

# Data reduction for weighted and outlier-resistant clustering

Dan Feldman\*

Leonard J. Schulman†

## Abstract

Statistical data frequently includes outliers; these can distort the results of estimation procedures and optimization problems. For this reason, loss functions which deemphasize the effect of outliers are widely used by statisticians. However, there are relatively few algorithmic results about clustering with outliers.

For instance, the *k*-median with outliers problem uses a loss function  $f_{c_1, \dots, c_k}(x)$  which is equal to the minimum of a penalty  $h$ , and the least distance between the data point  $x$  and a center  $c_i$ . The loss-minimizing choice of  $\{c_1, \dots, c_k\}$  is an outlier-resistant clustering of the data. This problem is also a natural special case of the *k*-median with penalties problem considered by [Charikar, Khuller, Mount and Narasimhan SODA'01].

The essential challenge that arises in these optimization problems is data reduction for the *weighted k*-median problem. We solve this problem, which was previously solved only in one dimension ([Har-Peled FSTTCS'06], [Feldman, Fiat and Sharir FOCS'06]). As a corollary, we also achieve improved data reduction for the *k*-line-median problem.

## 1 Introduction

**1.1 Weighted optimization problems** We show how to perform data reduction for a variety of problems in optimization and statistical estimation. The problems are of the following form: a metric space  $(M, \text{dist})$  and a family of functions  $F$  are specified. Then, given a set  $P$  of  $n$  points in  $M$ , the optimization problem is to find an  $f$  which is a  $(1 + \varepsilon)$ -approximate minimizer of  $f(P)$  among all  $f \in F$ , where  $f(P) = \sum_{p \in P} f(p)$ . We focus particularly on families  $F$  which are appropriate for minimization of a loss function (e.g., max likelihood estimation) in statistical inference. A key case we treat pertains to the problem of clustering with outliers:

***k*-median with outliers in  $\ell_2^d$ :** ( $\ell_2^d$  is  $\mathbb{R}^d$  with the Euclidean metric.) Here one is interested in modeling data as being distributed about  $k$  centers, with points that are beyond a threshold distance  $h$  being considered

as outliers.

$$F_{d,k}^{\text{out}} : \mathbb{R}^d \rightarrow \mathbb{R}_+ \quad \text{where } \mathbb{R}_+ = \text{nonnegative reals}$$
$$F_{d,k}^{\text{out}} = \{f_{C,h}\}_{C \subseteq \mathbb{R}^d, |C|=k, h>0}$$
$$f_{C,h}(p) = \min\{h, \min_{c \in C} \|p - c\|\}$$

Note that, conditional on a point being treated as an outlier (assigned value  $h$ ), it has no further effect on the cost-minimizing choice of the centers  $C$ . This way of formalizing the treatment of outliers is a slightly simplified form of the *Tukey biweight* loss function used by statisticians to perform outlier-resistant estimation. It is also a special case of the *k*-median with penalties problem considered by Charikar, Khuller, Mount and Narasimhan [13] (the distinction being that they allow each point  $p$  a custom penalty  $h(p)$  for being an outlier).

*Data reduction* means the replacement of the input  $P$  (implicitly, the uniform probability distribution on  $P$ ) by a probability distribution  $\nu$  on a much smaller set  $A$ , such that for all  $f \in F$ ,

$$f(\nu) = (1 \pm \varepsilon)f(P)/n, \quad \text{where } f(\nu) = \sum_{p \in A} \nu(p)f(p).$$

$(A, \nu)$  is often called an  $\varepsilon$ -approximation or core-set for the data  $P$  w.r.t. the family of functions  $F$ . In recent years a substantial body of work has gone into providing (or showing non-existence of) data reduction for various problems, because one can then run a relatively inefficient optimization algorithm (possibly even exhaustive search) on the core-set.

**1.2 Our Results** We show that for any constant  $k$ , we can efficiently construct core-sets of cardinality  $O((\Delta \log^2 n)/\varepsilon^2)$  for certain types of  $k$ -clustering problems;  $\Delta$  is a Vapnik-Chervonenkis-type measure of the combinatorial complexity of the clustering problem. (If  $P$  is in  $\ell_2^d$  then  $\Delta \in O(d)$ ; if  $P$  is in a finite metric space then  $\Delta \in O(\log n)$ .)

These clustering problems include the well-known *k*-median, *k*-means, etc., but go beyond these to include treatment of outliers, and “*M*-estimators” in robust statistics such as the Huber and Tukey loss functions.

The key obstacle we overcome, which has not been overcome previously except in one dimension, is that

\*Caltech, Pasadena CA 91125. Email: dannyf@caltech.edu.

†Caltech, Pasadena CA 91125. Email: schulman@caltech.edu.

Work supported in part by the NSF.

of handling “weights” on the clustering centers. The ability to handle weights and non-increasing distance functions is what allows us to provide core-sets for the outliers family  $F_{d,k}^{\text{out}}$  defined above.

**Construction time.** Our algorithm is randomized and runs in time linear in the input size. As observed in [21], core-set construction can be derandomized using deterministic algorithms for computing  $\varepsilon$ -approximations. Matoušek provided a general algorithm for computing such approximations in time linear in  $n$  but exponential in  $\Delta$  [38]. Indeed, these techniques can also be applied to our core-sets; see [21] for more details. For the special case  $d = 1$  (points on a line) we can construct deterministic  $\varepsilon$ -approximations by choosing a single representative from every subset of  $\varepsilon n$  consecutive input points. This yields corresponding core-sets of size only linear in  $\log^2(n)/\varepsilon$  for the case  $d = 1$ , which is smaller than the output of the randomized algorithms.

**Streaming and distributed computations** Using the map-reduce technique [28], our core-sets imply polylogarithmic space and polylogarithmic update-time algorithms for clustering streaming data with outliers; this is apparently the first result of this type. Similarly, the techniques can be adapted to parallelization [23, 5].

We summarize our results in Tables 1 and 2 and in the following paragraphs. The following first example explains what we mean by weights and is also, from the mathematical point of view, perhaps the central example to keep in mind:

**Weighted  $k$ -median in  $\ell_2^d$ :** Here one is interested in modeling possible heterogeneity among cluster centers. This is natural, for example, in the context of *mixture models* in which the components of the mixture have varying standard deviations.

$$F_{d,k}^{\text{wght}} : \mathbb{R}^d \rightarrow \mathbb{R}_+$$

$$F_{d,k}^{\text{wght}} = \{f_C\}_{C \subseteq \mathbb{R}^d \times \mathbb{R}_+, |C|=k}$$

$$f_C(p) = \min_{(c,w) \in C} w \cdot \|p - c\|$$

Our approach is more general than is implied by these two examples, but is somewhat technical and is deferred to Section 2; our main theorem is Theorem 5.1 and example application of it for robust statistics can be found in Corollary 5.2.

**$k$ -line median in  $\ell_2^d$ .** Heterogeneity also arises naturally from geometric considerations, in the reduction of the (unweighted)  $k$ -line-median problem to weighted  $k$ -median on a line ( $d = 1$ ) [27, 21]. Hence, the above core-sets for weighted centers immediately imply the smallest known core-sets for  $k$  lines in  $\mathbb{R}^d$ .

**Robust Statistics ( $M$ -estimators).** As described, our approach provides core-sets for (at least) two of the most important outlier-resistant statistical estimators. Huber’s estimator is used very widely [25, 32]; Zhang [50] writes that “this estimator is so satisfactory that it has been recommended for almost all situations”. Hardin et al. [29] write that “Tukey’s biweight has been well established as a resistant measure of location and scale for multivariate data [44, 32, 31]”. Both of these estimators are types of  $M$ -estimators; very little is known about the computational complexity of optimizing  $M$ -estimators [42, 41, 44, 29, 25, 32] and our paper shows how to make improvement in this direction.

Formally, for a distance threshold  $h$  for outliers, a point  $p \in P$  and a given query center  $c \in \mathbb{R}^d$ , with  $x = \|p - c\|$ , we define the cost functions

$$f_c^{h,\text{Huber}}(p) = \min\{x^2/2, h(x - h/2)\}$$

$$f_c^{h,\text{Tukey}}(p) = \min\{h^2/6, \frac{3h^4x^2 - 3h^2x^4 + x^6}{6h^4}\}$$

The corresponding core-sets for the families of these functions are of size  $O(\Delta\varepsilon^{-2} \log^2 n)$ . See Sec. 5 for the more general version with a set of  $k$ -weighted centers.

**Robust  $k$ -median with  $m$  outliers.** Our method also enables an approach to the Robust  $k$ -median with  $m$  outliers problem: discarded outliers can be treated simply as infinite-weight centers, so our core-sets can handle a constant number of discarded outliers. This causes an exponential dependence of the size and runtime on  $m$ , but is still the only known near-linear-time (in  $n$ )  $(1 + \varepsilon)$ -approximation for the problem, even for  $k = 1$  and  $d = 2$ .

### 1.3 Literature

**Sampling literature.** Data reduction by uniform sampling goes back to the foundations of statistics; the most relevant line of work for our purpose is that initiated by Vapnik and Chervonenkis [49, 43, 30, 12, 48, 40, 33, 8, 9, 3, 7] (and see [45, 47]). However, for estimation of general nonnegative (esp., unbounded) loss functions, and for the design of approximation algorithms for (related) optimization problems, it is essential to design *weighted* sampling methods. This is a more recent line of work, beginning at least with [10, 15, 34, 4, 46]. There are also methods for deterministic data reduction [39, 18, 21] but the results are generally weaker and we shall not emphasize this aspect of the problem in the paper.

**Clustering literature.** The  $k$ -median problem was shown to be NP-hard by a reduction from dominating set [36]. This problem is a special case of  $k$ -clustering problems with various exponents  $r > 0$ , with loss

Optimization Problem	Metric	$\ell_r$ loss	Approx.	Time	Ref.
$k$ -Median with penalties	Arbitrary	$r = \infty$	3	$O(n^3)$	[13]
$k$ -Median with penalties	Arbitrary	$r = 1$	4	$n^{3+O(1)}$	[13]
$k$ -Median with penalties	Arbitrary	$r \in O(1)$	$O(1)$	$k^{O(k)}n \log(n) + nk^{O(k)} \log^2 n$	★★
$k$ -Median with penalties	$\mathbb{R}^d$	$r \in O(1)$	$1 + \varepsilon$	$ndk^{O(k)} + (k\varepsilon^{-2} \log(n))^k$	★★
M-Estimators	Arbitrary	$r \in O(1)$	Heuristics	?	e.g. [29]
M-Estimators	$\mathbb{R}^d$	$r \in O(1)$	$1 + \varepsilon$	$ndk^{O(k)} + (k\varepsilon^{-2} \log(n))^k$	★★
M-Estimators	Arbitrary	$r \in O(1)$	$O(1)$	$k^{O(k)}n \log(n) + nk^{O(k)} \log^2 n$	★★
Robust $k$ -Median with $m$ outliers	Arbitrary	$r = 1$	$O(1)$	$O(k^2(k+m)^2n^3 \log n)$	[14]
Robust 2-Median with $m$ outliers	$\mathbb{R}^2$	$r = \infty$	1	$O(nm^7 \log^3 n)$	[1]
Robust 4-Median with $m$ outliers	$\mathbb{R}^2$	$r = \infty$	1	$nm^{O(1)} \log n$	[1]
Robust 5-Median with $m$ outliers	$\mathbb{R}^2$	$r = \infty$	1	$nm^{O(1)} \log^5 n$	[1]
Robust $k$ -Median with $m$ outliers	Arbitrary	$r = \infty$	3	$O(n^3)$	[13]
Robust $k$ -Median with $m$ outliers	$\mathbb{R}^d$	$r \in O(1)$	$1 + \varepsilon$	$nd(m+k)^{O(m+k)} + (\varepsilon^{-1}k \log n)^{O(1)}$	★★
Robust $k$ -Median with $m$ outliers	Arbitrary	$r \in O(1)$	$O(1)$	$n \log(n)(m+k)^{O(m+k)} + (k \log n)^{O(1)}$	★★
$k$ -Line Median	$\mathbb{R}^d$	$r = 1, 2$	$1 + \varepsilon$	$nd(k/\varepsilon)^{O(1)} + d(\log d)^{(k/\varepsilon)^{O(1)}}$	[16]
$k$ -Line Median	$\mathbb{R}^d$	$r = 1$	$1 + \varepsilon$	$ndk^{O(1)} + (\varepsilon^{-d} \log n)^{O(dk)}$	[27]
$k$ -Line Median	$\mathbb{R}^d$	$r = 2$	$1 + \varepsilon$	$ndk^{O(1)} + (\varepsilon^{-d} \log n)^{O(dk^2)}$	[20]
$k$ -Line Median	$\mathbb{R}^d$	$r = 1, 2$	$1 + \varepsilon$	$ndk^{O(1)} + (\varepsilon^{-1} \log n)^{O(k)}$	[21]
$k$ -Line Median	$\mathbb{R}^d$	$r \in O(1)$	$1 + \varepsilon$	$ndk^{O(1)} + \varepsilon^{-2} \log(n) \cdot k^{O(k)}$	★★

Table 1: Approximation Algorithms. The input is a set  $P$  of  $n$  points in  $\mathbb{R}^d$  or in an arbitrary metric space. The results of this paper are marked with ★★.

Core-set	Metric	$\ell_r$ loss	Size	Ref.
Weighted $k$ -median	$\mathbb{R}^1$	$r = 1$	$(\varepsilon^{-1} \log n)^{O(k)}$	[27]
Weighted $k$ -median	$\mathbb{R}^1$	$r = 2$	$(\varepsilon^{-1} \log n)^{O(k^2)}$	[20]
Weighted $k$ -median	$\mathbb{R}^1$	$r = \infty$	$(k/\varepsilon)^{O(k)}$	[2]
Weighted $k$ -median	$\mathbb{R}^d$	$r = \infty$	$O(k!/\varepsilon^{dk})$	[26]
Weighted $k$ -median	$\mathbb{R}^d$ /Arbitrary	$r \in O(1)$	$k^{O(k)}(\varepsilon^{-1}d \log n)^2$	★★
$k$ -Line median	$\mathbb{R}^d$	$r = 1, 2$	$dk\varepsilon^{-2} + (\varepsilon^{-1} \log n)^{O(k^2)}$	[21]
$k$ -Line median	$\mathbb{R}^d$	$r \in O(1)$	$dk\varepsilon^{-2} + k^{O(k)}(\varepsilon^{-1} \log n)^2$	★★
$k$ -Median with penalties	$\mathbb{R}^d$ /Arbitrary	$r \in O(1)$	$k^{O(k)}(\varepsilon^{-1}d \log n)^2$	★★
Robust $k$ -median with $m$ outliers	$\mathbb{R}^d$ /Arbitrary	$r \in O(1)$	$(k+m)^{O(k+m)}(\varepsilon^{-1}d \log n)^2$	★★
M-Estimators	$\mathbb{R}^d$ /Arbitrary	$r \in O(1)$	$k^{O(k)}(\varepsilon^{-1}d \log n)^2$	★★

Table 2: Core-sets. The input is a set  $P$  of  $n$  points in  $\mathbb{R}^d$  or in an arbitrary metric space. New results of this paper are marked with ★★. We denote  $d = O(\log n)$  for the case of an arbitrary metric space.

function  $f_C(p) = \min_{c \in C} \text{dist}(p, c)^r$  for centers  $C = \{c_1, \dots, c_k\}$ . The  $k$ -means problem (exponent  $r = 2$ ) is NP-hard even for  $k = 2$  [17] or in the Euclidean plane [37]. The case  $r = \infty$  refers to the  $k$ -center problem  $f_C(P) = \max_{p \in P} \min_{c \in C} \text{dist}(p, c)$ ; it is NP-hard to approximate this to within a factor of 1.822 even in the Euclidean plane [19].

The current best approximation guarantee for  $k$ -median in general metrics is  $(3 + \varepsilon)$  [6]. When  $k$  is fixed, [22] provided a weak core-set of size independent of  $d$  for  $k$ -means that yields an algorithm that takes time  $O(nd) + (k/\varepsilon^2)^{O(k/\varepsilon)}$ . (A weak core-set is sufficient for optimization but not for evaluation of general queries.) Recently, this result was generalized and improved for any constant  $r \geq 1$ , with weak core-sets of size only linear in  $k$  [21]. Strong core-sets of size  $(dk)^{O(1)}$  for the  $k$ -median problem for any constant  $r > 0$  were provided in [35].

In the  $k$ -median with penalties problem [13], for each input point we may decide to either provide service and pay the service cost to its nearest center, or to pay the penalty. Setting all penalties to 1 gives the standard notion, which has also been studied earlier in the context of TSP and Steiner trees, see [24, 11] and references therein. As mentioned above, this is precisely our  $F_{d,k}^{\text{out}}$  problem; it is also very close to clustering with Tukey loss, see Sec. 5.

An alternative approach to handling outliers is the *robust  $k$ -median with  $m$  outliers problem* due to [13]. Here there is, besides the usual  $k$ -median formulation, an additional parameter  $m$  which is the number of points we are allowed to “discard”. The problem is to place the  $k$  centers so as to minimize the sum, over the best set of  $n - m$  data points, of the distance to the closest center. This is a less “continuous” way of treating outliers and, correspondingly,  $m$  enters significantly into the time complexities of algorithms. Our weighted- $k$ -median algorithm can be used to address this problem, see Sec. 5. [13] also considered relaxing the number of discarded points, and provided a polynomial time algorithm that outputs a  $k$ -clustering serving  $(1 - \varepsilon)(n - m)$  points with cost within  $4(1 + 1/\varepsilon)$  times the optimum cost (for  $n - m$  points). Recently, [21] improved the running time for this problem to linear in  $n$  by showing that an  $\varepsilon$ -approximation of  $P$  for  $k$  balls (in particular, a small uniform random sample) is a core-set for this problem.

The *weighted  $k$ -median* problem was introduced in Har-Peled [27]; that paper provided an  $O((\log n)^k)$ -size core-set for this problem in one dimension, and posed the construction of core-sets in higher dimension as an open problem. The same paper proved a lower bound

of

$$\Omega(\max\{(k/\varepsilon) \log(n/k), 2^k\})$$

for the size of a core-set for weighted  $k$ -medians, even in one dimension. (We do not know a stronger lower bound in higher dimension.) Thus our results are optimal, as a function of  $n$ , up to a log factor.

In the  $k$ -line-median problem, the “centers” are actually lines in  $\mathbb{R}^d$ . This problem can be reduced to the *weighted  $k$ -median* problem in one dimension [27, 21]. Our core-set for this problem, of size  $O((\varepsilon^{-1} \log n)^2)$ , improves on the best previous  $O((\varepsilon^{-1} \log n)^{O(k)})$ . Our method also considerably simplifies, even for the one-dimensional version of the problem, the construction in [27] (which both these papers depend upon).

## 2 Preliminaries

**2.1 Loss Functions** Let  $(M, \text{dist})$  be the metric space in which our points (or data items) lie. Our framework depends upon a distortion (or “loss”) function

$$D : \mathbb{R}_+ \rightarrow \mathbb{R}_+.$$

We impose the following requirements on  $D$ :

1.  $D$  is monotone non-decreasing.
2. *Log-Log Lipschitz Condition:* There is a constant  $0 < r < \infty$  such that for all  $x, \delta > 0$ ,

$$(2.1) \quad D(xe^\delta) \leq e^{r\delta} D(x).$$

We use  $\tilde{D}$  to denote a bivariate function as follows: for  $p, q \in M$ ,  $\tilde{D}(x, y) := D(\text{dist}(p, q))$ .

LEMMA 2.1. *The conditions above imply*

- (i) For  $\phi = (4r)^r$ ,

$$(2.2) \quad \tilde{D}(p, c) - \tilde{D}(q, c) \leq \phi \tilde{D}(p, q) + \frac{\tilde{D}(p, c)}{4}$$

- (ii) *(Weak triangle inequality)* For  $\rho = \max\{2^{r-1}, 1\}$ ,

$$(2.3) \quad \tilde{D}(p, q) \leq \rho(\tilde{D}(p, c) + \tilde{D}(c, q)).$$

*Proof.* (i) Let  $x = \text{dist}(p, c), y = \text{dist}(q, c), z = \text{dist}(p, q)$ . So we are to show  $D(x) - D(y) \leq \phi D(z) + D(x)/4$ . We suppose that  $x > y$  and  $D(x) > \phi D(z)$ , otherwise the lemma is immediate. So by Eqn 2.1,  $x > z\phi^{1/r}$ .

An equivalent form of Eqn 2.1 is that for  $\delta > 0$ ,  $D(xe^{-\delta}) \geq e^{-r\delta} D(x)$ . So  $D(x) - D(y) \leq D(x) \cdot (1 - (y/x)^r)$ .

Note that for  $u \geq 0$ ,  $1 - u^r \leq r(1 - u)$ ; this follows because, viewing each side as a function of

$u$ , the two functions are tangent at  $u = 1$ , and the LHS is convex-cap while the RHS is linear. Applying this we have  $D(x) - D(y) \leq D(x) \cdot r \cdot (x - y)/x$ . Applying the triangle inequality  $x - y \leq z$  we have that  $D(x) - D(y) \leq D(x) \cdot r \cdot z/x$ . By our earlier bound this is  $< D(x) \cdot r \cdot \phi^{-1/r}$ . Plugging in  $\phi = (4r)^r$  implies Eqn 2.2.

(ii) By the triangle inequality and Eqn 2.1, for any  $0 < p < 1$ ,

$$\begin{aligned} \tilde{D}(p, q) &\leq p\tilde{D}(p, c) \left( \frac{\text{dist}(p, c) + \text{dist}(c, q)}{\text{dist}(p, c)} \right)^r \\ &\quad + (1 - p)\tilde{D}(c, q) \left( \frac{\text{dist}(p, c) + \text{dist}(c, q)}{\text{dist}(c, q)} \right)^r \\ &= (\text{dist}(p, c) + \text{dist}(c, q))^r \\ &\quad \times \left( \frac{p\tilde{D}(p, c)}{\text{dist}(p, c)^r} + \frac{(1 - p)\tilde{D}(c, q)}{\text{dist}(c, q)^r} \right). \end{aligned}$$

Substituting  $p = \text{dist}(p, c)^r / (\text{dist}(p, c)^r + \text{dist}(c, q)^r)$  we have  $\tilde{D}(p, q) \leq (\tilde{D}(p, c) + \tilde{D}(c, q)) \frac{(\text{dist}(p, c) + \text{dist}(c, q))^r}{\text{dist}(p, c)^r + \text{dist}(c, q)^r}$ . By convexity considerations, for  $r \geq 1$  the factor is maximized with  $\text{dist}(p, c) = \text{dist}(c, q)$  and for  $r \leq 1$  it is maximized with  $\text{dist}(c, q) = 0$ , yielding Eqn 2.3.

## 2.2 Tractable $(M, \tilde{D})$ Problems

**DEFINITION 2.1. (TRACTABLE  $(M, \tilde{D})$  PROBLEMS)**  
Let  $(M, \text{dist})$  be a metric space. Let  $\tilde{D}$  be a function from  $M \times M$  to  $[0, \infty)$ . We call the problem  $(M, \tilde{D})$  tractable if inequalities (2.2) and (2.3) hold for some constants  $\phi, \rho \in (0, \infty)$ .

Thus the conditions we imposed on  $D$  in Sec. 2.1 imply that  $(M, \tilde{D})$  is tractable with  $\phi = (4r)^r$ ,  $\rho = \max\{2^{r-1}, 1\}$ .

In Theorem 4.1 we show how to perform data reduction for tractable  $(M, \tilde{D})$  problems, conditional on a shatter function (essentially, VC dimension) bound.

Let  $P \subseteq M$  be a finite set of points. For  $B \subseteq M$ , we denote by  $\text{closest}(P, B, \gamma)$  the set that consists of the  $\lceil \gamma|P| \rceil$  points  $p \in P$  with the smallest values of  $\min_{q \in B} \tilde{D}(p, q)$ . For  $p \in M$  and a set  $C \subseteq M \times \mathbb{R}_+$  define  $\tilde{D}_W(p, C) = \min_{(c, w) \in C} w \cdot \tilde{D}(p, c)$ . Each  $(c, w) \in C$  is called a *weighted center*. For integer  $k \geq 1$  write  $[k] := \{1, \dots, k\}$ .

We show how to perform data reduction for a variety of statistical problems by considering appropriate choices of  $M$  and  $\tilde{D}$  and showing that the above properties are satisfied. The families of functions we consider have the following description:

$$\begin{aligned} F_{M, \tilde{D}} &: M \rightarrow \mathbb{R}_+ \\ F_{M, \tilde{D}} &= \{f_C\}_{C \subseteq M \times \mathbb{R}_+, |C|=k} \\ f(p) &= \tilde{D}_W(p, C) \end{aligned}$$

In this notation, the weighted  $k$ -median problem is  $(M = \mathbb{R}^d, \text{dist} = \text{Euclidean metric}, \text{and } \tilde{D} = \text{dist})$ ; the weighted  $k$ -mean problem is  $(M = \mathbb{R}^d, \text{dist} = \text{Euclidean metric}, \text{and } \tilde{D} = \text{dist}^2)$ ; and the  $k$ -median with outliers problem is (a special case of) the problem  $(M = \mathbb{R}^d, \text{dist} = \text{Euclidean metric}, \text{and } \tilde{D} = \min\{\text{dist}, 1\})$ .

As established in [35, 21], a sufficient condition for data reduction is that the *total sensitivity*  $\mathcal{T} = \mathcal{T}(F_{M, \tilde{D}})$  be small, and that we be able to effectively compute good upper bounds  $s(p)$  for the sensitivities of the points of  $P$ <sup>1</sup>; the cardinality of the resulting set  $A$  is then approximately  $\mathcal{T}^2 d / \varepsilon^2$ , where  $d$  is a Vapnik-Chervonenkis measure of the combinatorial complexity of the family  $F_{M, \tilde{D}}$ .

Before showing how to compute bounds on the sensitivities of points we need two more definitions.

**DEFINITION 2.2.** For a finite set  $Q \subseteq M$  and  $\gamma \in [0, 1]$ , define

$$\tilde{D}^*(Q, \gamma) := \min_{c \in M} \sum_{p \in \text{closest}(Q, c, \gamma)} \tilde{D}(p, c).$$

A point  $c$  which achieves the above minimum is, in a sense, a median of a densest region of the data. (One may also think of it as a good “median with outliers” for the data.) In what follows it would be very useful to have a subroutine to compute such a point, but this is a nearly circular request (though not quite as hard as the full goal of the paper). Instead we will be able to achieve our results using a subroutine which produces a point with the following weaker property.

**DEFINITION 2.3. (ROBUST MEDIAN)** For  $\gamma \in [0, 1]$ ,  $\tau \in (0, 1)$  and  $\alpha > 0$ , the point  $q \in M$  is a  $(\gamma, \tau, \alpha)$ -median of the finite set  $Q \subseteq M$  if

$$(2.4) \quad \sum_{p \in \text{closest}(Q, \{q\}, (1-\tau)\gamma)} \tilde{D}(p, q) \leq \alpha \cdot \tilde{D}^*(Q, \gamma).$$

## 3 Bounding point sensitivities

### 3.1 Sensitivity bound for weighted medians

The key technical advance in this paper lies in the following lemma, which shows how to translate the new definitions of the previous section into good upper bounds on the sensitivities of data points. This lemma is what enables us to handle weighted clustering problems.

In each application one needs only to ensure that the problem is “tractable” as in definition 2.1, and that the appropriate shatter function ( $\sim$  VC dimension) is bounded.

<sup>1</sup>For a family  $F$  and  $n$  data points  $P$ , the sensitivity of  $p \in P$  is  $s(p) = \sup_{f \in F} f(p) / ((1/n) \sum_{q \in P} f(q))$ ; the total sensitivity  $\mathcal{T}(F)$  is  $\sup_P \sum_{p \in P} s(p)$ .

LEMMA 3.1. Let  $(M, \tilde{D})$  be tractable and let  $P \subseteq M$  be a finite set. Suppose that  $(q_k, Q_k)$  is the output of the algorithm *Recursive-Robust-Median* $(P, k)$ . Then for every set  $C = \{(c_1, w_1), \dots, (c_k, w_k)\} \subseteq M \times [0, \infty)$  and  $p \in Q_k$  such that  $\tilde{D}_W(p, C) > 0$ , we have

$$\frac{\tilde{D}_W(p, C)}{\sum_{q \in P} \tilde{D}_W(q, C)} \leq \frac{O(k)}{|Q_k|}.$$

---

**Algorithm 1: Recursive-Robust-Median** $(P, k)$

---

**Input:** A set  $P \subseteq M$ , an integer  $k \geq 1$ .

**Output:** A pair  $(q_k, Q_k)$  that satisfies Lemma 3.1.

```

1  $Q_0 \leftarrow P$ 
2 for  $i = 1$  to  $k$  do
3   Compute a  $(1/k, \tau, \alpha)$ -median  $q_i \in M$  of
    $Q_{i-1}$  for some constants  $\tau \in (0, 1)$  and
    $\alpha \in (0, \infty)$  /* See Definition 2.3 of
    $(\gamma, \tau, \alpha)$ -median, and Algorithm 2 in
   Section 2.3 for a suggested
   implementation. */
4    $Q_i \leftarrow \text{closest}(Q_{i-1}, \{q_i\}, (1 - \tau)/(2k))$ 
5 return  $(q_k, Q_k)$ 

```

---

*Proof of Lemma 3.1:* Consider the variables  $Q_0, \dots, Q_k$  and  $q_1, \dots, q_k$  that are computed during the execution of *Recursive-Robust-Median* $(P, k)$ . A point  $p \in P$  is served by the weighted center  $(c, w) \in C$  if  $\tilde{D}_W(p, C) = w \cdot \tilde{D}(p, c)$ . For every  $i \in [k + 1]$ , let  $(c_i, w_i) \in C$  denote a center that serves at least  $|Q_{i-1}|/k$  points from  $Q_{i-1}$ . Let  $P_i$  denote the points of  $P$  that are served by  $(c_i, w_i)$ . For every  $i \in [k]$ , let

$$Q'_i := \text{closest}(Q_{i-1}, \{q_i\}, (1 - \tau)/k),$$

and  $\tilde{D}_i^* = \sum_{q \in Q'_i} \tilde{D}(q, q_i)$ .

Since  $|P_i \cap Q_{i-1}| \geq |Q_{i-1}|/k \geq |Q'_i|$ , we have by Definition 2.2,

$$(3.5) \quad \sum_{q \in P_i \cap Q_{i-1}} \tilde{D}(q, c_i) \geq \tilde{D}^*(Q_{i-1}, 1/k).$$

We prove the lemma using the following case analysis.

**Case (i):** There is an  $i \in [k]$  such that

$$(3.6) \quad \tilde{D}(p, c_i) \leq \frac{16\phi\rho\alpha \cdot \tilde{D}_i^*}{|Q'_k|}.$$

**Case (ii):** Otherwise.

**Proof of Case (i):** By (3.6) we have

$$(3.7) \quad \begin{aligned} \frac{\tilde{D}_W(p, C)}{\sum_{q \in P} \tilde{D}_W(q, C)} &\leq \frac{w_i \cdot \tilde{D}(p, c_i)}{w_i \sum_{q \in P_i} \tilde{D}_W(q, C)} \\ &\leq \frac{\tilde{D}(p, c_i)}{\sum_{q \in P_i} \tilde{D}(q, c_i)} \\ &\leq \frac{16\phi\rho\alpha \cdot \tilde{D}_i^*/|Q'_k|}{\sum_{q \in P_i \cap Q_{i-1}} \tilde{D}(q, c_i)}. \end{aligned}$$

By Definition 2.3, we have  $\tilde{D}^*(Q_{i-1}, 1/k) \geq \tilde{D}_i^*/\alpha$ . Using this with (3.5) yields

$$\sum_{q \in P_i \cap Q_{i-1}} \tilde{D}(q, c_i) \geq \tilde{D}_i^*/\alpha.$$

By the last inequality and (3.7) we obtain

$$\frac{\tilde{D}_W(p, C)}{\sum_{q \in P} \tilde{D}_W(q, C)} \leq \frac{16\phi\rho\alpha \cdot \tilde{D}_i^*/|Q'_k|}{\tilde{D}_i^*/\alpha} \leq \frac{16\phi\rho\alpha^2}{|Q_k|}.$$

**Proof of Case (ii):** By the pigeonhole principle,  $c_i = c_j$  for some  $i, j \in [k + 1]$ ,  $i < j$ . Put  $q \in P_j \cap Q_{j-1}$  and note that  $p \in Q_k \subseteq Q_{j-1}$ . Using the Markov inequality,

$$\tilde{D}(q, q_{j-1}), \tilde{D}(p, q_{j-1}) \leq \frac{2\tilde{D}_{j-1}^*}{|Q'_{j-1}|}.$$

By this, the symmetry of  $\tilde{D}(\cdot, \cdot)$  and (2.3),

$$\tilde{D}(p, q) \leq \rho(\tilde{D}(p, q_{j-1}) + \tilde{D}(q_{j-1}, q)) \leq \frac{4\rho \cdot \tilde{D}_{j-1}^*}{|Q'_{j-1}|}.$$

Using the last inequality with (2.2) yields

$$\begin{aligned} \tilde{D}(p, c_j) - \tilde{D}(q, c_j) &\leq \phi\tilde{D}(p, q) + \frac{\tilde{D}(p, c_j)}{4} \\ &\leq \frac{4\phi\rho \cdot \tilde{D}_{j-1}^*}{|Q'_{j-1}|} + \frac{\tilde{D}(p, c_j)}{4} \\ &\leq \frac{4\phi\rho\alpha \cdot \tilde{D}_i^*}{|Q'_k|} + \frac{\tilde{D}(p, c_j)}{4}. \end{aligned}$$

Since Case (i) does not hold, we have

$$16\phi\rho\alpha \cdot \tilde{D}_i^*/|Q'_k| < \tilde{D}(p, c_i) = \tilde{D}(p, c_j).$$

Combining the last two inequalities yields

$$\begin{aligned} \tilde{D}(p, c_j) - \tilde{D}(q, c_j) &< \frac{\tilde{D}(p, c_j)}{4} + \frac{\tilde{D}(p, c_j)}{4} \\ &= \frac{\tilde{D}(p, c_j)}{2}. \end{aligned}$$

That is,  $\tilde{D}(q, c_j) > \tilde{D}(p, c_j)/2$ . Hence,

$$\begin{aligned} \frac{\tilde{D}_W(p, C)}{\sum_{q \in P} \tilde{D}_W(q, C)} &\leq \frac{\tilde{D}(p, c_j)}{\sum_{q \in P_j \cap Q_{j-1}} \tilde{D}(q, c_j)} \\ &< \frac{2\tilde{D}(p, c_j)}{\tilde{D}(p, c_j) \cdot |P_j \cap Q_{j-1}|} \\ &\leq \frac{2k}{|Q_{j-1}|} \\ &\leq \frac{2k}{|Q_k|}. \end{aligned}$$

□

#### 4 Data reduction for tractable $(M, \tilde{D})$ problems

**DEFINITION 4.1.** ( $\dim(M, \tilde{D}, k)$  [49]) *Let  $(M, \tilde{D})$  be tractable. For every  $r \geq 0$  and  $C \subseteq M \times [0, \infty)$  of size  $|C| = k$ , let  $\text{ball}(C, r) = \{p \in P \mid \tilde{D}_W(p, C) \leq r\}$ .*

*Let*

$$\text{balls} = \{\text{ball}(C, r) \mid C \subseteq M \times [0, \infty), |C| = k, r \geq 0\}.$$

*The dimension  $\dim(M, \tilde{D}, k)$  is the smallest integer  $d$  such that for every finite  $S \subseteq M$  we have*

$$\left| \{S \cap \beta \mid \beta \in \text{balls}\} \right| \leq |S|^d.$$

The following is a corollary of [21, Theorem 13.1].

**COROLLARY 4.1.** *Let  $(M, \tilde{D})$  be tractable, and  $P \subseteq M$  be a finite set of points. Let  $\varepsilon \in (0, 1/4)$ . Let  $s : P \rightarrow [0, \infty)$  be a function on  $P$  such that*

$$s(p) \geq \max_{C \in M \times [0, \infty), |C|=k} \frac{\tilde{D}_W(p, C)}{\sum_{q \in P} \tilde{D}_W(q, C)}.$$

*Let  $\mathcal{T} = \sum_{p \in P} s(p)$ , and  $b$  be a sufficiently large constant. Pick a (multi)-set  $A$  of  $b\mathcal{T}^2(\dim(M, \tilde{D}, k) + \log(1/\delta))/\varepsilon^2$  points from  $P$  by repeatedly, i.i.d., selecting  $p \in P$  with probability  $s(p)/\mathcal{T}$ . For  $p \in A$  let  $\nu(p) = \mathcal{T}/(|A| \cdot s(p))$ . Then, with probability at least  $1 - \delta$ :*

*For all  $C \in M \times [0, \infty)$  and  $|C| = k$ :*

$$\left| \sum_{p \in P} \tilde{D}_W(p, C) - \sum_{p \in A} \nu(p) \tilde{D}_W(p, C) \right| \leq \varepsilon \sum_{p \in P} \tilde{D}_W(p, C).$$

*Proof.* Let  $X = (M \times [0, \infty))^k$ . For every  $p \in P$  and  $C \in X$ , let  $f_p(C) = \tilde{D}_W(p, C)$ ,  $s_{f_p} = f'_p = f_p$ ,  $m(f_p) = n \cdot s(p)/\mathcal{T}$ , and  $g_{f_p}(C) = f_p(C)/m(f_p)$ . Let  $G_{f_p}$  consists of  $m(f_p)$  copies of  $g_f$  and let  $G = \bigcup_{p \in P} G_{f_p}$ . Hence,  $S = \{g_{f_p} \mid p \in A\}$  is a uniform random sampling from  $G$ . By [21, Theorem 6.9], for a sufficiently large  $b$ ,  $S$

is an  $(\varepsilon/(2\mathcal{T}))$ -approximation of  $G$ , with probability at least  $1 - \delta$ . Assume that this event indeed occurs. Let  $U = \{g \cdot |G|/|S| \mid g \in S\}$ . By [21, Theorem 13.1], we obtain that for all  $C \in X$ ,

$$(4.8) \quad \left| \sum_{p \in P} f_p(C) - \sum_{f \in U} f(C) \right| \leq \frac{\varepsilon}{\mathcal{T}} \max_{p \in P} \frac{f_p(C)}{m(f_p)} \sum_{p \in P} m(f_p).$$

We have

$$\sum_{p \in P} m(f_p) = \sum_{p \in P} \frac{ns(p)}{\mathcal{T}} = n,$$

and, for every  $C \in X$

$$\max_{p \in P} \frac{f_p(C)}{m(f_p)} = \frac{\mathcal{T}}{n} \max_{p \in P} \frac{\tilde{D}_W(p, C)}{s(p)} \leq \frac{\mathcal{T}}{n} \sum_{p \in P} \tilde{D}_W(p, C).$$

For every  $f = g_{f_p} \cdot |G|/|S| \in U$  we have

$$\begin{aligned} f(C) &= \frac{g_{f_p}(C) \cdot |G|}{|S|} \\ &= \frac{f_p(C) \cdot n}{m(f_p)|A|} \\ &= \frac{\tilde{D}_W(p, C) \cdot n}{m(f_p)|A|} \\ &= \frac{\tilde{D}_W(p, C) \cdot \mathcal{T}}{s(p) \cdot |A|} \\ &= \nu(p) \tilde{D}_W(p, C). \end{aligned}$$

Substituting the last three inequalities in (4.8) yields

$\forall C \in X :$

$$\begin{aligned} &\left| \sum_{p \in P} f_p(C) - \sum_{p \in A} \nu(p) \tilde{D}_W(p, C) \right| \\ &\leq \frac{\varepsilon}{\mathcal{T}} \cdot \frac{\mathcal{T}}{n} \sum_{q \in P} \tilde{D}_W(q, C) \cdot n \\ &= \varepsilon \sum_{p \in P} \tilde{D}_W(p, C). \end{aligned}$$

**THEOREM 4.1.** *Let  $(M, \tilde{D})$  be tractable. Let  $P \subseteq M$  be a set of  $n$  points, and  $(\varepsilon, \delta) \in (0, 1/10)$ . Let  $A$  and  $\nu$  be the output of the procedure **Core-Set** $(P, k, \varepsilon, \delta)$ ; see **Algorithm 2**. Then the following hold:*

(i)

$$|A| = \frac{k^{O(k)} (\log n)^2}{\varepsilon^2} \cdot (\dim(M, \tilde{D}, k) + \log(1/\delta)).$$

(ii) With probability at least  $1 - \delta$ ,

$$\begin{aligned} \forall C \subseteq M \times [0, \infty), |C| = k : \\ \left| \sum_{p \in P} \tilde{D}_W(p, C) - \sum_{p \in A} \nu(p) \cdot \tilde{D}_W(p, C) \right| \\ \leq \varepsilon \sum_{p \in P} \tilde{D}_W(p, C). \end{aligned}$$

---

**Algorithm 2: Core-Set**( $P, k, \varepsilon, \delta$ )

---

**Input:** A set  $P \subseteq M$ , an integer  $k \geq 1$ , and  $\tau, \delta \in (0, 1/10)$  where  $(M, \tilde{D})$  is tractable.

**Output:** A set  $A$  and a probability measure  $\nu$  on  $A$  that satisfy Theorem 4.1.

```

1  $b \leftarrow$  a constant that can be determined from the
  proof of Theorem 4.1
2  $Q_0 \leftarrow P$ 
3 while  $|Q_0| > b$  do
4    $(q_k, Q_k) \leftarrow$ 
     Recursive-Robust-Median( $Q_0, k$ )
     /* See Algorithm 1. */
5   for each  $p \in Q_k$  do  $s(p) \leftarrow \frac{bk}{|Q_k|}$ 
6    $Q_0 \leftarrow Q_0 \setminus Q_k$ 
7 for each  $p \in Q_0$  do  $s(p) \leftarrow 1$ 
8  $\mathcal{T} \leftarrow \sum_{p \in P} s(p)$ 
9 Pick a (multi)-set  $A$  of
   $b\mathcal{T}^2(\dim(M, \tilde{D}, k) + \log(1/\delta))/\varepsilon^2$  points from  $P$ 
  by repeatedly, i.i.d., selecting  $p \in P$  with
  probability  $s(p)/\mathcal{T}$ 
10 for each  $p \in A$  do  $\nu(p) \leftarrow \frac{\mathcal{T}}{|A| \cdot s(p)}$ 
11 return  $(A, \nu)$ 

```

---

The structure of the algorithm is this: Algorithm 1 is run repeatedly to identify a “dense” cluster in the data (Line 4). Due to Lemma 3.1 the sensitivity of each point in this cluster is bounded by some constant divided by the current number of points. The cluster is then removed, and we repeat.

*Proof.* (i) For every  $i \in [k]$ , let  $Q_i^{(j)}$  denote the value of  $Q_i$  at Line 3 of Algorithm 1 during the  $j$ th “while” iteration. Let  $J$  denote the total number of “while” iterations. By Line 4 of Algorithm 1, we have that  $|Q_i^{(j)}| \geq |Q_{i-1}^{(j)}|/(bk)$ . Hence,

$$|Q_k^{(j)}| \geq \frac{|Q_0^{(j)}|}{(bk)^k} = \frac{|Q_0^{(j)}|}{k^{O(k)}}.$$

By the last equation and Line 6 of Alg. 2, for every  $j \in [J - 1]$  we have

$$\begin{aligned} (4.9) \quad |Q_0^{(j+1)}| &= |Q_0^{(j)}| - |Q_k^{(j)}| \leq |Q_0^{(j)}| - \frac{|Q_0^{(j)}|}{k^{O(k)}} \\ &= |Q_0^{(1)}| \left(1 - \frac{1}{k^{O(k)}}\right)^j = n \left(1 - \frac{1}{k^{O(k)}}\right)^j. \end{aligned}$$

Since  $|Q_0^{(J)}| \geq 1$ , substituting  $j = J - 1$  in the previous inequality we conclude that

$$(4.10) \quad J \leq k^{O(k)} \log n.$$

By Line 3, the size of  $Q_0$  during the execution of Line 7 is  $O(1)$ . By the definition of  $s(p)$  in Lines 5 and 7 we have

$$\begin{aligned} \mathcal{T} &= \sum_{p \in P} s(p) = \sum_{j \in [J]} \sum_{p \in Q_k} s(p) + O(1) \\ &= \sum_{j \in [J]} \sum_{p \in Q_k} \frac{bk}{|Q_k|} + O(1) = J \cdot bk + O(1). \end{aligned}$$

Together with (4.10) we obtain  $\mathcal{T} \leq k^{O(k)} \log n$ . By this and Line 9 we thus have

$$\begin{aligned} |A| &= \frac{b\mathcal{T}^2 (\dim(M, \tilde{D}, k) + \log(1/\delta))}{\varepsilon^2} \\ &= \frac{k^{O(k)} (\log n)^2}{\varepsilon^2} \cdot (\dim(M, \tilde{D}, k) + \log(1/\delta)). \end{aligned}$$

(ii) The pair  $(q_k^{(j)}, Q_k^{(j)})$  satisfies Lemma 3.1 for every  $j \in [J]$ . Hence, for the value  $s(p)$  that is defined in Line 5 of Algorithm 2, and an appropriate  $b$ ,

$$s(p) = \frac{bk}{|Q_k^{(j)}|} \geq \frac{\tilde{D}_W(p, C)}{\sum_{q \in Q_0^{(j)}} \tilde{D}_W(q, C)} \geq \frac{\tilde{D}_W(p, C)}{\sum_{q \in P} \tilde{D}_W(q, C)}.$$

By Corollary 4.1, with probability at least  $1 - \delta/b$  we have

$$\forall C \in M \times [0, \infty), |C| = k :$$

$$\left| \sum_{p \in P} \tilde{D}_W(p, c) - \sum_{p \in A} \nu(p) \tilde{D}_W(p, c) \right| \leq \varepsilon \sum_{p \in P} \tilde{D}_W(p, c).$$

## 5 Efficient implementation

Algorithm 2 calls the procedure `Recursive-Robust-Median` in Line 4 which is in turn required to compute a robust median in its third line. In this section we show efficient computation of such a robust median that succeeds with high probability. In Theorem 5.1 we then bound the running time of Algorithm 2.



LEMMA 5.1. ([21]) Let  $(M, \tilde{D})$  be tractable. Let  $Q \subseteq M$  be a finite set of points. Let  $\gamma \in (0, 1)$ , and  $\tau, \delta \in (0, 1/10)$ . Pick uniformly, i.i.d., a (multi)-set  $S$  of

$$s = \frac{b}{\tau^4 \gamma^2} \cdot \log\left(\frac{1}{\delta}\right)$$

points from  $Q$ , where  $b$  is a sufficiently large universal constant. With probability at least  $1 - \delta$ , any  $((1 - \tau)\gamma, \tau, 2)$ -median for  $S$  is a  $(\gamma, 4\tau, 2)$ -median for  $Q$ .

*Proof.* For every  $p \in P$  and  $c \in M$  let  $f(c) = \tilde{D}(p, c)$ . Let  $\mathcal{D}(S) = S$  for every  $S \subset P$ . Using the (weak) triangle inequality, we have that one of the points of  $S$  is a constant factor approximation for the median of  $S$ . The theorem now follows from [21, Lemma 9.6].

LEMMA 5.2. Let  $(M, \tilde{D})$  be tractable. Let  $Q \subseteq M$  be a finite set of points,  $k \geq 1$  an integer, and  $\delta \in (0, 1)$ . Let  $q \in M$  be the output of a call to `Median`( $Q, k, \delta$ ); See Algorithm 3. Then, with probability at least  $1 - \delta$ , the point  $q$  is a  $(1/k, 1/4, 2)$ -median for  $Q$ . The running time of the algorithm is  $tb^2k^4 \log^2(1/\delta)$  where  $O(t)$  is the time it takes to compute  $\tilde{D}(p, q)$  for  $p, q \in M$ s.

---

**Algorithm 3: Median**( $Q, k, \delta$ )

---

**Input:** A finite set  $Q \subseteq M$ , an integer  $k \geq 1$ , and  $\delta \in (0, 1/10)$ .

**Output:** A point  $q \in M$  that satisfies Theorem 5.2.

- 1  $b \leftarrow$  a universal constant that can be determined from the proof of Theorem 5.2
  - 2 Pick a uniform i.i.d. sample  $S$  of  $bk^2 \log(1/\delta)$  points from  $Q$
  - 3  $q \leftarrow$  a point that minimizes  $\sum_{p \in \text{closest}(S, \{q\}, 15/(16k))} \tilde{D}(p, q)$  over  $q \in S$
  - 4 **return**  $q$
- 

*Proof.* We consider the variables  $b$  and  $S$  as defined in Algorithm 3.

Running time: For a set  $Q \subseteq M$ , the running time of `Median`( $Q, k, \delta$ ) is dominated by Line 3 which can be implemented in  $t|S|^2 = tb^2k^4 \log^2(1/\delta)$  time by computing the distance between every pair of points in  $S$  and using order statistics.

Correctness: Put  $\tau = 1/16$  and  $\gamma = 1/k$ . Let  $q^* \in M$  be a  $((1 - \tau)\gamma, 0, 1)$ -median of  $S$ . Let  $q$  be the closest point to  $q^*$  in  $S$ . By (2.3), for every  $p \in M$  we have

$$\tilde{D}(p, q) \leq \rho(\tilde{D}(p, q^*) + \tilde{D}(q^*, q)) \leq 2\rho \cdot \tilde{D}(p, q^*).$$

Summing this over every  $p \in \text{closest}(S, \{q^*\}, (1 - \tau)\gamma)$  yields

$$\sum_{p \in \text{closest}(S, \{q^*\}, (1 - \tau)\gamma)} \tilde{D}(p, q) \leq 2\rho \tilde{D}^*(S, (1 - \tau)\gamma).$$

Hence,  $q$  is a  $((1 - \tau)\gamma, 0, 2)$ -median of  $S$ , which is also a  $((1 - \tau)\gamma, \tau, 2)$ -median of  $S$ . The theorem now follows from Lemma 5.1.

THEOREM 5.1. Let  $(M, \tilde{D})$  be tractable. Let  $P \subseteq M$  be a set of  $n$  points, and  $(\varepsilon, \delta) \in (0, 1/10)$ . Then a set  $A$  and a function  $\nu$  can be computed in time

$$ntk^{O(k)} + tk^{O(k)} \log(n) \log^2(\log(n)/\delta) + \frac{k^{O(k)}(\log n)^3}{\varepsilon^2} \cdot (\dim(M, \tilde{D}, k) + \log(1/\delta)).$$

such that the following hold:

(i)

$$|A| = \frac{k^{O(k)}(\log n)^2}{\varepsilon^2} \cdot (\dim(M, \tilde{D}, k) + \log(1/\delta)).$$

(ii) With probability at least  $1 - \delta$ ,

$$\forall C \subseteq M \times [0, \infty), |C| = k :$$

$$\left| \sum_{p \in P} \tilde{D}_W(p, C) - \sum_{p \in A} \nu(p) \cdot \tilde{D}_W(p, C) \right| \leq \varepsilon \sum_{p \in P} \tilde{D}_W(p, C).$$

*Proof.* Let  $A$  and  $\nu$  be the output of Algorithm 2. Properties (i) and (ii) then follow from Theorem 4.1. It is left to bound the construction time. The running time of Algorithm 2 is dominated by the call to `Recursive-Robust-Median` in Line 4 with  $Q_0^{(j)}$  during the  $j$ th “while” iteration, which we bound as follows.

Line 3 of `Recursive-Robust-Median` compute a median  $q_i$  of  $Q_{i-1}^{(j)}$  for every  $i \in [k]$ . By replacing  $\delta$  with  $\delta' = \delta/(bJk)$  in Lemma 5.2, we have that, with probability at least  $1 - \delta'$ , such a  $q_i$  can be computed in time  $tb^2k^4 \log^2(1/\delta')$ . Hence, with probability at least  $1 - k\delta' = 1 - \delta/(bJ)$ , the desired  $q_i$  is computed for every  $i \in [k]$ . Line 4 of `Recursive-Robust-Median` can be computed in  $t \cdot |Q_{i-1}^{(j)}| \leq t \cdot |Q_0^{(j)}|$  using order statistics.

We conclude that with probability at least  $1 - \delta/b$ , Line 4 of Algorithm 2 computes the desired pair  $(q_k, Q_k)$  during all the  $J$  “while” iterations. Using (4.9), the

overall execution time of Line 4 is

$$\begin{aligned} & \sum_{j \in [J]} t \cdot |Q_0^{(j)}| + \sum_{j \in [J]} tk^2 k^5 \log^2(1/\delta^j) \\ & \leq tk \sum_{j \in [J]} |Q_0^{(j)}| + Jtk^5 \log^2(J/\delta) \\ & \leq ntk \sum_{j \in [J]} \left(1 - \frac{1}{k^{O(k)}}\right)^j \\ & \quad + tk^{O(k)} \log(n) \log^2(\log(n)/\delta) \\ & \leq ntk^{O(k)} + tk^{O(k)} \log(n) \log^2(\log(n)/\delta). \end{aligned}$$

Line 9 of Algorithm 2 can be implemented using a binary tree in time

$$\log(n) \cdot |A| \leq \frac{k^{O(k)}(\log n)^3}{\varepsilon^2} \cdot (\dim(M, \tilde{D}, k) + \log(1/\delta)).$$

The last two inequalities prove the theorem.

The following main theorem includes both  $F_{d,k}^{\text{wght}}$  and  $F_{d,k}^{\text{out}}$  as special cases by taking, respectively,  $h = \infty$  or all weights equal.

**THEOREM 5.2.** *Let  $P$  be a set of points in a metric space  $(M, \text{dist})$ ,  $r \in (0, \infty)$  be a constant and  $h \in (0, \infty)$ . Let  $k \geq 1$  be an integer, and  $\varepsilon \in (0, 1/10)$ . A set  $A$  of size*

$$|A| = \frac{k^{O(k)}(\log n)^2}{\varepsilon^2} \cdot (\dim(M, \tilde{D}, k) + \log(1/\delta))$$

and a weight function  $\nu : A \rightarrow \mathbb{R}_+$  can be computed in time

$$\begin{aligned} & ntk^{O(k)} + tk^{O(k)} \log(n) \log^2(\log(n)/\delta) \\ & + \frac{k^{O(k)}(\log n)^3}{\varepsilon^2} \cdot (\dim(M, \tilde{D}, k) + \log(1/\delta)), \end{aligned}$$

such that, with probability at least  $1 - \delta$ ,

$$\begin{aligned} & \forall C \subseteq M \times [0, \infty), |C| = k : \\ & \left| \sum_{p \in P} f_C(p) - \sum_{p \in A} \nu(p) f_C(p) \right| \leq \varepsilon \sum_{p \in P} f_C(p), \end{aligned}$$

where  $f_C(p) := \min_{(c,w) \in C} (w \min\{h, \text{dist}^r(p, c)\})$ , and  $t$  is the time it takes to compute  $\tilde{D}(p, q)$  for  $p, q \in M$ .

*Proof.* Define  $D : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  as  $D(x) = \min\{h, x^r\}$ , and  $\tilde{D}(p, q) = D(\text{dist}(p, q))$ . Using Theorem 4.1 and Lemma 2.1, it suffices to note that  $D$  is monotone non-decreasing and to show that

$$D(xe^\delta) \leq e^{r\delta} D(x).$$

Indeed, for every  $x, \delta \geq 0$  we have

$$\begin{aligned} (5.11) \quad D(xe^\delta) &= \min\{h, x^r e^{r\delta}\} \\ &\leq \min\{e^{r\delta} h, x^r e^{r\delta}\} \\ &= e^{r\delta} \min\{h, x^r\} \\ &= e^{r\delta} D(x). \end{aligned}$$

All the applications that are mentioned in the introduction are straightforward results of the last theorem. For example, the following loss (or distortion) functions, known as *M-estimators*, are popular with statisticians performing robust (i.e., outlier-resistant) estimation [32, 50]. By defining  $x = \text{dist}(p, q)$  and  $h$  as the distance threshold for outliers, we obtain:

$$\begin{aligned} D^{h, \text{Huber}}(\text{dist}(p, q)) &= D^{h, \text{Huber}}(x) \\ &= \min\{x^2/2, h(x - h/2)\} \\ D^{h, \text{Tukey}}(\text{dist}(p, q)) &= D^{h, \text{Tukey}}(x) \\ &= \min\{h^2/6, \frac{3h^4 x^2 - 3h^2 x^4 + x^6}{6h^4}\} \end{aligned}$$

It is straightforward to show that the families  $(M, \tilde{D}^{1, \text{Huber}})$  and  $(M, \tilde{D}^{1, \text{Tukey}})$  are tractable, with a proof very similar to that in (5.11). The shatter function is again that of balls in the metric space. Consequently, Theorem 5.2 shows that we can perform data reduction for  $k$ -clustering (with weights) for these loss functions.

## References

- [1] P. Agarwal and J. Phillips. An efficient algorithm for 2D Euclidean 2-center with outliers. *Algorithms-ESA 2008*, pages 64–75, 2008.
- [2] P. K. Agarwal, C. M. Procopiuc, and K. R. Varadarajan. Approximation algorithms for k-line center. In *Proc. 10th Ann. European Symp. on Algorithms (ESA)*, volume 2461 of *Lecture Notes in Computer Science*, pages 54–63. Springer, 2002.
- [3] N. Alon, S. Ben-David, N. Cesa-Bianchi, and D. Haultsler. Scale-sensitive dimensions, uniform convergence, and learnability. *J. ACM*, 44(4):615–631, July 1997.
- [4] N. Alon and B. Sudakov. On two segmentation problems. *J. of Algorithms*, 33:173–184, 1999.
- [5] V. Angadi and S. Patel. Hough transformation. Technical report, Rochester Institute of Technology, 2010.
- [6] V. Arya, N. Garg, R. Khandekar, A. Meyerson, K. Munagala, and V. Pandit. Local search heuristics for k-median and facility location problems. *SIAM Journal on Computing*, 33(3):544–562, 2004.
- [7] P. L. Bartlett and P. M. Long. Prediction, learning, uniform convergence, and scale-sensitive dimensions. *J. Comput. System Sci.*, 56:174–190, 1998.

- [8] P. L. Bartlett, P. M. Long, and R. C. Williamson. Fat-shattering and the learnability of real-valued functions. *J. Comput. System Sci.*, 52:434–452, 1996.
- [9] S. Ben-David, N. Cesa-Bianchi, D. Haussler, and P. M. Long. Characterizations of learnability for classes of  $\{0, \dots, n\}$ -valued functions. *J. of Comput. Syst. Sci.*, 50(1):74–86, 1995.
- [10] A. Benczúr and D. Karger. Approximating s-t minimum cuts in  $\tilde{O}(n^2)$  time. In *Proc. 28th STOC*, pages 47–55, 1996.
- [11] D. Bienstock, M. X. Goemans, D. Simchi-Levi, and D. P. Williamson. A note on the prize collecting traveling salesman problem. *Math. Program*, 59:413–420, 1993.
- [12] A. Blumer, A. Ehrenfeucht, D. Haussler, and M. K. Warmuth. Learnability and the Vapnik-Chervonenkis dimension. *J. ACM*, 36(4):929–965, October 1989.
- [13] M. Charikar, S. Khuller, D. M. Mount, and G. Narasimhan. Algorithms for facility location problems with outliers. In *Proceedings of the twelfth annual ACM-SIAM symposium on Discrete algorithms, SODA '01*, pages 642–651, Philadelphia, PA, USA, 2001. Society for Industrial and Applied Mathematics.
- [14] K. Chen. A constant factor approximation algorithm for k -median clustering with outliers. In Shang-Hua Teng, editor, *Proceedings of the Nineteenth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2008, San Francisco, California, USA, January 20-22, 2008*, pages 826–835. SIAM, 2008.
- [15] W. Fernandez de la Vega and C. Kenyon. A randomized approximation scheme for metric max-cut. In *39th Ann. Symp. on Foundations of Computer Science*, pages 468–471. IEEE, 1998.
- [16] A. Deshpande and K. R. Varadarajan. Sampling-based dimension reduction for subspace approximation. In *Proc. 39th Ann. ACM Symp. on Theory of Computing (STOC)*, pages 641–650, 2007.
- [17] P. Drineas, A. Frieze, R. Kannan, S. Vempala, and V. Vinay. Clustering large graphs via the singular value decomposition. In *Machine Learning*, volume 56, pages 9–33, 2004.
- [18] M. Effros and L. J. Schulman. Deterministic clustering with data nets. <http://eccc.hpi-web.de/eccc-reports/2004/TR04-050/index.html>, 2004. ECCC TR04-050.
- [19] T. Feder and D. Greene. Optimal algorithms for approximate clustering. In *Proceedings of the twentieth annual ACM symposium on Theory of computing, STOC '88*, pages 434–444, New York, NY, USA, 1988. ACM.
- [20] D. Feldman, A. Fiat, and M. Sharir. Coresets for weighted facilities and their applications. In *Proc. 47th IEEE Ann. Symp. on Foundations of Computer Science (FOCS)*, pages 315–324, 2006.
- [21] D. Feldman and M. Langberg. A unified framework for approximating and clustering data. In L. Fortnow and S. P. Vadhan, editors, *STOC*, pages 569–578. ACM, 2011. Full version at arXiv:1106.1379; theorem citations follow the full version.
- [22] D. Feldman, M. Monemizadeh, and C. Sohler. A PTAS for k-means clustering based on weak coresets. In *Proc. 23rd ACM Symp. on Computational Geometry (SoCG)*, pages 11–18, 2007.
- [23] B. Foster, S. Mahadevan, and R. Wang. A GPU-based Approximate SVD Algorithm. Technical report, U. Massachusetts, 2010.
- [24] M. X. Goemans and D. P. Williamson. A general approximation technique for constrained forest problems. *SIAM Journal on Computing*, 24(2):296–317, 1995.
- [25] F. Hampel, C. Hennig, and E. Ronchetti. A smoothing principle for the Huber and other location M-estimators. *Computational Statistics & Data Analysis*, 55(1):324–337, 2011.
- [26] S. Har-Peled. No coreset, no cry. In *Proc. 24th Int. Conf. Foundations of Software Technology and Theoretical Computer Science (FSTTCS)*, volume 3328 of *Lecture Notes in Computer Science*, pages 324–335. Springer, 2004.
- [27] S. Har-Peled. Coresets for discrete integration and clustering. In *Proc. 26th Int. Conf. Foundations of Software Technology and Theoretical Computer Science (FSTTCS)*, volume 4337 of *Lecture Notes in Computer Science*, pages 33–44. Springer, 2006.
- [28] S. Har-Peled and S. Mazumdar. On coresets for k-means and k-median clustering. In *Proc. 36th Ann. ACM Symp. on Theory of Computing (STOC)*, pages 291–300, 2004.
- [29] J. Hardin, A. Mitani, L. Hicks, and B. VanKoten. A robust measure of correlation between two genes on a microarray. *BMC bioinformatics*, 8(1):220, 2007.
- [30] D. Haussler and E. Welzl. Epsilon-nets and simplex range queries. *Discrete Comput. Geom.*, 2:127–151, 1987.
- [31] D. C. Hoaglin, F. Mosteller, and J. W. Tukey. *Understanding robust and exploratory data analysis*, volume 3. Wiley New York, 1983.
- [32] P. J. Huber. *Robust Statistics*. Wiley, 1981.
- [33] M. J. Kearns and R. E. Schapire. Efficient distribution-free learning of probabilistic concepts. *J. Comput. System Sci.*, 48:464–497, 1994.
- [34] J. Kleinberg, C. Papadimitriou, and P. Raghavan. Approximation algorithms for segmentation problems. In *Proc. 30th STOC*, 1998.
- [35] M. Langberg and L. J. Schulman. Universal epsilon-approximators for integrals. In M. Charikar, editor, *SODA*, pages 598–607. SIAM, 2010.
- [36] J. H. Lin and J. S. Vitter.  $\epsilon$ -approximations with minimum packing constraint violation (extended abstract). In *Proceedings of the twenty-fourth annual ACM symposium on Theory of computing*, pages 771–782. ACM, 1992.
- [37] M. Mahajan, P. Nimbhorkar, and K. Varadarajan. The planar k-means problem is NP-hard. In S. Das and R. Uehara, editors, *Proc. WALCOM, LNCS 5431*, pages 274–285. Springer-Verlag, 2009.
- [38] J. Matoušek. Approximations and optimal geometric

- divide-and-conquer. *J. Comput. Syst. Sci*, 50(2):203–208, 1995.
- [39] J. Matoušek. On approximate geometric k-clustering. *Discrete & Computational Geometry*, 24(1):61–84, 2000.
- [40] B. K. Natarajan. On learning sets and functions. *Machine Learning*, 4:67–97, 1989.
- [41] S. Negahban, P. Ravikumar, M. J. Wainwright, and B. Yu. A unified framework for high-dimensional analysis of  $m$ -estimators with decomposable regularizers. *Arxiv preprint arXiv:1010.2731*, 2010.
- [42] M. Owen. Tukey’s Biweight Correlation and the Breakdown. Master’s thesis, Pomona College, 2010.
- [43] D. Pollard. *Convergence of stochastic processes*. Springer, 1984.
- [44] P. J. Rousseeuw, A. M. Leroy, and J. Wiley. *Robust regression and outlier detection*, volume 3. Wiley Online Library, 1987.
- [45] N. Sauer. On the density of families of sets. *J. Comb. Theory, Ser. A*, 13(1):145–147, 1972.
- [46] L. J. Schulman. Clustering for edge-cost minimization. In *Proc. 32nd STOC*, pages 547–555, 2000.
- [47] S. Shelah. A combinatorial problem, stability and order for models and theories in infinitary languages. *Pacific J. Math.*, 41:247–261, 1972.
- [48] V. N. Vapnik. Inductive principles of the search for empirical dependences. In *Proc. 2nd Annual Workshop on Computational Learning Theory*, 1989.
- [49] V. N. Vapnik and A. Y. Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability and its Applications*, 16(2):264–280, 1971.
- [50] Z. Zhang. M-estimators. <http://research.microsoft.com/en-us/um/people/zhang/INRIA/Publis/Tutorial-Estim/node20.html>. [Online; accessed July 2011].