

# Clustering Data Streams

Sudipto Guha <sup>\*</sup>    Nina Mishra <sup>†</sup>    Rajeev Motwani <sup>‡</sup>    Liadan O’Callaghan <sup>§</sup>

## Abstract

*We study clustering under the data stream model of computation where: given a sequence of points, the objective is to maintain a consistently good clustering of the sequence observed so far, using a small amount of memory and time. The data stream model is relevant to new classes of applications involving massive data sets, such as web click stream analysis and multimedia data analysis. We give constant-factor approximation algorithms for the  $k$ -Median problem in the data stream model of computation in a single pass. We also show negative results implying that our algorithms cannot be improved in a certain sense.*

## 1 Introduction

A data stream is an ordered sequence of points that can be read only once or a small number of times. Formally, a data stream is a sequence of points  $x_1, \dots, x_i, \dots, x_n$  read in increasing order of the indices  $i$ . The performance of an algorithm that operates on data streams is measured by the number of passes the algorithm must make over the stream, when constrained in terms of available memory, in addition to the more conventional measures. The data stream model is motivated by emerging application involving massive data sets, e.g., customer click streams, telephone records, large sets of web pages, multimedia data, and sets of retail chain transactions can be modeled as data streams. These data sets are far too

large to fit in main memory and are typically stored in secondary storage devices, making access, particularly random access, very expensive. Data stream algorithms access the input only via linear scans without random access and only require a few (hopefully, one) such scans over the data. Furthermore, since the amount of data far exceeds the amount of space (main memory) available to the algorithm, it is not possible for the algorithm to “remember” too much of the data scanned in the past. This scarcity of space necessitates the design of a novel kind of algorithm that stores only a *summary* of past data, leaving enough memory for the processing of future data. We remark that this is not the same as the model of online algorithms.

Clustering has recently been widely studied across several disciplines, but only a few of the techniques developed scale to support clustering of very large data sets. A common formulation of clustering is the  $k$ -Median problem: find  $k$  centers in a set of  $n$  points so as to minimize the sum of distances from data points to their closest cluster centers. Most algorithms for  $k$ -Median have large space requirements and involve random access to the input data. We give constant-factor approximation algorithms for the  $k$ -Median problem that naturally fit into this data stream setting. Our algorithms make a single pass over the data and use small space. We first give a randomized constant-factor approximation algorithm for  $k$ -Median, which makes one pass over the data using  $n^\epsilon$  memory (for  $\epsilon < 1$ ) and requires only  $\tilde{O}(nk)$  time. We also prove that any deterministic  $k$ -Median algorithm that achieves a constant-factor approximation cannot run in time less than  $\Omega(nk)$ . Finally, we give a deterministic  $\tilde{O}(nk)$ -time,  $\text{polylog}(n)$ -approximation single-pass algorithm that uses  $n^\epsilon$  space, for  $\epsilon < 1$ .

**Related Work on Data Streams** One of the first results in data streams was the result of Munro and Paterson [16], where they studied the space requirement of selection and sorting as a function of the number of passes over the data. The model was formalized by Henzinger, Raghavan, and Rajagopalan [7], who gave several algorithms and complexity results re-

---

<sup>\*</sup>Department of Computer Science, Stanford University, CA 94305. Email: [sudipto@cs.stanford.edu](mailto:sudipto@cs.stanford.edu). Research supported by IBM Research Fellowship and NSF Grant IIS-9811904.

<sup>†</sup>Hewlett Packard Laboratories, Palo Alto, CA 94304, Email: [nmishra@hpl.hp.com](mailto:nmishra@hpl.hp.com)

<sup>‡</sup>Department of Computer Science, Stanford University, CA 94305. Email: [rajeev@cs.stanford.edu](mailto:rajeev@cs.stanford.edu). Research supported in part by NSF Grant IIS-9811904.

<sup>§</sup>Department of Computer Science, Stanford University, CA 94305. Email: [loc@cs.stanford.edu](mailto:loc@cs.stanford.edu). Research supported in part by an NSF Graduate Fellowship, ARO MURI Grant DAAH04-96-1-0007, and NSF Grant IIS-9811904.

lated to graph-theoretic problems and their applications. Other recent results on data streams can be found in [4, 13, 14, 6].

**Related Work on Clustering** In this paper we shall consider models in which clusters have a distinguished point, or “center.” In the  $k$ -Median problem, the objective is to minimize the average distance from data points to their closest cluster centers. The 1-median problem was first posed by Weber [17]. In the  $k$ -Center problem, the objective is to minimize the maximum radius of a cluster. The above problems are all NP-hard, so we will be concerned with approximation algorithms. We will assume that the domain space of points is discrete, i.e., the cluster centers must be among the input points. The continuous case is related to the discrete problem by small factors (see Theorem 2.1). Throughout the paper we also assume that the input points are drawn from a metric space.

In the recent past, several approximation algorithms have been proposed for the  $k$ -Median problem [3, 10, 2]. These algorithms require  $O(n^2)$  space to compute the dual variables or primal constraints. We will be interested in algorithms which use more than  $k$  medians but run in linear space [12, 2, 9].

Charikar, Chekuri, Feder, and Motwani [1] gave a constant-factor algorithm for the incremental  $k$ -Center problem, which is also a single-pass algorithm requiring  $O(nk \log k)$  time and  $O(k)$  space. There is a large difference, however, between the  $k$ -Center and the  $k$ -Median problem since a set of  $k + 1$  suitably separate points provides a lower bound for the  $k$ -Center problem. These points can be thought of as a proof of the goodness of the clustering. For the  $k$ -Median problem, allowing weighted points, no such succinct proof exist and the optimization problem takes on a more global character.

**Our Results** We begin by giving an algorithm that requires small space, and then later address the issue of clustering in one pass. In Section 2 we give a simple algorithm based on divide-and-conquer that achieves a constant-factor approximation in small space. Elements of the algorithm and its analysis form the basis for the constant-factor algorithm given in Section 3. This algorithm runs in time  $O(n^{1+\epsilon})$ , uses  $O(n^\epsilon)$  memory, and makes a single pass over the data. Next, in Section 4, using randomization, we show how to reduce the running time to  $O(nk)$  without requiring more than a single pass. In Section 5 we show it is not possible to obtain any bounded approximation ratio in deterministic  $o(nk)$  time; we also show how to achieve a

poly-log  $n$  approximation ratio in a single pass in deterministic  $\tilde{O}(nk)$  time.

## 2 Clustering in Small Space

One of the first requisites of clustering a data stream is that the computation be carried out in small space. Our first goal will be to show that clustering can be carried out in small ( $n^\epsilon$  for  $n$  data points) space, without being concerned with the number of passes. Subsequently we will see how to implement the algorithm in one pass.

In order to cluster in small space, we investigate algorithms that examine the data in a piecemeal fashion. In particular, we study the performance of a divide-and-conquer algorithm, called Small-Space, that divides the data into pieces, clusters each of these pieces, and then again clusters the centers obtained (where each center is weighted by the number of points closer to it than to any other center). We show that this piecemeal approach is good, in that: if we had a constant-factor approximation algorithm, running it in divide-and-conquer fashion would still yield a (slightly worse) constant-factor approximation. We then propose another algorithm (Smaller-Space) that is similar to the piecemeal approach except that instead of reclustering only once, it repeatedly reclusters weighted centers. For this algorithm, we prove that if we recluster a constant number of times, a constant-factor approximation is still obtained, although, as expected, the constant factor worsens with each successive reclustering. The advantage of Small(er)-Space is that we sacrifice somewhat the quality of the clustering approximation to obtain an algorithm uses much less memory.

### 2.1 Simple Divide-and-Conquer and Separability Theorems

We start with the version of the algorithm that reclusters only once. Elements of the algorithm and its analysis will be used in a black-box manner in the algorithms in the rest of the paper.

#### Algorithm Small-Space(S)

Divide  $S$  into  $l$  disjoint pieces  $\chi_1, \dots, \chi_l$ .

For each  $i$ , find  $O(k)$  centers in  $\chi_i$ . Assign each point in  $\chi_i$  to its closest center.

Let  $\chi'$  be the  $O(lk)$  centers obtained in (2), where each center  $c$  is weighted by the number of points assigned to it.

Cluster  $\chi'$  to find  $k$  centers.

Since we are interested in clustering in small space,  $l$  will be set so that both  $S$  and  $\chi'$  fit in main memory, if possible. If  $S$  is very large, no such  $l$  may exist – we will address this issue later.

Before analyzing algorithm Small-Space, we describe the relationship between the discrete and continuous clustering problem. The following is folklore and is included for completeness.

**Theorem 2.1** *Given an instance of the  $k$ -median problem with a solution of cost  $C$ , where the medians may not belong to the set of input points, there exists a solution of cost  $2C$  where all the medians belong to the set of input points.*

**Proof:** Consider the solution of cost  $C$ , and let the points  $j_1, \dots, j_q$  be assigned to median  $i$ . Since median  $i$  may not be in the input, consider the point  $j_l$  which is closest to  $i$  as the median (instead of  $i$ ). Thus the assignment distance of every point  $j_r$  at most doubles, since  $c_{j_r, j_l}$  can be bounded by  $c_{j_r, i} + c_{j_r, i}$  (where  $c_{xy}$  denotes the distance from  $x$  to  $y$ ). Over all  $n$  points in the original set, the assignment distance can at most double, summing to at most  $2C$ .  $\square$

The following separability theorem sets the stage for a divide-and-conquer algorithm. This theorem carries over to other clustering metrics such as the sum of squared distances.

**Theorem 2.2** *Consider any set of  $n$  points arbitrarily partitioned into disjoint sets  $\chi_1, \dots, \chi_\ell$ . The sum of the optimum solution values for the  $k$ -median problem on the  $\ell$  sets of points is at most twice the cost of the optimum  $k$ -median problem solution for all  $n$  points.*<sup>1</sup>

**Proof:** Consider the medians used for the optimum  $k$ -median solution. If each partition uses these medians, the cost of the solution will be exactly the cost of the optimal solution. This follows since the objective function for  $k$ -median is the sum of distances to the nearest median for every point. However the set of medians chosen by the optimum solution need not be present in a partition. But in the case where the medians points can be arbitrary points in the space, the above theorem is proved.

In case we have to choose the medians from the given set of points, the medians used by the optimum solution will not be available to every partition. In this

<sup>1</sup>The factor 2 is avoided in the Euclidean case if we allow that medians can be arbitrary points in space, rather than requiring that they be points from the original data set.

case use Theorem 2.1 to construct a solution which is at most 2 times the cost of the optimum solution.  $\square$

Next we show that the new instance, where all the points  $i$  that have median  $i'$  shift their weight to the point  $i'$  (i.e., the weighted  $O(lk)$  centers  $S'$  in Step 2 of Algorithm Small-Space), has a good feasible clustering solution. Notice that the set of points in the new instance is much smaller and may not even contain the medians of the optimum solution.

**Theorem 2.3** *If the sum of the costs of the  $l$  optimum  $k$ -median solutions for  $\chi_1, \dots, \chi_l$  is  $C$  and if  $C^*$  is the cost of the optimum  $k$ -median solution for the entire set  $S$ , then there exists a solution of cost at most  $2(C + C^*)$  to the new weighted instance  $\chi'$ .*<sup>2</sup>

**Proof:** As in the proof of the previous theorem, we will consider the  $k$  medians in the optimum continuous solution.

Let the median to which  $i'$  is assigned to in the optimum continuous solution for  $\chi'$  be  $\tau(i')$ . Further, let  $d_{i'}$  be the number of points assigned to the median  $i'$ . The cost of  $\chi'$  can be expressed as  $\sum_{i'} c_{i'\tau(i')} d_{i'}$  (where again  $c_{xy}$  is the distance from  $x$  to  $y$ ). Each point  $i'$  in the new instance  $\chi'$  can be viewed as a collection of points, namely those points  $i$  assigned to the median  $i'$ . Thus the cost of  $\chi'$  can also be expressed as  $\sum_i c_{i'\tau(i')}$ .

Let the median to which  $i$  is assigned to in the optimum continuous solution for  $S$  be  $\sigma(i)$ . The cost of the new instance  $\chi'$  is no more than  $\sum_i c_{i'\sigma(i)}$  since  $\tau$  is optimum for  $\chi'$ . This sum is in turn bounded by  $\sum_i (c_{i'i} + c_{i\sigma(i)})$ . The first term summed over all points  $i$  evaluates to  $C$  and the second term evaluates to  $C^*$ .

Thus we showed an assignment to the medians of the optimal solution at cost  $C + C^*$ . Using Theorem 2.1, the theorem follows. (Note that the theorem can also be shown to hold when the original points in  $S$  are weighted.)  $\square$

We now show that if we run a bicriteria  $(a, b)$ -approximation algorithm (where at most  $ak$  medians are output with cost at most  $b$  times the optimum  $k$ -Median solution) in Step 2 of Algorithm Small-Space and we run a  $c$ -approximation algorithm in Step 4, then the resulting approximation by Small-Space can be suitably bounded.

**Theorem 2.4** *The algorithm Small-Space has an approximation factor of  $2c(1 + 2b) + 2b$ .*

**Proof:** Let the optimal  $k$ -median solution be of cost  $C^*$ . Then the cost of the solution  $C$  at the end of the

<sup>2</sup>Again, the factor 2 is avoided if we use the Euclidean distance and allow medians to be arbitrary points.

first stage is at most  $2bC^*$ . This is true due to Theorem 2.2, since we are adding the cost of the solutions to each partition, each of which is a  $b$ -approximation for that partition. Now by Theorem 2.3, there exists a solution to the  $k$ -median problem on the modified instance of cost  $2(C + C^*)$ . Since we have a  $c$ -approximation, we have a solution of cost  $2c(1 + 2b)C^*$  to the modified instance. The theorem is obtained by summing the two costs.  $\square$

The black-box nature of this algorithm will allow us to devise a new divide-and-conquer algorithm.

## 2.2 Divide-and-Conquer Strategy

We now generalize Small-Space so that the algorithm recursively calls itself on a successively smaller set of weighted centers.

### Algorithm Smaller-Space(S,i)

Divide  $S$  into  $l$  disjoint pieces  $\chi_1, \dots, \chi_l$ .

For each  $i$ , find  $O(k)$  centers in  $\chi_i$ . Assign each point in  $\chi_i$  to its closest center.

Let  $\chi'$  be the  $O(lk)$  centers obtained in (2), where each center  $c$  is weighted by the number of points assigned to it.

Call Algorithm Smaller-Space( $\chi', i - 1$ ).

We can claim the following.

**Theorem 2.5** *For constant  $i$ , Algorithm Smaller-Space( $S, i$ ) gives a constant-factor approximation to the  $k$ -Median problem.*

**Proof:** Assume that the approximation factor for the  $j$ th level is  $A_j$ . From Theorem 2.2 we know that the cost of the solution of the first level is  $2b$  times optimal. From Theorem 2.4 we get that the approximation factor  $A_j$  would satisfy a simple recurrence,

$$A_j = 2A_{j-1}(2b + 1) + 2b$$

The solution of the recurrence is  $c \cdot (2(2b + 1))^j$ . This is  $O(1)$  given  $j$  is a constant.  $\square$

Since the intermediate medians in  $\chi'$  must be stored in memory, the number of subsets  $l$  that we partition  $S$  into is limited. In particular, if the size of main memory is  $M$ , then we would need to partition  $S$  into  $l$  subsets so that each subset fits in main memory, i.e.,  $(n/l) \leq M$  and so that the weighted  $lk$  centers in  $\chi'$  also fit in main memory, i.e.,  $lk \leq M$ . Such an  $l$  may not always exist.

In the next section we will see a way to get around this problem. In fact we will be able to implement the hierarchical scheme more cleverly and obtain a clustering algorithm for an interesting model of computation. We have two themes to develop this idea. The first is to do away with the storage of the intermediate medians, and the second is to design a more interesting recursive algorithm. We take up the former and relegate the second to a later section.

## 3 The Data Stream Model

Under the Data Stream Model, computation takes place within bounded space  $M$  and the data can only be accessed via linear scans (i.e., a data point can be seen only once in a scan, and points must be viewed in order).

In this section we will modify the multi-level algorithm to operate on data streams. We will present a one-pass,  $O(1)$ -approximation in this model assuming that the bounded memory  $M$  is not too small, more specifically  $n^\epsilon$  where  $n$  denotes the size of the stream.

This model and the line of analysis have similarities to incremental clustering and online models. However our approach will be a bit different. We will maintain a forest of assignments. We will complete this to  $k$  trees, and all the nodes in a tree will be assigned to the median denoted by the root of the tree. First we will show how to solve the problem of storing intermediate medians. Next we will inspect the space requirements and running time.

**Data Stream Algorithm** To achieve this, we will modify our multi-level algorithm slightly. The algorithm will be the following:

1. Input the first  $m$  points; use a bicriterion algorithm to reduce these to  $O(k)$  (say  $2k$ ) points. As usual, the weight of each intermediate median is the number of points assigned to it in the bicriterion clustering. (Assume  $m$  is a multiple of  $2k$ .) This requires  $O(f(m))$  space, which for a primal dual algorithm can be  $O(m^2)$ . We will see a  $O(mk)$ -space algorithm later.
2. Repeat the above till we have seen  $m^2/(2k)$  of the original data points. At this point we have  $m$  intermediate medians.
3. Cluster these  $m$  first-level medians into  $2k$  second-level medians and proceed.
4. In general, maintain at most  $m$  level- $i$  medians, and, on seeing  $m$ , generate  $2k$  level- $i + 1$  medians, with the weight of a new median as the sum of

the weights of the intermediate medians assigned to it.

5. When we have seen all the original data points (or we want to have a clustering of the points we have seen so far) we cluster all the intermediate medians into  $k$  final medians.

Note that this algorithm is identical to the multi-level algorithm described before.

The number of levels required by this algorithm is at most  $O(\log(n/m)/\log(m/k))$ . If we have  $k \ll m$  and  $m = O(n^\epsilon)$  for some constant  $\epsilon < 1$ , we have an  $O(1)$ -approximation. Using linear programming or primal dual algorithms we will have  $m = \sqrt{M}$  where  $M$  is the memory size (ignoring factors due to maintaining intermediate medians of different levels). We argued that the number of levels would be a constant when  $m = n^\epsilon$  and hence when  $M = n^{2\epsilon}$  for some  $\epsilon < 1/2$ .

**Linear Space Clustering** The approximation quality which we can prove (and intuitively the actual quality of clustering obtained on an instance) will depend heavily on the number of levels we have. From this perspective it is profitable to use a space-efficient algorithm.

We can use the local search algorithm in [2] to provide a bicriterion approximation in *space linear in  $m$* , the number of points clustered at a time. The advantage of this algorithm is that it maintains only an assignment and therefore uses linear space. However the complication is that for this algorithm to achieve a bounded bicriterion approximation, we need to set a “cost” to each median used, so that we penalize if many more than  $k$  medians are used. The algorithm solves a facility location problem after setting the cost of each median to be used. However this can be done by guessing this cost in powers of  $(1 + \gamma)$  for some  $0 < \gamma < 1/6$  and choosing the best solution with at most  $2k$  medians. In the last step, to get  $k$  medians we use a two step process to reduce the number of medians to  $2k$  and then use [10, 2] to reduce to  $k$ . This allows us to cluster with  $m = M$  points at a time provided  $k^2 \leq M$ .

**The Running Time** The running time of this clustering is dominated by the contribution from the first level. The local search algorithm is quadratic and the total running time is  $O(n^{1+\epsilon})$  where  $M = n^\epsilon$ . We argued before, however, that  $\epsilon$  will not be very small and hence the approximation quality which we can prove will remain small.

We therefore claim the following theorem,

**Theorem 3.1** *We can solve the  $k$ -Median problem on a data stream with time  $O(n^{1+\epsilon})$  and space  $\theta(n^\epsilon)$  up to a factor  $2^{O(\frac{1}{\epsilon})}$ .*

We have two avenues to pursue. The running time will be lower-bounded by the space we require, and we improve this bottleneck to get linear space clustering, but first, to achieve scalability, our goal will be to get clustering in time  $\tilde{O}(nk)$ . This will mean an amortized update of  $O(k \text{ polylog}(n))$ . In the next section we will motivate how to achieve this, and provide evidence that ours is a hard bound for the running time of a clustering algorithm.

The second issue is to present an algorithm with approximation guarantee which is polynomial in  $\frac{1}{\epsilon}$ . We will show how to achieve this in Section 5.

## 4 Clustering Data Streams in $\tilde{O}(nk)$ Time

Let us recall the algorithm we have developed so far. We have  $k^2 \ll M$ , and we are applying an alternate implementation of a multi-level algorithm.

We are clustering  $m = O(M)$  (assuming  $M = O(n^\epsilon)$  for constant  $\epsilon > 0$ ) points and storing  $2k$  medians to “compress” the description of these data points. We use the local search-based algorithm in [2]. We keep repeating this procedure till we see  $m$  of these descriptors or intermediate medians and compress them further into  $2k$ . Finally, when we are required to output a clustering, we compress all the intermediate medians (over all the levels there will be at most  $O(M)$  of them) and get  $O(k)$  penultimate medians which we cluster into exactly  $k$  using the primal dual algorithm as in [10, 2].

### 4.1 Earlier Work on Clustering in $\tilde{O}(nk)$ Time

We will use the results in [9] on metric space algorithms that are subquadratic. The algorithm as defined will consist of two passes and will have constant probability of success. For high probability results, the algorithm will make  $O(\log n)$  passes. As stated, the algorithm will only work if the original data points are *unweighted*. Consider the following algorithm:

1. Draw a sample of size  $s = \sqrt{nk}$ .
2. Find  $k$  medians from these  $s$  points using the primal dual algorithm in [10].
3. Assign each of the  $n$  original points to its closest median.
4. Collect the  $n/s$  points with the largest assignment distance.

5. Find  $k$  medians from among these  $n/s$  points.
6. We have at this point  $2k$  medians.

**Theorem 4.1** [9] *The above algorithm gives an  $O(1)$  approximation with  $2k$  medians with constant probability.*

The above algorithm<sup>3</sup> provides a constant-factor approximation for the  $k$ -Median problem (using  $2k$  medians) with constant probability. Repeat the above experiment  $O(\log n)$  times for high probability. We will not run this algorithm by itself, but as a substep in our algorithm. The algorithm requires  $\tilde{O}(nk)$  time and space. Using this algorithm with the local search tradeoff results in [2] reduces the space requirement to  $O(\sqrt{nk})$ .

Alternate sampling-based results exist for the  $k$ -Median measure that do extend to the weighted case [15], however these results assume Euclidean space.

## 4.2 Extension to the Weighted Case

We need this sampling-based algorithm to work on weighted input. It is necessary to draw a random sample based on the weights of the points; otherwise the medians with respect to the sample do not convey much information. The simple idea of sampling points with respect to their weights does not help. The philosophy of the above method is that a random sample will be reasonable for most points, that there will not be many outliers (at most  $n$  divided by the sample size, up to constants), and that in the second phase it is sufficient to account for these outliers.

If the points have weights, however, in the first step we may only eliminate  $k$  points. Therefore sampling according to weights does not carry through. Contrast this with the algorithm in [5] where the points were in Euclidean space and the measure was sum of squares of distances. Both these facts were crucial for their algorithm.

We suggest the following modification. The basic idea is scaling. We can round the weights to the nearest power of  $(1 + \epsilon)$  for  $\epsilon > 0$ . In each group we can ignore the weight and lose a  $(1 + \epsilon)$  factor. Since we have an  $\tilde{O}(nk)$  algorithm, summing over all groups, the running time is still  $\tilde{O}(nk)$ . The correct way to implement this is to compute the exponent values of the weights and use only those groups which exist, otherwise the running time will depend on the largest weight.

<sup>3</sup>The algorithm presented here, without the last step, is essentially the same as in [9], however the primal dual algorithm which requires  $O(n^2)$  time to solve  $k$ -Median problem was not known when the result was published. The result proved therein was using  $O(n^2k^2)$  local search algorithm in [12] which was a bicriterion approximation.

## 4.3 The Full Algorithm

We will use this sampling-based scheme to develop a one-pass and  $O(nk)$ -time algorithm that requires only  $O(n^\epsilon)$  space.

- Input the first  $O(M/k)$  points, and use the randomized algorithm above to cluster this to  $2k$  intermediate median points.
- Use a local search algorithm to cluster  $O(M)$  intermediate medians of level  $i$  to  $2k$  medians of level  $i + 1$ .
- Use the primal dual algorithm of Jain and Vazirani [10] to cluster the final  $O(k)$  medians to  $k$  medians.

Notice that the algorithm remains one pass, since the  $O(\log n)$  iterations of the randomized subalgorithm just add to the running time. Thus, over the first phase, the contribution to the running time is  $\tilde{O}(nk)$ . Over the next level, we have  $\frac{nk}{M}$  points, and if we cluster  $O(M)$  of these at a time taking  $O(M^2)$  time, the total time for the second phase is  $O(nk)$  again. The contribution from the rest of the levels decreases geometrically, so the running time is  $\tilde{O}(nk)$ . As shown in the previous sections, the number of levels in this algorithm is  $O(\log_{\frac{M}{k}} n)$ , and so we have a constant-factor approximation for  $k \ll M = \theta(n^\epsilon)$  for some small  $\epsilon$ .<sup>4</sup>

Thus we claim the following theorem,

**Theorem 4.2** *The  $k$ -Median problem has a constant-factor approximation algorithm running in time  $O(nk \log n)$ , in one pass over the data set, using  $n^\epsilon$  memory, for small  $k$ .*

## 5 Lower Bounds and Deterministic Algorithms

In this section we explore whether our algorithms could be speeded up further and whether randomization is needed. For the former, note that we have a clustering algorithm that requires time  $\tilde{O}(nk)$  and a natural question is could we have done better? We'll show that we couldn't have done much better since a deterministic lower bound for  $k$ -Median is  $\Omega(nk)$ . Thus, modulo randomization, our time bounds pretty much match the lower bound. For the latter, we show one way to get rid of randomization that yields a single pass, small memory  $k$ -Median algorithm that is a

<sup>4</sup>We could have used the sampling-based algorithm in the intermediate steps as well, however such a recursive, sampling-based algorithm will have greater errors, in theory and very likely in practice.

poly-log  $n$  approximation. Thus we do also have a deterministic algorithm, but with more loss of clustering quality.

## 5.1 Lower Bounds

We now show that any constant-factor deterministic approximation algorithm requires  $\Omega(nk)$  time. We measure the running time by the number of times the algorithm queries the distance function.

We consider a restricted family of sets of points where there exists a  $k$ -clustering with the property that the distance between any pair of points in the same cluster is 0 and the distance between any pair of points in different clusters is 1. Since the optimum  $k$ -clustering has value 0 (where the *value* is the distance from points to nearest centers), any algorithm that doesn't discover the optimum  $k$ -clustering does not find a constant-factor approximation.

Note that the above problem is equivalent to the following Graph  $k$ -Partition Problem: Given a graph  $G$  which is a complete  $k$ -partite graph for some  $k$ , find the  $k$ -partition of the vertices of  $G$  into independent sets. The equivalence can be easily realized as follows: The set of points  $\{s_1, \dots, s_n\}$  to be clustered naturally translates to the set of vertices  $\{v_1, \dots, v_n\}$  and there is an edge between  $v_i, v_j$  iff  $\text{dist}(s_i, s_j) > 0$ . Observe that a constant-factor  $k$ -clustering can be computed with  $t$  queries to the distance function iff a graph  $k$ -partition can be computed with  $t$  queries to the adjacency matrix of  $G$ .

Kavraki, Latombe, Motwani, and Raghavan [8] show that any deterministic algorithm that finds a Graph  $k$ -Partition requires  $\Omega(nk)$  queries to the adjacency matrix of  $G$ . This result establishes a deterministic lower bound for  $k$ -Median.

**Theorem 5.1** *A deterministic  $k$ -Median algorithm must make  $\Omega(nk)$  queries to the distance function to achieve a constant-factor approximation.*

## 5.2 Deterministic Algorithms Requiring $\tilde{O}(nk)$ Time

One natural question we can ask is what we can achieve without randomization. We have already seen how to get an  $O(n^{1+\epsilon})$ -time clustering algorithm that uses  $n^\epsilon$  space and gives a constant-factor approximation. However this constant factor grows as  $2^{\frac{1}{\epsilon}}$ , and if we were to ask for an  $\tilde{O}(nk)$ -time algorithm we would have an approximation factor polynomial in  $(n/k)$ . Modifying our approach slightly, we can show the following:

**Theorem 5.2** *In  $\tilde{O}(nk)$  deterministic time, we have a poly-log  $n$  approximation for the  $k$ -Median problem in  $n^\epsilon$  space and a single pass.*

**Proof:** First we will have to construct an algorithm that runs in time  $\tilde{O}(nk)$ . Then we can reduce the space required in the same way as for the previously described randomized algorithm.

Consider the primal-dual algorithm that gives a constant-factor (say  $c$ ) approximation for the  $k$ -Median problem. This algorithm takes time (and space)  $an^2$  for some constant  $a$ . Consider the following algorithm, which we will call  $A_1$ : partition the  $n$  original points into  $p_1$  equal-size subsets, apply the primal-dual algorithm to each of these subsets, and then apply it to the  $p_1 k$  weighted points so obtained, to get  $k$  final medians. If we choose  $p_1 = (n/k)^{\frac{2}{3}}$ , the running time of  $A_1$  is  $2an^{\frac{4}{3}}k^{\frac{2}{3}}$ , and the space required is  $2an^{\frac{4}{3}}k^{\frac{2}{3}}$  also. By Theorem 2.4 we have an approximation of  $4c^2 + 4c$ .

Now define  $A_2$  to split the dataset into  $p_2$  partitions and apply  $A_1$  on each of them and on the resulting intermediate medians (notice we can easily ensure an implementation to get a one-pass algorithm). Solving to minimize the running time will yield  $p_2 = (n/k)^{4/5}$ . Therefore the running time and space required both become  $4an^{\frac{16}{15}}k^{\frac{14}{15}}$ .

If we continue this process so that  $A_i$  calls  $A_{i-1}$  on  $p_i$  partitions, we can prove without much difficulty that the running time and the space required by the algorithm will both be  $a2^i n^{\left(1 + \frac{1}{2^{2^i-1}-1}\right)} k^{\left(1 - \frac{1}{2^{2^i-1}-1}\right)}$ . However the approximation factor  $c_i$  grows as  $c_i = 4c_{i-1}^2 + 4c_{i-1}$ .

To get the exponent of  $n$  in the running time to be 1, it is sufficient to have  $i = \theta(\log \log \log n)$ . This makes the running time  $nk$  (hiding poly  $\log \log n$  factors) and gives approximation  $O(\log^p n)$  since the approximation factor is  $4^{2^i}$ . Thus we have a poly-log  $n$  approximation in  $\tilde{O}(nk)$  space and time. Now we can use this in our previous algorithm to get an  $O(\log^p n)$  approximation in  $n^\epsilon$  space and  $\tilde{O}(nk)$  time, without using randomization.  $\square$

The above actually shows that we have an  $O(n^{1+\epsilon})$ -time clustering with approximation guarantee polynomial in  $\frac{1}{\epsilon}$ . Combining this with Theorem 3.1 we get the following,

**Theorem 5.3** *The  $k$ -Median problem can be approximated in time  $\tilde{O}(n^{1+\epsilon\delta})$  and space  $\theta(n^\delta)$  up to a factor of  $O(\text{poly}(\frac{1}{\epsilon})2^{\frac{1}{\epsilon}})$ .*

## Acknowledgments

We thank Umesh Dayal, Aris Gionis, Meichun Hsu, Piotr Indyk, Dan Oblinger, and Bin Zhang for numerous fruitful discussions.

## References

- [1] M. Charikar, C. Chekuri, T. Feder and R. Motwani. Incremental clustering and dynamic information retrieval. In *Proceedings of the 29th Annual ACM Symposium on Theory of Computing*, 1997.
- [2] M. Charikar, and S. Guha. Improved Combinatorial Algorithms for the Facility Location and  $k$ -Median Problems. In *Proceedings of the 40th Annual IEEE Symposium on Foundations of Computer Science*, pages 378-388, 1999.
- [3] M. Charikar, S. Guha, É. Tardos and D. B. Shmoys. A constant factor approximation algorithm for the  $k$ -Median problem. in *Proceedings of the 31st Annual ACM Symposium on Theory of Computing*, pages 1-10, 1999.
- [4] P. Flajolet and G. N. Martin. Probabilistic Counting In *Proceedings of 24th Annual IEEE Symposium on Foundations of Computer Science*, pages 76-82, 1983.
- [5] A. Frieze, R. Kannan, and S. Vempala. Fast Monte Carlo algorithms for finding low rank approximation. In *Proceedings of the 39th Annual IEEE Symposium on Foundations of Computer Science*, pages ,1998.
- [6] J. Feigenbaum, S. Kannan, M. Strauss, and M. Vishwanathan. An approximate  $L^1$ -difference algorithm for massive data sets. In *Proceedings of 40th Annual IEEE Symposium on Foundations of Computer Science*, pages 501-511, 1999.
- [7] M. R. Henzinger, P. Raghavan, and S. Rajagopalan. Computing on Data Streams Technical Report 1998-011, Digital Equipment Corporation, Systems Research Center, May 1998.
- [8] L. E. Kavradi, J. C. Latombe, R. Motwani, and P. Raghavan. Randomized query processing in robot path planning. In *Journal of Computer and System Sciences* Special issue, vol 57, pages 50-60, 1998.
- [9] P. Indyk, Sublinear time algorithms for metric space problems. In *Proceedings of the 31st Annual ACM Symposium on Theory of Computing*, pages 428-434, 1999.
- [10] K. Jain and V. Vazirani, Primal-Dual Approximation Algorithms for Metric Facility Location and  $k$ -Median Problems. In *Proceedings of the 40th Annual IEEE Symposium on Foundations of Computer Science*, pages 1-10, 1999.
- [11] V. Kann, S. Khanna, J. Lagergren, and A. Panconesi. On the hardness of approximating MAX  $k$ -cut and its dual.
- [12] M. R. Korupolu, C. G. Plaxton, and R. Rajaraman. Analysis of a local search heuristic for facility location problems. In *Proceedings of the 9th Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 1-10, 1998.
- [13] G. S. Manku, S. Rajagopalan, and B. Lindsay. Approximate medians and other quantiles in one pass with limited memory. In *Proceedings of the 1998 ACM SIGMOD International Conference on Management of Data*, pages 426-435, 1998.
- [14] G. S. Manku, S. Rajagopalan, and B. Lindsay. Random sampling techniques for space efficient online computation of order statistics of large databases. In *Proceedings of the 1999 ACM SIGMOD International Conference on Management of Data*, pages 251-262, 1999.
- [15] N. Mishra, D. Oblinger, and L. Pitt. Way-Sublinear Time Approximate (PAC) Clustering. Manuscript, 2000.
- [16] J. I. Munro and M. S. Paterson. Selection and Sorting with Limited Storage. *Theoretical Computer Science*, vol 12, pages 315-323, 1980.
- [17] A. Weber. Ueber den Standort der Industrien. Erster Teil. Reine Theorie der Standorte. Mit einem mathematischen Anhang von G.PICK. (in German). Verlag, J. C. B. Mohr, Tbingen, Germany, 1909.