
Online Facility Location with Multiple Advice

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Clustering is a central topic in unsupervised learning and its online formulation has
2 received a lot of attention in recent years. In this paper, we study the classic facility
3 location problem in the presence of multiple machine-learned advice. We design
4 an algorithm with provable performance guarantees such that, if the advice is good,
5 it outperforms the best-known online algorithms for the problem, and if it is bad it
6 still matches their performance. We complement our theoretical analysis with an
7 in-depth study of the performance of our algorithm, showing its effectiveness on
8 synthetic and real-world data sets.

9 1 Introduction

10 Clustering is a central topic in unsupervised learning [Jain and Dubes, 1988, Aggarwal and Reddy,
11 2014, Gan et al., 2020] In the past few years, its online version has gained a lot of attention [Zhang
12 et al., 2004, Liberty et al., 2016, Lattanzi and Vassilvitskii, 2017, Cohen-Addad et al., 2021]. In this
13 paper we are interested in the fundamental question of whether machine-learned advice, especially
14 when coming from multiple sources, can boost the performance of clustering algorithms in the online
15 setting. In particular, we are interested in designing online algorithms with access to multiple advice
16 that exhibit strong theoretical guarantees together with a good experimental performance.

17 To explore the interplay between online algorithms and multiple advice we focus on facility location,
18 a classic clustering problem with a long and rich history in computer science and operations research
19 [Korte and Vygen, 2018], which can also be seen as the Lagrangian relaxation of k -median clustering.
20 In the online setting of k -median clustering, it is easy to see that if the requirement of having k
21 clusters is a hard constraint then the cost of solutions is bound to be arbitrarily far from the optimum,
22 which makes the straightforward online formulation rather uninteresting algorithmically. There are
23 two standard ways in which the constraint can be relaxed: bi-criteria solutions and the Lagrangian
24 relaxation. Facility location captures the latter.

25 In the online version of facility location, points arrive in sequence and, as soon as a point arrives, an
26 irrevocable decision must be made: we either assign it to an existing opened facility or we open a new
27 one to serve it. Each facility, therefore, induces a cluster, consisting of the points assigned to it. The
28 cost of this sequence of decisions is given by the sum of the service costs (*i.e.* the distance between a
29 point and its facility) and the total facility cost (opening each facility has a cost, which in this paper
30 we assume to be uniform). The goal is to find as cheap a clustering as possible. This online version
31 was first introduced by Meyerson [2001], who gave an elegant $O(\log(n))$ -approximation algorithm,
32 and studied extensively thereafter [Fotakis, 2003, 2005, 2011, Cygan et al., 2018, Cohen-Addad et al.,
33 2019, Ahmed et al., 2020, Guo et al., 2020]. One of the main reasons for this heightened attention
34 is that even in dynamic settings the resulting clustering is very stable. This is particularly useful in
35 real-world systems where clusters are served directly to users or when they are used in a downstream
36 machine learning model, a situation where modifying the clustering would incur too high a cost.

37 In real-world applications, one typically has side information that could help solve the online
 38 clustering problem. This natural scenario is a recurrent theme in machine learning that focuses on
 39 how to take advantage of expert advice. In the context of online algorithms, Lykouris and Vassilvitskii
 40 [2018] and Mahdian et al. [2012] introduced a formal framework for online problems within which
 41 mathematically rigorous performance guarantees can be given in terms of the quality of the advice. In
 42 particular, they showed how to incorporate the advice in a *robust* way: the advice is exploited when it
 43 is good and disregarded when it is bad. All this happens automatically: the robust algorithm does
 44 not need to know in advance which is which. Remarkably, when the advice is bad the performance
 45 of the best online algorithm is matched. The approach combines mathematical rigor with practical
 46 effectiveness, adding the desirable dimension of robustness. The solution we develop for facility
 47 location is in the same spirit, with a focus on multiple advice which adds an important new dimension.
 48 Thus, the goal of this paper is to show how to take advantage of multiple advice within this rigorous
 49 algorithmic framework in the case of online facility location.

In our set-up the advice comes in the form of a family of sets, each suggesting a list of facilities to
 open. As we show with our experiments, multiple advice of this form can be easily and efficiently
 produced on the basis of past data. The main technical result of this paper is the following. Suppose
 we have a family of sets $\mathcal{S}_1, \dots, \mathcal{S}_k$ which intuitively represent sets of suggested facilities to open,
 and let \mathcal{S} be their union. We develop an online randomized algorithm, called TAKEHEED, that
 computes a solution whose expected cost is

$$\mathcal{O}(\log(|\mathcal{S}|) \cdot \text{OPT}(\mathcal{S})),$$

50 where $\text{OPT}(\mathcal{S})$ is the best solution that can be obtained using only facilities in \mathcal{S} . Notice that the
 51 $\text{OPT}(\mathcal{S})$ factor is no worse than any of the suggested sets in isolation since $\text{OPT}(\mathcal{S}) \leq \min_i \text{OPT}(\mathcal{S}_i)$.
 52 In fact, it can be much better. For instance, every set could contain just a few good points and a lot of
 53 noise. Or, one of the sets could be good but its identity unknown. And yet, TAKEHEED separates the
 54 wheat from the chaff and delivers a near-optimal solution. This result (which is the main technical
 55 step toward the solution) can be seen as a solution to the version of facility location in which the
 56 facilities are constrained to be in a specific set, as opposed to the entire input metric space. The above
 57 result is essentially the best possible. In Section 4, we prove that the $\log |\mathcal{S}|$ factor is near-optimal.
 58 We also show that the best-known algorithms (see [Meyerson, 2001, Fotakis, 2003]) cannot match
 59 the same bound, for they can only attain an $\Omega(|\mathcal{S}|)$ approximation factor in general. This is a strong
 60 indication that techniques different from those employed by those state-of-the-art algorithms are
 61 required for this problem, providing a compelling motivation for our approach which is based on the
 62 theory of Hierarchically Separated Trees (HST's) introduced by Bartal [1996] in a seminal paper.

On the practical side, our experimental analysis with real and synthetic data shows that it is easy and
 inexpensive to compute multiple advice of the form described, allowing TAKEHEED to outperform
 state-of-the-art algorithms such as those in Meyerson [2001], Fotakis [2003], Anagnostopoulos et al.
 [2004]. Furthermore, we also show that we can make our algorithm robust to erroneous advice. By
 combining TAKEHEED with known results of Mahdian et al. [2012] we can get a robust version of
 TAKEHEED called WARY, whose performance guarantee is

$$\mathcal{O}\left(\min\left\{\log(|\mathcal{S}|) \cdot \text{OPT}(\mathcal{S}), \frac{\log(n)}{\log \log(n)} \cdot \text{OPT}\right\}\right).$$

63 Note that the second term matches the best possible bound, in view of the lower bound of Fotakis
 64 [2003]. Our experiments confirm the good behavior in practice of this robust version. Indeed, it
 65 outperforms other state-of-the-art robust solutions.

One of the appealing features of TAKEHEED is that it leverages multiple advice in a simple way.
 The algorithm however could suffer from an “adversarial” attack of sorts. If very large bad sets
 of suggested facilities are given as part of the advice, the theoretical guarantee, which depends on
 the cardinality of the union of the advice sets, becomes uninteresting. To deal with this we design
 ROBUST-TAKEHEED, an algorithm that is resilient to the presence of such large, bad advice. If the
 best advice set is \mathcal{S}_{i^*} , the expected cost of the solution produced by this algorithm is

$$\mathcal{O}(\log \log(n) \cdot (\log(k) + \log(|\mathcal{S}_{i^*}|)) \cdot \text{OPT}(\mathcal{S}_{i^*})).$$

66 This result is mainly of theoretical interest. Indeed, one of the strengths of TAKEHEED is that in
 67 practice small advice sets are easy to obtain and the term $\log(|\mathcal{S}|)$ is going to be quite small, making
 68 the overall approach very actionable.

69 **Related Work.** Facility location and online algorithms have been extensively studied in the algo-
70 rithm and machine learning literature. Here we give a high-level overview of the main results.

71 *Facility Location.* The offline metric (uncapacitated) facility location problem has been extensively
72 studied, and various algorithms and heuristics have been formulated (see Korte and Vygen [2018]).
73 The problem is known to be NP-hard, and it cannot be approximated within a factor of 1.463 unless
74 $\text{NP} \subseteq \text{DTIME}(n^{\log \log(n)})$ [Guha and Khuller, 1999]. The best-known polynomial-time algorithm
75 yields a 1.488 approximation ratio [Li, 2013].

76 *Online Setting.* We cast our results in the standard setting of online competitive analysis [Borodin
77 and El-Yaniv, 2005]. The online version of the metric facility location problem was introduced by
78 Meyerson [2001], together with a randomized algorithm providing an expected $O(\log(n))$ approx-
79 imation of the optimal solution. Subsequently, Fotakis [2003] provided a deterministic algorithm
80 with a $O(\log(n)/\log \log(n))$ competitive ratio and proved that no randomized algorithm can obtain
81 an asymptotically better competitive ratio against an oblivious adversary, even if the points lie on
82 a line. Anagnostopoulos et al. [2004] also provided an online algorithm for facility location with
83 restricted facilities: their algorithm works only in the euclidean plane, and our procedure combining
84 the advice can be seen as a generalization of this result in a generic metric space. Mahdian et al.
85 [2012] showed how to combine two online algorithms for facility location into a single algorithm
86 achieving a constant approximation of the costs of the best algorithm. In fact, the mixing procedure
87 can be also used to combine an algorithm with a suggested set of facilities, or two sets of suggested
88 facilities. Applied inductively, it can also be used to mix multiple advice. In asymptotic terms,
89 the resulting bound is incomparable to ours (see the Supplementary Material) and, in practice, our
90 approach appears to be preferable, as documented in the experimental section. Other formulations
91 of online clustering have been proposed: e.g., Mettu and Plaxton [2003] allow to chose the order
92 by which processing the input, which is provided in advance; Guo et al. [2020] consider the online
93 model with recourse, in which one is allowed to close some facilities and reassign the clients to new
94 ones. Guo et al. [2020] also use HSTs for their clustering problem, but in a different way to ours.

95 *Online algorithms with learned advice.* There have been various previous attempts to improve online
96 algorithms with external information. Common assumptions are related to the distribution of the
97 input data. The most popular assumption is that data arrive according to the random arrival model,
98 i.e. the input is a random permutation of a fixed unknown set of elements. This leads to major
99 improvements in the competitive ratio for many problems [McGregor, 2014]. Even Meyerson’s
100 algorithm is shown to have a constant expected competitive ratio in this setting [Meyerson, 2001]. In
101 many other scenarios, it is common to assume that input data are sampled from a fixed distribution,
102 like for online matching [Mirrokni et al., 2012] or job scheduling [Mitzenmacher, 2020]. Recently,
103 Lykouris and Vassilvitskii [2018] formulated the Online with Machine Learned Advice (OMLA)
104 model, which is a theoretical framework for combining online algorithms with machine learning
105 predictions. They also designed an algorithm for the caching problem in this model, which has
106 been improved by Rohatgi [2020]. Various other problems have been tackled in this model with
107 encouraging outcomes, like ski-rental [Purohit et al., 2018] and scheduling with restricted assignments
108 [Lattanzi et al., 2020] (see Mitzenmacher and Vassilvitskii [2020] for a survey). In a slightly different
109 setting, Medina and Vassilvitskii [2017] designed an algorithm for reserve price optimization in
110 auctions. Learned predictions have been also used to improve classical data structures [Kraska et al.,
111 2018, Mitzenmacher, 2018, Hsu et al., 2019]. A different model is the advice model, in which there
112 is an oracle knowing the exact input and providing a small part of it to the online algorithm (see
113 Boyar et al. [2017] for a survey). Note that this oracle is supposed to be perfect and completely
114 reliable, while in our setting, the predictions can be extremely inaccurate, nevertheless the online
115 algorithm needs to be robust to them. Another different setting is the one of online learning with
116 experts [Choromanska and Monteleoni, 2012]: here the advice, as the proposed clustering decisions,
117 can change at every time step, while in our setting they cannot change over time.

118 2 Preliminaries

119 Let us begin by defining our problem: UNIFORM ONLINE FACILITY LOCATION. We are given a
120 sequence of points \mathcal{P} , called *clients*, possibly with repetitions, from a metric space X with distance d ,
121 and $f > 0$ called the uniform *facility cost*. The points are presented to us one after the other. As soon
122 as a point p arrives, we must select a point $c(p) \in X$, called an *opened facility*, and assign p to it. $c(p)$
123 can be an already opened facility if it exists, or a brand new point. At the end of this procedure, every

124 client in $p \in \mathcal{P}$ has a corresponding facility $c(p)$. Let k be the number of opened facilities at the end.
 125 The cost of the assignment is defined as $k \cdot f + \sum_{p \in \mathcal{P}} d(p, c(p))$. If we consider facility location as a
 126 Lagrangian relaxation of k -median, then f plays the role of Lagrangian multiplier (the larger the f ,
 127 the fewer the clusters). To identify an instance of the problem is sufficient to specify the pair (\mathcal{P}, f)
 128 without explicit mention of the underlying metric space (X, d) , which is fixed and whose existence
 129 is clear from the context. We adopt this convention from now on. We denote as $\text{OPT}(P)$, or simply
 130 OPT when the context is clear, the cheapest possible facility assignment for P . Sometimes the set of
 131 possible facilities will be restricted to be in a finite set S . In this case, the optimum assignment will
 132 be denoted as $\text{OPT}(P, S)$.

133 We make extensive use of the elegant theory of *Hierarchically Separated Trees* (HST's). The theory
 134 was introduced in a seminal paper by Bartal [1996]. In what follows we mainly use the definitions
 135 and results of Fakcharoenphol et al. [2004]. An *HST family* for a finite metric space (X, d) consists of
 136 a family of rooted trees \mathcal{T} and a probability distribution \mathcal{D} over the trees of \mathcal{T} satisfying the following
 137 properties. The leaves of every $T \in \mathcal{T}$ are the points of the metric space X . The edges of every tree
 138 in the family are weighted according to a simple scheme. First, the edges from a node to its children
 139 have the same weight. Second, if we follow any path from the root to a leaf the edge weights decrease
 140 exponentially. In this paper we only consider 2-HSTs, so they always decrease by a factor of at least
 141 2. Each tree T induces a simple metric d_T among the leaves: the distance between two leaves u and
 142 v is the sum of the edge weights along the unique path between them in T . Fakcharoenphol et al.
 143 [2004] show that, given any finite metric space (X, d) it is possible to construct in polynomial-time a
 144 2-HST family $(\mathcal{T}, \mathcal{D})$ with this structure which moreover satisfies the following:

145 (i) for every $T \in \mathcal{T}$, and for all $u, v \in X$, $d(u, v) \leq d_T(u, v)$

146 (ii) for all $u, v \in X$, $\mathbb{E}_{T \sim \mathcal{D}}[d_T(u, v)] \leq 8 \log(|X|) \cdot d(u, v)$.

147 In what follows, an HST family will always be a pair $(\mathcal{T}, \mathcal{D})$ with the above properties and tree
 148 structure. All logarithms in this paper are to the base 2.

149 3 The Online Algorithm and Its Analysis

150 Recall our setting. We have an instance of the online facility location and multiple advice in the form
 151 of a family of sets, each suggesting a set of facilities to open. We would like to take advantage of
 152 these suggestions when they are good and ignore them when they are bad. The problem is that we do
 153 not know in advance which is which. The main technical tool toward the solution is the following.

154 **Theorem 3.1.** *Let (\mathcal{P}, f) be an instance of uniform online facility location, let $\mathcal{S}_1, \dots, \mathcal{S}_k$ be sets
 155 of facilities (i.e. finite sets of points from the underlying metric space), and let $\mathcal{S} := \bigcup_{i=1}^k \mathcal{S}_i$. Then,
 156 there is a randomized online algorithm that assigns facilities of \mathcal{S} to points in \mathcal{P} whose expected cost
 157 is $O(\log(|\mathcal{S}|) \cdot \text{OPT}(\mathcal{P}, \mathcal{S}))$.*

158 Intuitively, the sets $\mathcal{S}_1, \dots, \mathcal{S}_k$ play the role of multiple machine-learned advice. We shall refer to
 159 each \mathcal{S}_i as a *suggested set*. The first observation is that, $\text{OPT}(\mathcal{P}, \mathcal{S}) \leq \min_{i \in [k]} \text{COST}(\mathcal{P}, \mathcal{S}_i)$, where
 160 $\text{COST}(\mathcal{P}, X)$ denotes the cost incurred by using all and only the facilities of the set X . Thus, we can
 161 approximate well the best suggested set. In fact, no suggested set by itself needs to be particularly
 162 good. As long as their union contains a good subset, a good solution can be found! For instance,
 163 every suggested set could contain a few good points and a lot of noise. Yet, the algorithm would
 164 be able to extract a good solution. Notice that, even if the advice comes from many sources, we
 165 combine it in a very simple way: we take the union \mathcal{S} of the suggested sets. We believe this is quite
 166 an appealing feature of our approach. However, the $\log(|\mathcal{S}|)$ factor shows that we could suffer from
 167 the presence of a few large and bad advice sets. Even if this is not the case in practical settings, in the
 168 Supplementary Material we show how to deal with this problem.¹

169 If the union \mathcal{S} contains no good subset, the solution can be bad. However, we can obviate the problem
 170 and attain robustness thanks to known results of Mahdian et al. [2012] and Meyerson [2001].

Corollary 3.1.1. *Let \mathcal{P} and \mathcal{S} be as in the statement of Theorem 3.1, and let $n = |\mathcal{P}|$. Then, there is
 an online randomized algorithm whose expected cost is at most,*

$$O\left(\min\left\{\log(|\mathcal{S}|) \cdot \text{OPT}(\mathcal{P}, \mathcal{S}), \frac{\log(n)}{\log \log(n)} \cdot \text{OPT}\right\}\right).$$

¹We develop another Algorithm, ROBUST-TAKEHEED, which does not merge together all the advice sets and produces a solution with quality not depending on $|\mathcal{S}|$ but on the size of the best advice. See the Supplementary Material for more details.

171 Corollary 3.1.1 is a direct application of the results of Mahdian et al. [2012], and shows that our
 172 algorithm is comparable to the true optimum of the problem. Therefore, we focus on Theorem 3.1
 173 from now on and we present the proof of Corollary 3.1.1 with the full description of the algorithm in
 174 the Supplementary Material.

175 Intuitively, we are going to implement the following reduction. The initial problem is to find good
 176 facilities, taken from the set \mathcal{S} , for an online sequence of points $\mathcal{P} := \{p_1, \dots, p_n\}$. Instead, we
 177 will find facilities for the sequence of “nearby” points $Q := \{q_{p_1}, \dots, q_{p_n}\}$, where q_p is the point
 178 of \mathcal{S} closest to p . Recall that, since \mathcal{S} is a metric space, there is an HST family $(\mathcal{T}, \mathcal{D})$ for it. By
 179 construction, the leaves of every $T \in \mathcal{T}$ are the facilities in \mathcal{S} . To find good facilities for Q , we
 180 will pick a tree $T \in \mathcal{T}$ randomly according to \mathcal{D} once and for all and use it to find good facilities
 181 for the entire sequence Q . When a request for q_p arrives, we return a (hopefully) good leaf of T .
 182 In choosing the leaf, we work with the distance d_T induced by the tree T which is (in expectation)
 183 a good approximation of the distance d in the metric space. And so, a good solution in the metric
 184 induced by T is also going to be (in expectation) a good solution in the original metric space. The
 185 advantage of using T relies on the strong structural properties of HSTs. Note that the sequence Q
 186 may contain repeated points even when the input sequence \mathcal{P} does not. Known online algorithms do
 187 not seem to be able to cope with this situation in a satisfactory way. Indeed, in Section 4 we prove a
 188 lower bound that rules out Meyerson’s well-known algorithm, motivating our HST-based approach.
 189 The core of the solution is algorithm PLUCK which finds a good facility for p ’s proxy q_p . To find
 190 facilities for the whole input sequence, PLUCK is called repeatedly by procedure TAKEHEED. The
 191 input to TAKEHEED are the sequence of clients \mathcal{P} from a metric space (X, d) , the set \mathcal{S} of suggested
 192 facilities, and an HST family $(\mathcal{T}, \mathcal{D})$ for the metric space (\mathcal{S}, d) induced by \mathcal{S} .

Algorithm 1 Algorithm TAKEHEED

- 1: Pick T at random from the HST family $(\mathcal{T}, \mathcal{D})$.
 - 2: When a point $p \in \mathcal{P}$ arrives, let $q_p \in \mathcal{S}$ be its closest suggested facility.
 - 3: Call PLUCK with input q_p, T and $\text{root}(T)$; assign the returned facility ℓ_p to p .
-

193 It remains to define PLUCK. Its input consists of a point $q_p \in Q \subseteq \mathcal{S}$, an HST tree T as above, and a
 194 node $v \in T$, initially set to the root of T (so, effectively, we can say that its input consists of q_p and
 195 T). The output is a leaf $\ell \in T$. Starting from the root, PLUCK goes down the tree T recursively in
 196 search of a good leaf for q_p . Once the leaf is found it is *opened* and returned. The tree T is passed
 197 “by reference” so that: *once a leaf is opened, it remains open during all future invocations of PLUCK*.
 198 Initially, no leaf is open. Every node in the tree T has a potential, initially set to zero. The potentials
 199 too are modified “by reference” and are preserved from one invocation of PLUCK to the next. In
 200 Algorithm 2 we give the pseudo-code for PLUCK. The main goal of PLUCK is to serve the newly
 201 arrived client with a facility within distance f from q_p , with the additional guarantee that not too
 202 many facilities are opened by successive executions of the algorithm. To achieve this goal PLUCK
 203 operates as follows. If there is no open facility within distance f from q_p , PLUCK opens a new facility
 204 within distance $f/3$ from q_p . In so doing, it selects the facility that is able to serve the largest number
 205 of future clients (see subroutine SELECTHEAVIESTCHILD). Otherwise, PLUCK opens a new facility
 206 if and only if the potential of an internal node is higher than the cost of opening a facility. Intuitively,
 207 this guarantees that we do not open too many facilities because: in the former case we reduce the
 208 number of possible points without facilities within distance f substantially; in the latter, we can
 209 charge the opening cost of the facility to the serving cost of the instance.

Algorithm 2 Algorithm PLUCK(T, v, q)

Input: HST tree T , a node $v \in T$, a point $q \in \mathcal{S}$ Output: a leaf of T 1: if v is a leaf of T then 2: open v (if already opened, do nothing) 3: return v 4: Let $T(v)$ be the sub-tree rooted at v 5: if $T(v)$ has no opened leaf within dist. f from q then 6: $w = \text{SELECTHEAVIESTCHILD}(T, v, q)$ 7: return PLUCK(T, w, q) 8: else 9: Let x be the child of v s.t. q is a leaf of $T(x)$	10: if $T(x)$ contains an opened leaf then 11: return PLUCK(T, x, q) 12: else 13: Select the child y of v s.t. $T(y)$ contains the opened leaf ℓ closest to q (ties broken arbitrarily) 14: Increase x ’s potential by $d_T(q, \ell)$ 15: if the new potential exceeds the cost f of opening a facility then 16: return PLUCK(T, x, q) 17: else return ℓ
---	--

210 PLUCK is well-defined and it always terminates (see the Supplementary Material for a formal proof).

211 **Notation 3.1.** For every q_p and tree T , PLUCK returns a leaf, denoted from now on as ℓ_p , which is
 212 the facility associated to the point p (recall that the leaves of T are the set \mathcal{S} of suggested facilities).

213 Furthermore, note that ℓ_p is always at tree distance no greater than f from q_p , i.e. $d_T(q_p, \ell_p) \leq f$. And,
 214 if ℓ_p is opened on input q_p by the invocation of SELECTHEAVIESTCHILD, then $d_T(q_p, \ell_p) \leq f/3$.

Algorithm 3 Algorithm SELECTHEAVIESTCHILD(T, v, q)

Input: an HST tree T , a node $v \in T$, and a point $q \in \mathcal{S}$

Output: a child of v in T

1: **return** w such that:

- (i) w is a child of v ;
 - (ii) $T(w)$ has a leaf at distance no greater than $f/3$ from q ;
 - (iii) the number of leaves of $T(w)$ is maximum (among those satisfying (ii), breaking ties arbitrarily)
-

215 3.1 Algorithm Analysis

216 All the proofs have been moved to the Supplementary Material for lack of space. However, here
 217 we recreate the workflow of the proof of Theorem 3.1. To establish Theorem 3.1 we fix an optimal
 218 solution restricted to \mathcal{S} , $\mathcal{F}^* := \{c_1^*, \dots, c_{k^*}^*\} \subseteq \mathcal{S}$, and compare it to the solution returned by
 219 TAKEHEED. The optimal solution naturally induces a clustering of the clients $\mathcal{C}^* := \{C_1^*, \dots, C_{k^*}^*\}$
 220 where $C_i^* := \{p \in \mathcal{P} : c_i^* \text{ is the facility in } \mathcal{F}^* \text{ closest to } p\}$.

221 **Notation 3.2.** Let c_p^* denote from now on the facility in the optimal solution closest to point p . And
 222 let k^* denote the number of facilities opened by the optimal solution.

223 In what follows, we establish a series of bounds that hold for every tree $T \in \mathcal{T}$ in the HST family.
 224 Once this is done, we invoke linearity of expectation to obtain the desired bound. Recall that d is
 225 the distance in the original metric space while d_T is the metric induced by a tree $T \in \mathcal{T}$. We start
 226 with two useful lemmas. The first says that we can bound the service cost of a point p with that of its
 227 proxy q_p in the tree metric d_T . The second relies on the properties of HSTs.

228 **Lemma 3.2.** Let $T \in \mathcal{T}$. Then, $d(p, \ell_p) \leq d(p, c_p^*) + d_T(q_p, \ell_p)$.

229 **Lemma 3.3.** $\mathbb{E}_{T \sim \mathcal{D}}[d_T(q_p, c_p^*)] \leq 16 \log(|\mathcal{S}|)d(p, c_p^*)$.

230 Hence, to bound the cost $d(p, \ell_p)$ for point p , we need to find a bridge between $d_T(q_p, \ell_p)$ and
 231 $d_T(q_p, c_p^*)$. We now partition the clients \mathcal{P} into two groups, those that are ‘‘far’’ from the optimal
 232 centers and those that are ‘‘near’’: $\mathcal{P}_{\text{NEAR}}^T := \{p \in \mathcal{P} : d_T(q_p, c_p^*) \leq f/6\}$ and $\mathcal{P}_{\text{FAR}}^T := \mathcal{P} \setminus \mathcal{P}_{\text{NEAR}}^T$.

233 **Notation 3.3.** Let ℓ_p be the opened facility returned by PLUCK on input q_p . If this is the first time
 234 the leaf was opened, we assign the entire facility cost f to p . Otherwise, we assign zero. Let $f^T(p)$
 235 denote the facility cost assigned to p by this mechanism.

236 Given that T is fixed in our analysis and there is no ambiguity, we will drop the superscript T
 237 for notational convenience. The cost incurred by far away points can be easily bounded through
 238 point-wise bounds, following by the definition of \mathcal{P}_{FAR} .

Lemma 3.4.

$$\sum_{p \in \mathcal{P}_{\text{FAR}}} \left(d_T(q_p, \ell_p) + f(p) \right) \leq 12 \sum_{p \in \mathcal{P}_{\text{FAR}}} d_T(q_p, c_p^*).$$

239 We can focus on $\mathcal{P}_{\text{NEAR}}$ from now on, for which the proofs get more involved. In Lemma 3.2 we
 240 bounded the service cost incurred by the algorithm for each point p with its optimal cost and the
 241 distance between the proxy q_p and the facility ℓ_p . We now do the same for the facility cost.

242 **Lemma 3.5.** Let $T \in \mathcal{T}$.

$$\sum_{p \in \mathcal{P}_{\text{NEAR}}} f(p) \leq 3 \sum_{p \in \mathcal{P}} d_T(q_p, \ell_p) + fk^*$$

243 In view of Lemmas 3.2, 3.3 and 3.5, to close the circle we need to relate $d_T(q_p, \ell_p)$ to $d_T(q_p, c_p^*)$.
 244 Technically this is the challenging step, achieved by the next lemma which requires careful use of the
 245 properties of PLUCK.

Lemma 3.6.

$$\sum_{p \in \mathcal{P}_{\text{NEAR}}} d_T(q_p, \ell_p) \leq 2fk^* \log(|\mathcal{S}|) + 4 \sum_{p \in \mathcal{P}_{\text{NEAR}}} d_T(q_p, c_p^*).$$

246 By combining Lemmas 3.4 – 3.6 a similar bound for the facility costs can be established.

Lemma 3.7. *Let $T \in \mathcal{T}$. Then,*

$$\sum_{p \in \mathcal{P}_{\text{NEAR}}} f(p) \leq (6 \log(|\mathcal{S}|) + 1) f k^* + 18 \sum_{p \in \mathcal{P}} d_T(q_p, c_p^*).$$

247 All the pieces can be assembled to prove Theorem 3.1 (see the Supplementary Material).

248 4 Lower Bounds

249 In this section we present two lower bounds whose proof is deferred to the Supplementary Material.
250 The first shows the approximation guarantee of Theorem 3.1 to be almost optimal. The second shows
251 that, in the worst case, Meyerson [2001]’s well-known algorithm computes a solution whose cost is
252 $\Omega(|\mathcal{S}| \text{OPT})$, as opposed to the $O(\log(|\mathcal{S}|) \text{OPT})$ bound of Theorem 3.1. The same is true for Fotakis
253 [2003] algorithm too. This provides strong motivation for our HST approach.

254 **Theorem 4.1.** *Let $f = 1$ be the uniform facility cost. For every $k \geq 1$, there exist a metric space,
255 k sets of facilities $\mathcal{S}_i, i \in [k]$, (with $|\mathcal{S}| = |\cup_{i \in [k]} \mathcal{S}_i| = k$), and a probability distribution over the
256 input sequences such that, when an input sequence \mathcal{P} is generated according to this distribution:
257 (i) with probability 1, $\min_{i \in [k]} \text{COST}(\mathcal{P}, \mathcal{S}_i) = O(1)$ (i.e. one of the sets suggests a solution of cost
258 $O(1)$) and yet, (ii) each online algorithm provides a solution of expected cost at least $\Omega\left(\frac{\log |\mathcal{S}|}{\log \log |\mathcal{S}|}\right)$.*

259 We now proceed by defining a general class of algorithms for online facility location.

260 **Definition 4.1.** *An algorithm for online facility location is said to be “Meyerson-like” if it only opens
261 facilities at previously seen input points.*

262 Notice that Meyerson [2001] algorithm and Fotakis [2003] algorithms are Meyerson-like according
263 to Definition 4.1, while the Algorithm from Anagnostopoulos et al. [2004] is not.

264 **Theorem 4.2.** *Let $f = 1$ be the uniform facility cost. For every integer $k > 0$, there is a finite metric
265 space (\mathcal{S}, d) of cardinality $|\mathcal{S}| = k$ and input sequences such that: (i) The optimal solution has cost
266 $O(1)$, but (ii) Any Meyerson-like algorithm computes solutions of cost $\Omega(|\mathcal{S}|) = \Omega(k)$.*

267 Notice that $|\mathcal{S}| = k$ is not the size of the input, which could contain repeated points of \mathcal{S} , so these
268 lower bounds have nothing to do with the logarithmic lower bound from Fotakis [2003]. Theorem 4.2
269 further validates our approach, since it rules out any Meyerson-like algorithm. On the other side, we
270 manage to bypass this lower bound by considering as candidate facilities also points from the advice.

271 5 Experiments

272 We studied the performance of TAKEHEED with real-world and synthetic datasets by comparing
273 it against the state of the art online algorithms of Meyerson [2001], Fotakis [2003], and Anagnos-
274 topoulos et al. [2004]. For all algorithms above we used our own implementation (available in the
275 Supplementary Material). Note that the first algorithm is randomized, as it is TAKEHEED. To generate
276 the trees of the HST families we implemented the algorithm of [Fakcharoenphol et al., 2004].

277 As a baseline to compare against we use the 1.52-approximation algorithm for offline facility location
278 by Mahdian et al. [2002], referred to as THEBASELINE in the sequel. For it, we used the Max Planck
279 Institute implementation Hoefer [2002].

280 To test robustness, we use the algorithm of Mahdian et al. [2012]. This algorithm accepts in input
281 two algorithms and mixes them in a way as to ensure that cost of its output is always of the same
282 order of magnitude of the better of the two. The algorithm can also mix a set of points (i.e. advice)
283 with an algorithm, or two sets of points.

284 For synthetic instances, we sampled points from a mixture of Gaussians (from 10 to 100) having
285 centers in $[0, 1]^2$ and variance scaled by the number of mixtures. The metric is the euclidean distance.
286 For real real-world datasets, we consider *Gowalla* and *Brightkite*, from the SNAP Dataset Collection
287 Leskovec and Krevl [2014], and *Uber FiveThirtyEight* [2015]. Each of these contains latitude,
288 longitude, and timestamp of some activity: check-in from the location-based social networks Cho
289 et al. [2011] for the first two (restricted to the US), and pickups in New York City for Uber. For
290 each dataset, an input instance consists of all such activities that occurred in a given day, sorted

Table 1: For the real datasets, an instance consists of the activities of one day. To obtain an input for the algorithm an instance is paired with facility cost.

Dataset	# Instances	Median size	Size range
Mixtures	750	5,000	1,000 - 10,000
Brightkite	901	3,119	5 - 7,140
Gowalla	388	9,354	39 - 25,425
Uber	183	24,550	10,202 - 43,205

291 by timestamp. The metric used is the great-circle distance in kilometers. Table 1 gives additional
 292 information on the datasets.

293 For every input instance, all randomized algorithms were run 10 times. The plots report the average
 294 and variance of the cost of their solutions. All experiments ran on a desktop computer.

295 **Instances.** For synthetic instances, the input sequence and the advice were generated as follows.
 296 First, a set of N points was sampled from the mixture of Gaussians. The input sequence consisted
 297 of these points in random order. The advice was obtained as follows. THEBASELINE was given
 298 the random set in input and produced a set of M facilities in output. Then, M additional random
 299 points were added. Finally, the $2M$ points were randomly partitioned into 5 sets of equal size. In the
 300 experiments, for every mixture, 50 random instances were generated. For real-world datasets, we
 301 divided up every day into four consecutive time windows of 6 hours each, starting from midnight.
 302 An input sequence consisted of all the points (activities) of a given day. The multiple advice was
 303 obtained by running THEBASELINE on each time window of the two previous days, for a total of
 304 8 sets of suggested facilities. In this way, we were able to obtain multiple advice containing both
 305 bad and good suggestions. For every input sequence, we tried different values of f , the facility cost.
 306 Changing f did not affect the outcome in terms of performance comparison. The results we show
 307 exemplify the general picture. Likewise, we ran all the randomized algorithms using different seeds.
 308 Again, this had no discernible consequence on the performance.

309 **Raw advice.** The previous discussion describes how multiple advice was generated. Algo-
 310 rithm RAWADVICE provides a useful baseline to quantify its quality. It operates as follows. Given the
 311 union of the facility sets comprising the advice, each client in the input sequence is assigned to the
 312 closest facility, which is opened. The cost of this solution gives a useful quantitative indication of how
 313 good the multiple advice is. In principle, the algorithm of Mahdian et al. [2012] gives another way
 314 to mix the sets. By generalizing the mixing process, it is possible to work with multiple algorithms,
 315 switching to the best-performing one for each arriving client. In practice, however, just taking the
 316 union gives the best result. This is true for both the normal and robust versions of the algorithm. So
 317 we focus only on the approach merging all the suggestions in what follows.

318 **Preprocessing.** The performance of TAKEHEED hinges upon the tree sampled from the HST family
 319 (recall that we used the algorithm from Fakcharoenphol et al. [2004]). Recall that the points of the
 320 multiple advice induce a metric space for which an HST family exists. To improve performance, we
 321 sampled several trees (10 and 100, respectively, for real and synthetic datasets) and picked the one
 322 that minimized distortion.² The distortion was computed naively, as the maximum over all the pairs
 323 of points, with no attempt at optimization. Even so, the preprocessing time was negligible. We do not
 324 report exact figures because here we are only concerned with the cost of the solutions.

325 **Outcome.** In the sequel we use the term *competitive ratio* in the following sense: the competitive ratio
 326 of algorithm X is the ratio between the cost of the solution given by X and that of THEBASELINE. This
 327 gives a uniform measure of comparison. For Fotakis and ABUV algorithms, which are parametrized,
 328 we tried different values for the parameters and report here only the best ones.

329 In Figure 1 we report the results on synthetic instances. Note that the number of opened facilities
 330 for large values of f – the opening cost – becomes small and the difference in behavior between
 331 TAKEHEED and Advice reduces. In presence of noisy advice, TAKEHEED consistently outperforms
 332 both the other online algorithms and the naive use of the advice. Similar results are obtained for
 333 different numbers of Gaussian mixtures.

334 **Real-world instances.** In Figure 2 we compare TAKEHEED with the competitors for the Brightkite
 335 dataset. Similar results were obtained with all datasets (see the Supplementary Material). The

²the ratio between the metric space distance and the tree distance, for the worst pair of points

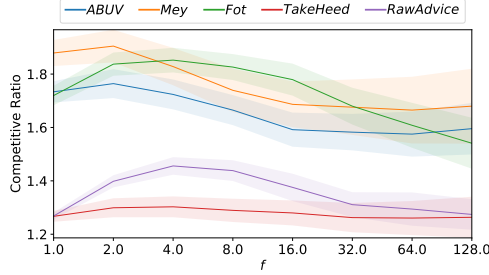


Figure 1: Average estimated competitive ratios on Gaussian mixtures (5000 points from 50 Gaussians) with noisy advice. The shaded area represents one standard deviation.

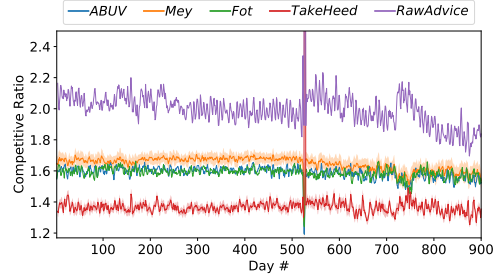


Figure 2: Average estimated competitive ratios on Brightkite with $f = 1000$. The shaded area represents one standard deviation.

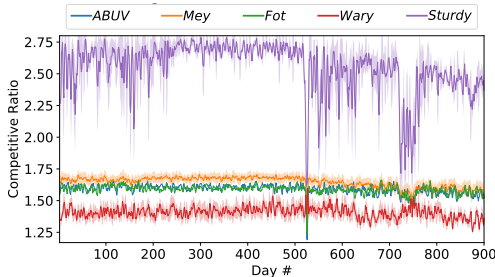


Figure 3: Average estimated competitive ratios of the robust variant on Brightkite with $f = 1000$. The shaded area represents one standard deviation.

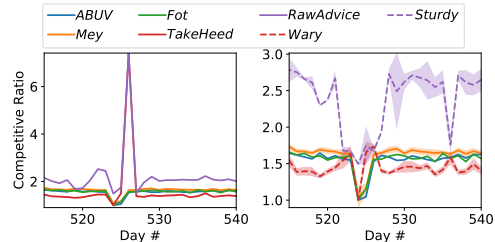


Figure 4: Detail of 2 and 3 around some extreme instances (not smooth). The robust versions manage to recover from bad advice.

336 facility cost was determined in such a way that THEBASELINE opened a number of facilities between
 337 1%–10% of the number of input points. TAKEHEED always outperforms all other algorithms, except
 338 for one outlier instance, to which we return (see the spike occurring around $day = 525$ of the plot).
 339 This good behavior was consistent across all datasets, illustrating the interesting fact that high-quality
 340 multiple advice might be cheaply and consistently produced.

341 **Robustness.** As mentioned, robust online algorithms can be obtained by using the mixing algorithm
 342 of Mahdian et al. [2012]. Algorithm WARY is obtained by mixing our TAKEHEED with Meyerson’s
 343 algorithm. While algorithm STURDY is obtained by simply mixing Meyerson’s algorithm with the
 344 advice set. The mixing algorithm has a parameter γ with which one can give more weight to one of
 345 the two components to be mixed. We show the outcome for $\gamma = 1.75$ which gave best results ($\gamma = 2$
 346 treats them equally, while $\gamma = 1.75$ gives less weight to Meyerson’s algorithm). An in-depth analysis
 347 of the role of this parameter is given in the Supplementary Material.

348 Figure 3 reports the outcome for the same Brightkite dataset, while Figure 4 zeroes in on the spike
 349 of Figure 2. The latter illustrates how the robustness plays out when the advice is bad. While
 350 TAKEHEED suffers, its robust counterpart WARY does as well as online algorithms using no advice,
 351 as predicted by the theory. Different values of the facility cost f produced the same outcome.

352 6 Conclusion

353 We introduced TAKEHEED, WARY, and ROBUST-TAKEHEED, three online algorithms for online
 354 facility location that take advantage in different ways of multiple advice – *i.e.* advice coming from
 355 disparate sources possibly containing noise and misleading information. The algorithms exhibit
 356 good theoretical guarantees as well as good practical performance. Our approach also illustrates that
 357 high-quality multiple advice might be easy and inexpensive to obtain in practical situations, giving
 358 rise to algorithms that can outperform online algorithms that heed no advice. In those seemingly rare
 359 occasions in which good advice is not possible to obtain, WARY, the robust version of TAKEHEED,
 360 has been shown to be able to disregard the poor advice and perform as well as the best online
 361 algorithms. We consider our results as a case study of a more general philosophy that appears to be a
 362 promising approach to tackle a large variety of interesting problems.

363 **References**

- 364 Charu C Aggarwal and Chandan K Reddy. Data clustering. *Algorithms and applications*. Chapman&Hall/CRC Data mining and Knowledge Discovery series, Londra, 2014.
365
- 366 Abu Reyan Ahmed, Md Saidur Rahman, and Stephen Kobourov. Online facility assignment. *Theoretical Computer Science*, 806:455–467, 2020.
367
- 368 Aris Anagnostopoulos, Russell Bent, Eli Upfal, and Pascal Van Hentenryck. A simple and deterministic competitive algorithm for online facility location. *Information and Computation*, 194(2): 175–202, 2004.
369
370
- 371 Yair Bartal. Probabilistic approximation of metric spaces and its algorithmic applications. In *Proceedings of 37th Conference on Foundations of Computer Science*, pages 184–193. IEEE, 1996.
372
- 373 Allan Borodin and Ran El-Yaniv. *Online computation and competitive analysis*. cambridge university press, 2005.
374
- 375 Joan Boyar, Lene M Favrholt, Christian Kudahl, Kim S Larsen, and Jesper W Mikkelsen. Online algorithms with advice: A survey. *ACM Computing Surveys (CSUR)*, 50(2):19, 2017.
376
- 377 Eunjoon Cho, Seth A Myers, and Jure Leskovec. Friendship and mobility: user movement in location-based social networks. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1082–1090. ACM, 2011.
378
379
- 380 Anna Choromanska and Claire Monteleoni. Online clustering with experts. In *Artificial Intelligence and Statistics*, pages 227–235. PMLR, 2012.
381
- 382 Vincent Cohen-Addad, Niklas Oskar D Hjuler, Nikos Parotsidis, David Saulpic, and Chris Schwiegelshohn. Fully dynamic consistent facility location. In *Advances in Neural Information Processing Systems*, pages 3250–3260, 2019.
383
384
- 385 Vincent Cohen-Addad, Benjamin Guedj, Varun Kanade, and Guy Rom. Online k-means clustering. In *International Conference on Artificial Intelligence and Statistics*, pages 1126–1134. PMLR, 2021.
386
387
- 388 Marek Cygan, Artur Czumaj, Marcin Mucha, and Piotr Sankowski. Online facility location with deletions. In *26th Annual European Symposium on Algorithms (ESA 2018)*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2018.
389
390
- 391 Jittat Fakcharoenphol, Satish Rao, and Kunal Talwar. A tight bound on approximating arbitrary metrics by tree metrics. *Journal of Computer and System Sciences*, 69(3):485–497, 2004.
392
- 393 FiveThirtyEight. Uber pickups in new york city. <https://www.kaggle.com/fivethirtyeight/uber-pickups-in-new-york-city>, 2015.
394
- 395 Dimitris Fotakis. On the competitive ratio for online facility location. In *International Colloquium on Automata, Languages, and Programming*, pages 637–652. Springer, 2003.
396
- 397 Dimitris Fotakis. A primal-dual algorithm for online non-uniform facility location. In Panayiotis Bozanis and Elias N. Houstis, editors, *Advances in Informatics*, pages 47–56, Berlin, Heidelberg, 2005. Springer Berlin Heidelberg.
398
399
- 400 Dimitris Fotakis. Online and incremental algorithms for facility location. *SIGACT News*, 42(1): 97–131, March 2011. ISSN 0163-5700. doi: 10.1145/1959045.1959065.
401
- 402 Guojun Gan, Chaoqun Ma, and Jianhong Wu. *Data clustering: theory, algorithms, and applications*. SIAM, 2020.
403
- 404 Sudipto Guha and Samir Khuller. Greedy strikes back: Improved facility location algorithms. *Journal of Algorithms*, 31(1):228 – 248, 1999. ISSN 0196-6774. doi: <https://doi.org/10.1006/jagm.1998.0993>.
405
406

- 407 Xiangyu Guo, Janardhan Kulkarni, Shi Li, and Jiayi Xian. On the facility location problem in
408 online and dynamic models. In *Approximation, Randomization, and Combinatorial Optimization.*
409 *Algorithms and Techniques (APPROX/RANDOM 2020)*. Schloss Dagstuhl-Leibniz-Zentrum für
410 Informatik, 2020.
- 411 Martin Hoefer. Performance of heuristic and approximation algorithms for the uncapacitated facility
412 location problem. Research Report MPI-I-2002-1-005, Max-Planck-Institut für Informatik,
413 Stuhlsatzenhausweg 85, 66123 Saarbrücken, Germany, December 2002.
- 414 Chen-Yu Hsu, Piotr Indyk, Dina Katabi, and Ali Vakilian. Learning-based frequency estimation
415 algorithms. In *International Conference on Learning Representations*, 2019.
- 416 Anil K Jain and Richard C Dubes. *Algorithms for clustering data*. Prentice-Hall, Inc., 1988.
- 417 Bernhard Korte and Jens Vygen. *Facility Location*, pages 629–665. Springer Berlin Heidelberg,
418 Berlin, Heidelberg, 2018.
- 419 Tim Kraska, Alex Beutel, Ed H Chi, Jeffrey Dean, and Neoklis Polyzotis. The case for learned index
420 structures. In *Proceedings of the 2018 International Conference on Management of Data*, pages
421 489–504. ACM, 2018.
- 422 Silvio Lattanzi and Sergei Vassilvitskii. Consistent k-clustering. In *Proceedings of the 34th Interna-*
423 *tional Conference on Machine Learning-Volume 70*, pages 1975–1984. JMLR. org, 2017.
- 424 Silvio Lattanzi, Thomas Lavastida, Benjamin Moseley, and Sergei Vassilvitskii. Online scheduling
425 via learned weights. In *Proceedings of the Fourteenth Annual ACM-SIAM Symposium on Discrete*
426 *Algorithms*, pages 1859–1877. SIAM, 2020.
- 427 Jure Leskovec and Andrej Krevl. SNAP Datasets: Stanford large network dataset collection, June
428 2014.
- 429 Shi Li. A 1.488 approximation algorithm for the uncapacitated facility location problem. *Information*
430 *and Computation*, 222:45 – 58, 2013. ISSN 0890-5401. doi: [https://doi.org/10.1016/j.ic.2012.01.](https://doi.org/10.1016/j.ic.2012.01.007)
431 007. 38th International Colloquium on Automata, Languages and Programming (ICALP 2011).
- 432 Edo Liberty, Ram Sriharsha, and Maxim Sviridenko. An algorithm for online k-means clustering.
433 In *2016 Proceedings of the eighteenth workshop on algorithm engineering and experiments*
434 *(ALENEX)*, pages 81–89. SIAM, 2016.
- 435 Thodoris Lykouris and Sergei Vassilvitskii. Competitive caching with machine learned advice. *arXiv*
436 *preprint arXiv:1802.05399*, 2018.
- 437 Mohammad Mahdian, Yinyu Ye, and Jiawei Zhang. Improved approximation algorithms for metric
438 facility location problems. In *International Workshop on Approximation Algorithms for Combina-*
439 *torial Optimization*, pages 229–242. Springer, 2002.
- 440 Mohammad Mahdian, Hamid Nazerzadeh, and Amin Saberi. Online optimization with uncertain
441 information. *ACM Transactions on Algorithms (TALG)*, 8(1):1–29, 2012.
- 442 Andrew McGregor. Graph stream algorithms: a survey. *ACM SIGMOD Record*, 43(1):9–20, 2014.
- 443 Andrés Muñoz Medina and Sergei Vassilvitskii. Revenue optimization with approximate bid pre-
444 dictions. In *Proceedings of the 31st International Conference on Neural Information Processing*
445 *Systems*, pages 1856–1864. Curran Associates Inc., 2017.
- 446 Ramgopal R Mettu and C Greg Plaxton. The online median problem. *SIAM Journal on Computing*,
447 32(3):816–832, 2003.
- 448 A. Meyerson. Online facility location. In *Proceedings of the 42nd IEEE Symposium on Foundations of*
449 *Computer Science, FOCS '01*, page 426, USA, 2001. IEEE Computer Society. ISBN 0769513905.
- 450 Vahab S Mirrokni, Shayan Oveis Gharan, and Morteza Zadimoghaddam. Simultaneous approxima-
451 tions for adversarial and stochastic online budgeted allocation. In *Proceedings of the twenty-third*
452 *annual ACM-SIAM symposium on Discrete Algorithms*, pages 1690–1701. Society for Industrial
453 and Applied Mathematics, 2012.

- 454 Michael Mitzenmacher. A model for learned bloom filters and related structures. *arXiv preprint*
455 *arXiv:1802.00884*, 2018.
- 456 Michael Mitzenmacher. Scheduling with predictions and the price of misprediction. In *11th*
457 *Innovations in Theoretical Computer Science Conference (ITCS 2020)*. Schloss Dagstuhl-Leibniz-
458 Zentrum für Informatik, 2020.
- 459 Michael Mitzenmacher and Sergei Vassilvitskii. Algorithms with predictions. *arXiv preprint*
460 *arXiv:2006.09123*, 2020.
- 461 Manish Purohit, Zoya Svitkina, and Ravi Kumar. Improving online algorithms via ml predictions. In
462 *Advances in Neural Information Processing Systems*, pages 9661–9670, 2018.
- 463 Dhruv Rohatgi. Near-optimal bounds for online caching with machine learned advice. In *Proceedings*
464 *of the Fourteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 1834–1845. SIAM,
465 2020.
- 466 Jian Zhang, Zoubin Ghahramani, and Yiming Yang. A probabilistic model for online document
467 clustering with application to novelty detection. *Advances in neural information processing*
468 *systems*, 17:1617–1624, 2004.

469 Checklist

- 470 1. For all authors...
- 471 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
472 contributions and scope? [Yes] All the results are stated in Section 3, Section 4 and in
473 the Supplementary Material.
- 474 (b) Did you describe the limitations of your work? [Yes] See Section 1.
- 475 (c) Did you discuss any potential negative societal impacts of your work? [N/A] We study
476 a classical algorithmic problem, and we do not foresee any negative societal impact.
- 477 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
478 them? [Yes]
- 479 2. If you are including theoretical results...
- 480 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Section 2.
- 481 (b) Did you include complete proofs of all theoretical results? [Yes] See the Supplementary
482 Material.
- 483 3. If you ran experiments...
- 484 (a) Did you include the code, data, and instructions needed to reproduce the main exper-
485 imental results (either in the supplemental material or as a URL)? [Yes] See Supple-
486 mentary Material
- 487 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
488 were chosen)? [Yes] See Section 5 for details regarding instance generation
- 489 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
490 ments multiple times)? [Yes]
- 491 (d) Did you include the total amount of compute and the type of resources used (e.g., type
492 of GPUs, internal cluster, or cloud provider)? [Yes] See Section 5
- 493 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 494 (a) If your work uses existing assets, did you cite the creators? [Yes] See Section 5 for
495 dataset and code origin
- 496 (b) Did you mention the license of the assets? [Yes]
- 497 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
498 Code for our algorithms and some of the baselines (see Section 5 for details)
- 499 (d) Did you discuss whether and how consent was obtained from people whose data you’re
500 using/curating? [N/A]
- 501 (e) Did you discuss whether the data you are using/curating contains personally identifiable
502 information or offensive content? [N/A]

503
504
505
506
507
508
509

5. If you used crowdsourcing or conducted research with human subjects...
- (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]