1.03 (0.004)	1.05 (0.004)	1.01 (0.003)	1.02 (0.001)
	80	1.05 (0.001)	1.03 (0.005)	1.05 (0.002)	1.01 (0.006)	1.02 (0.001)
	100	1.04 (0.002)	1.03 (0.004)	1.05 (0.001)	1.01 (0.004)	1.02 (0.001)

■ **Table 4** Distortions of the evaluation algorithms on the data sets preprocessed with PCA. Each cell specify the mean distortion along with the standard deviation in brackets of 10 repetitions.

Data set	$k$	BICO	Group Sampling	Ray Maker	Sensitivity Sampling	StreamKM++
Caltech+PCA	10	1.16 (0.004)	1.01 (0.002)	1.19 (0.009)	1.00 (0.001)	1.02 (0.002)
	20	1.42 (0.016)	1.01 (0.002)	1.51 (0.008)	1.01 (0.001)	1.04 (0.002)
	30	1.73 (0.045)	1.02 (0.001)	1.83 (0.014)	1.01 (0.001)	1.06 (0.002)
	40	2.13 (0.083)	1.02 (0.002)	2.16 (0.016)	1.01 (0.002)	1.07 (0.002)
	50	2.38 (0.139)	1.02 (0.001)	2.49 (0.017)	1.01 (0.002)	1.09 (0.002)
Census+PCA	10	1.18 (0.012)	1.01 (0.005)	1.20 (0.016)	1.02 (0.009)	1.01 (0.002)
	20	1.44 (0.046)	1.02 (0.003)	1.45 (0.019)	1.02 (0.010)	1.04 (0.002)
	30	1.64 (0.044)	1.02 (0.004)	1.63 (0.021)	1.01 (0.007)	1.06 (0.003)
	40	1.77 (0.035)	1.02 (0.002)	1.76 (0.028)	1.01 (0.004)	1.07 (0.002)
	50	1.87 (0.068)	1.02 (0.002)	1.82 (0.017)	1.01 (0.004)	1.08 (0.002)
Covertime+PCA	10	1.11 (0.003)	1.03 (0.012)	1.19 (0.011)	1.02 (0.010)	1.02 (0.002)
	20	1.11 (0.003)	1.03 (0.006)	1.21 (0.010)	1.02 (0.007)	1.02 (0.002)
	30	1.10 (0.002)	1.03 (0.007)	1.22 (0.011)	1.01 (0.005)	1.01 (0.002)
	40	1.09 (0.002)	1.03 (0.006)	1.23 (0.010)	1.01 (0.005)	1.01 (0.001)
	50	1.07 (0.001)	1.03 (0.003)	1.23 (0.008)	1.01 (0.003)	1.01 (0.001)
NYTimes+PCA	10	1.01 (0.000)	1.00 (0.000)	1.01 (0.001)	1.00 (0.000)	1.00 (0.000)
	20	1.02 (0.000)	1.00 (0.000)	1.02 (0.001)	1.00 (0.000)	1.01 (0.000)
	30	1.03 (0.000)	1.00 (0.000)	1.03 (0.001)	1.00 (0.000)	1.01 (0.000)
	40	1.04 (0.001)	1.00 (0.000)	1.04 (0.001)	1.00 (0.002)	1.01 (0.000)
	50	1.04 (0.002)	1.00 (0.000)	1.05 (0.001)	1.00 (0.000)	1.02 (0.000)