# On Hierarchical Diameter-Clustering and the Supplier Problem

Aparna Das and Claire Kenyon-Mathieu

Brown University, Providence RI 02918, USA

**Abstract.** Given a data set in a metric space, we study the problem of hierarchical clustering to minimize the maximum cluster diameter, and the hierarchical $k$-supplier problem with customers arriving online.
We prove that two previously known algorithms for hierarchical clustering, one (offline) due to Dasgupta and Long and the other (online) due to Charikar, Chekuri, Feder and Motwani, output essentially the same result when points are considered in the same order. We show that the analyses of both algorithms are tight and exhibit a new lower bound for hierarchical clustering. Finally we present the first constant factor approximation algorithm for the online hierarchical $k$-supplier problem.

## 1 Introduction

Clustering is the partitioning of data points into disjoint clusters (or groups) according to similarity [1, 9]. For example if the data points are books, a 2-clustering might consist of the clusters fiction, and non-fiction. In this way clustering can provide a concise view of large amounts of data. In many application domains it is useful to build a partitioning of the data that starts with broad categories which are gradually refined thus allowing the data to be viewed simultaneously at different levels of conciseness. This calls for a *hierarchical* or nested clustering of the data where clusters have subclusters, these have sub-subclusters, and so on. For example a hierarchical clustering might first separate the books into clusters fiction and non-fiction, then separate the fiction cluster into classics and non-classics and the non-fiction cluster into math, science and history, and so on. More formally, a *hierarchical clustering* of $n$ data points is a recursive partitioning of the points into $1, 2, 3, 4, \ldots, n$ clusters such that the $(k+1)$th clustering is obtained by dividing one of the clusters of the $k$th clustering into two parts, thus making the clustering gradually more fine-grained ([5], Section 10.9). This framework has long been popular among statisticians, biologists (particularly taxonomists) and social scientists [10].

A criterion commonly used to measure of the quality of a clustering is the maximum cluster diameter, where the diameter of a cluster is the distance between the two farthest points in the cluster. The goal is to find clusterings which minimize the maximum cluster diameter, thus similar points are placed in the same cluster while dissimilar points are separated. In this paper, we focus on the hierarchical diameter-clustering problem: finding a hierarchical clustering where the value of the clustering is the maximum cluster diameter.

Every associated $k$-clustering of the hierarchical clustering should be close to the optimal $k$-clustering, where the optimal $k$-clustering is the one that minimizes the maximum cluster diameter. The competitive ratio of a hierarchical clustering algorithm $A$ is the supremum, over $n$ and over input sets $S$ of size $n$, of the quantity $\max_{k \in [1,n]} A_k(S)/\mathrm{OPT}_k(S)$, where $\mathrm{OPT}_k(S)$ is the value of the optimal $k$-clustering of $S$ and $A_k(S)$ is the value[1] of the $k$-clustering constructed by algorithm $A$. Thus a hierarchical clustering algorithm with a small competitive ratio, produces $k$-clusterings which are close to the optimal for all $1 \leq k \leq n$.

The hierarchical diameter-clustering problem was studied in work by Dasgupta and Long [4] and by Charikar, Chekuri, Feder and Motwani [2]. A simple and commonly used algorithm for this problem is the greedy "agglomerative" algorithm [5], which starts with $n$ singletons clusters and repeatedly merges the two clusters whose union has smallest diameter. However, it is proved in [4] that this algorithm has competitive ratio $\Omega(\log k)$. The authors then propose a better, constant-factor algorithm, inspired by the "divisive" $k$-clustering algorithm of Gonzalez [6]. The algorithm proposed in [2] is instead "coalescent" and may be partially inspired by a $k$-clustering algorithm by Hochbaum and Shmoys [8]. Superficially, the hierarchical $k$-clustering algorithms presented in the two papers look quite different. Quoting [4]: "the earlier work of [2] uses similar techniques for a loosely related problem, and achieves the same bounds". Indeed, both papers present a 8 competitive deterministic algorithm and a $2e$ competitive randomized variant. Additionally, the algorithm from [2] focuses on online clustering, where points arrive one by one in an arbitrary sequence. We refer to the algorithm from [2] as the *tree-doubling algorithm* and to the algorithm from [4] as the *farthest algorithm*. Here are the main results from [4, 2].

**Theorem 1.** *For the hierarchical diameter-clustering problem,*
*The farthest algorithm is 8-competitive, in its deterministic form and $2e$-competitive in its randomized form [4].*
*The tree doubling algorithm is 8-competitive, in its deterministic form and $2e$-competitive in its randomized form [2].*

Our first contribution is to formally relate the two algorithms. Their specification contains some non-deterministic choices: the farthest algorithm starts from an arbitrary point, and the tree-doubling algorithm considers the points in arbitrary order. Assuming some conditions which remove the non-determinism, we prove that both in the deterministic and in the randomized cases the clustering produced by the farthest algorithm is always a refinement of the clustering produced by the tree-doubling algorithm, where refinement is defined as follows:

**Definition 1.** *A partition $F_1, F_2, \ldots F_l$ is a refinement of a partition $D_1, D_2, \ldots D_k$ iff $\forall i \leq l, \exists j \leq k$ such that $F_i \subseteq D_j$.*

The farthest clustering is a refinement and not equivalent to the tree doubling clustering because the farthest algorithm always outputs a $k$-clustering with

---

[1] If $A$ is randomized, then $A_k(S)$ should be replaced by $E(A_k(S))$.

exactly $k$ clusters, where as the tree doubling algorithm outputs one with at most $k$ clusters.

**Theorem 2. (Refinement)** *Assume that the first two points labeled by the farthest algorithm have distance equal to the diameter of the input. Also assume that the tree-doubling algorithm considers points in the order in which they were labeled by the farthest algorithm. Moreover, in the randomized setting, assume that the two algorithms choose the same random value[2] $r$.*

*Then, for every $k$, the $k$-clustering produced by the farthest (deterministic or randomized) algorithm is a refinement of the $k$-clustering produced by the tree-doubling (deterministic or randomized) algorithm.*

With this interpretation, we see that the competitive ratio of the farthest algorithm can be seen as a corollary of the competitive ratio of the tree-doubling algorithm. Could it be that the farthest algorithm is actually better? We answer this question in the negative by proving that the analysis of the farthest algorithm in [4] is tight.

**Theorem 3. (Tightness)** *The competitive ratio[3] of the deterministic farthest algorithm is at least 8.*

This means that the 8 competitive ratio upper bound for the farthest algorithm is tight, and, by the refinement theorem, the 8 competitive ratio upper bound for the tree-doubling algorithm is also tight. Proving tightness of the randomized variants are open.

Can the competitive ratio be improved? We turn to the question of what is the best competitive ratio achievable for any hierarchical clustering algorithm with no computational restrictions. In other words, what is the best we can expect from a hierarchical clustering algorithm if it is allowed to have non-polynomial running time? We prove that no deterministic algorithm can achieve a competitive ratio better than 2, and no randomized algorithm can achieve competitive ratio better than $3/2$. (Note that the lower bounds proved in [2] apply to the online model only and thus are incomparable to our lower bounds.)

**Theorem 4. (Hierarchical lower bound)** *No deterministic (respectively randomized) hierarchical clustering algorithm can have competitive ratio better than 2 (respectively better than $3/2$), even with unbounded computational power.*

How general are these techniques? In our final contribution, we extend the tree-doubling algorithm to design the first constant factor approximation algorithm for the *online hierarchical supplier problem*.

In the standard (offline, non hierarchical) *k-supplier problem*, we are given a set $S$ of suppliers and a set $C$ of customers and the distances between all

---

[2] The significance of parameter $r$ will be explained in Section 2.

[3] Charikar et al. [2] present a lower bound of 8 for their clustering algorithm, however it applies to the online setting but not to the hierarchical setting and does not extend to the tree-doubling algorithm.

customer-supplier pairs. We wish to select a set $S_k$ of $k$ suppliers and an assignment of each customer $c$ to a supplier $f(c)$ in $S_k$ so as to minimize the maximum distance from any customer to its supplier, $\max_{c \in C} d(c, f(c))$. For example the task of segmenting customers into a small number, $k$, market segments can be modeled as a $k$-supplier problem where the suppliers are a set of fixed templates representing different markets and the goal is to match markets to customers based on their buying patterns (the distance measure). A 3-approximation algorithm for the $k$-supplier problem is mentioned in [7].

In the more difficult *online hierarchical* setting, the set $S$ of suppliers is known in advance but new customers arrive as time goes on, so $C$ is a *sequence* of customers. When a new customer arrives, it is either assigned to one of the existing open suppliers, or a new supplier is opened to serve this customer. If opening a new supplier results in more than $k$ open suppliers then two existing open suppliers merge their customer lists, and one of them closes. This requirement ensures that the hierarchical condition is satisfied, i.e that $S_{i-1} \subseteq S_i$ and that for each supplier $s \in S_i \setminus S_{i-1}$, all the customers assigned to $s$ are assigned to the same supplier in $S_{i-1}$. For example, suppose customers arrive over time to use resources and we would like to dynamically increase/decrease the total number of resources allocated without having to do extensive recomputation. Using the hierarchical supplier solution, this only requires splitting/merging the customers currently assigned to one of the resources. The online hierarchical model is an increasingly important framework for clustering problems, where there is a requirement to gather and categorize large amounts of data on the fly (see [12] for example).

Using the tree-doubling algorithm as a subroutine, we obtain a constant-factor approximation algorithm for the online hierarchical supplier problem. (Note that in the offline case, we could equivalently have used the farthest algorithm as a subroutine. In fact, we conjecture that in the offline case, a similar result may also be obtainable using methods from [3, 11].)

**Theorem 5. (Online hierarchical supplier)** *For the online hierarchical $k$-supplier problem, there exists a deterministic 17-approximation algorithm and a randomized $(1 + 4e) = 11.87$-approximation algorithm.*

## 2 Review of the hierarchical clustering algorithms

### 2.1 The Farthest Algorithm from [4]

The input is a set of $n$ points $\{x_1, \ldots x_n\}$ with associated distance metric $d$. The algorithm has three main steps:

**Labeling the points.** Take an arbitrary point and label it 1. Give label $i$ for $i \in \{2, \ldots, n\}$, to the point which is farthest away from the previously labeled points. Let $d_i$ denote the distance from $i$ to the previous $i-1$ labeled points, i.e $d_i = \min_{1 \leq j \leq i-1} d(i, j)$. Thus $d_2 = d(1, 2)$.

**Assigning levels to labelled points.** For labelled point 1, set level$(1) = 0$. For labelled point $i \in \{2, \ldots, n\}$, set level$(i) = \lfloor \log_2(d_2/d_i) \rfloor + 1$.

**Organizing labelled points into a tree.** Organize the points into a tree referred to as the $\Pi'$-tree. Place point 1 as the root of the $\Pi'$-tree. For each point $i \in \{2, \ldots, n\}$, define its parent, $\pi'(i)$, to be the point closest to $i$ among the points with level strictly less than level($i$). Later, in order to compare the farthest algorithm with the tree doubling algorithm we set a specific tie breaking scheme for choosing the parent for a node. Insert points $i > 1$ into the $\Pi'$-tree in order of increasing levels connecting each point $i$ with an edge to its parent $\pi'(i)$.

The hierarchical clustering is represented implicitly in the $\Pi'$-tree. To obtain the $k$-clustering (of the hierarchical clustering) remove edges $(i, \pi'(i))$, for $i \in \{2, \ldots, k\}$ from the $\Pi'$-tree. Deleting these $k - 1$ edges splits the $\Pi'$-tree into $k$ connected components such that points $\{1, \ldots, k\}$ are in separate components. The components are returned as the $k$ clusters.

It is easy to verify that this defines a hierarchical clustering, and [4] proves that it satisfies the following properties. The distances $(d_i)_i$ are a monotone non-increasing sequence, and the levels $(\text{level}(i))_i$ are a monotone non-decreasing sequence. The definition of levels imply the following bounds on $d_i$.

$$d_2/2^{\text{level}(i)} < d_i \leq d_2/2^{\text{level}(i)-1} \ . \tag{1}$$

In addition [4] proves that:

$$d(i, \pi'(i)) \leq d_2/2^{\text{level}(i)-1} \ . \tag{2}$$

Charikar et al. [4] also present a randomized variant of the farthest algorithm, where the only difference is in the definition of levels. A value $r$ is chosen uniformly at random from the interval $[0, 1]$, and the levels are now defined by: $\text{level}(1) = 0$ and $\text{level}(i) = \lfloor \ln(d_2/d_i) + r \rfloor + 1$. The monotonicity properties are unchanged; and the two inequalities are replaced by the following.

$$e^r d_2/e^{\text{level}(i)} < d_i \leq e^r d_2/e^{\text{level}(i)-1} \quad \text{and} \quad d(i, \pi'(i)) \leq e^r d_2/e^{\text{level}(i)-1} \ . \tag{3}$$

## 2.2 The Tree-Doubling Algorithm from [2]

Here the input consists of a *sequence* of $n$ points $\{x_1, \ldots x_n\}$ with associated distance metric $d$. Let $\Delta$ denote the diameter of the points. The algorithm considers the points one by one in an online fashion and maintains a certain infinite rooted tree which we refer to as the $T^+$ tree. Each node in $T^+$ is associated to a point, and the set of nodes associated to the same point forms an infinite path in the tree. The first point is placed at depth 0 as the root of $T^+$, and a copy of this point is placed at each depth $d > 0$ along with a parent edge to the copy at depth $d - 1$. When a new point $p$ arrives it is inserted at a depth $d_p$, as defined by the insertion rule given below. A copy of $p$ is placed at each depth $d > d_p$ with a parent edge to the copy of $p$ at depth $d - 1$.

**(Insertion rule)** *Find the largest depth $d$ with a point $q$ such that $dist(p, q) \leq \Delta/2^d$. Point $p$ is inserted into depth $d_p = d + 1$ with a parent edge to $q$.*

To obtain a $k$-clustering, find the maximum depth $d$ in $T^+$ which has at most $k$ nodes. Delete all tree nodes at depth less than $d$. This leaves $\leq k$ subtrees rooted at the points at depth $d$. Delete all multiple copies of points from the subtrees and return these as the clusters.

By [2], the following properties are maintained as nodes are added to $T^+$ :

*Property 1.* **(Close-parent property)** Points at depth $d$ are at distance at most $\Delta/2^{d-1}$ from their parents.

*Property 2.* **(Far-cousins property)** Points at depth $d$ are at distance greater than $\Delta/2^d$ from one another.

Charikar et al. [2] also presents a randomized variant, where the only difference is in the insertion rule. A value $r$ is chosen uniformly at random from the interval $[0, 1]$, and the insertion rule is now: *Find the largest depth $d$ that contains a point $q$ such that $dist(p, q) \leq e^r \Delta/e^d$. Point $p$ is inserted into depth $d_p = d + 1$ with a parent edge to $q$.*

The properties are replaced by the following.

*Property 3.* Points at depth $d$ are at distance at most $e^r \Delta/e^{d-1}$ from their parents , and at distance greater than $e^r \Delta/e^d$ from one another.

## 3   Proof of the Refinement Theorem
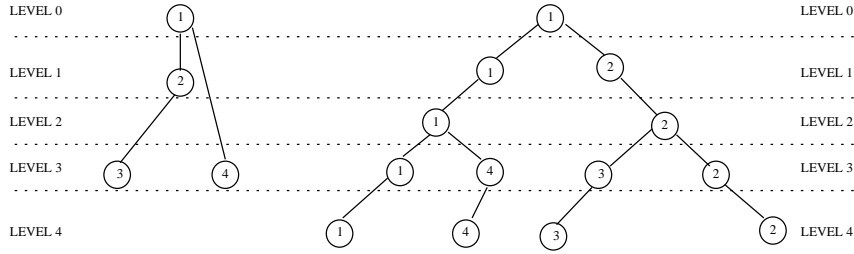
### 3.1   Proof of the Deterministic Version

To relate the farthest and tree doubling algorithms we first make some assumptions about their nondeterministic choices. The farthest algorithm starts its labelling at an arbitrary point. We will assume the first point labelled by the farthest algorithm is at distance $\Delta$ from the second point labelled by the algorithm and thus $d_2 = \Delta$. The tree doubling algorithm receives its input points in an arbitrary order. We assume that points arrive to the tree-doubling algorithm in the order they are labeled by the farthest algorithm. Lastly we assume that ties are broken in the same way by the two algorithms. Specifically if points $q, q'$ are tied to be a parent, the point with the larger label ( in the farthest algorithm) or the point which arrived first (in the tree doubling algorithm) is favored.

Our proof is based on the alternative construction of the tree doubling $T^+$ tree, given in algorithm 1, which builds a tree $T$ based on the farthest algorithm's $\Pi'$ tree. We prove that our construction is consistent with the tree doubling algorithm's insertion rule and hence $T$ could have legitimately been constructed by the tree-doubling algorithm. Finally we argue that the $k$-clustering defined by the $\Pi'$ tree is a refinement of the $k$-clustering defined by $T$. Given the $\Pi'$ tree of the farthest algorithm, the following algorithm constructs a $T^+$ tree. Fig. 1 shows the $\Pi'$ tree and the corresponding tree $T$ constructed by algorithm 1.

**Algorithm 1:** Given $\Pi'$ construct $T^+$
(1)       Let $T$ be an empty tree
(2)       Let $\ell = $ the maximum level of the points in $\Pi'$
(3)     **For** each level $i = 0, \ldots \ell$
(4)           Let $L_i$ denote the points with level $i$
(5)           Let $S = L_0 \cup L_1 \cup \ldots \cup L_i$
(6)           Insert each $p \in S$ at depth $i$ of $T$ with an edge to:
(7)               The copy of $\pi'(p)$ at depth $i-1$, if $level(p) = i$ or
(8)               The copy of $p$ at depth $i-1$, if $level(p) \neq i$
(9)       **Return** $T$

Let $T$ be the tree constructed by Algorithm 1.



**Fig. 1.** The $\Pi'$-Tree is shown on the left and the $T^+$ Tree is shown on the right.

**Lemma 1.** *$T$ satisfies the close-parent property.*

*Proof.* Let $p$ be any point at depth $d$ in $T$.

If $d \neq level(p)$, then by step (8) of algorithm 1, the parent of $p$ in $T$ is the copy of $p$ at depth $d-1$. Thus the close parent property follows trivially in this case since $d(p,p) = 0$.

Otherwise, $d = level(p)$. By step (7) of algorithm 1 parent$(p) = \pi'(p)$. Applying equation 2 with $d_2 = \Delta$, we have that, $d(p, \pi'(p)) \leq \Delta/2^{level(p)-1}$. Substituting $level(p) = d$ and $\pi'(p) = $ parent$(p)$, we get: $d(p, \text{parent}(p)) \leq \Delta/2^{d-1}$. $\square$

**Lemma 2.** *$T$ satisfies the insertion rule.*

*Proof.* By steps (3-6) of algorithm 1 if $level(p) = d$, then $p$ appears in $T$ for the first time at depth $d$ and its parent $q$ is $\pi'(p)$.

Since $level(\pi'(p)) < level(p)$, a copy of $\pi'(p)$ must be at depth $d-1$ in $T$. Since $T$ satisfies the close-parent property, $d(p, \pi'(p)) \leq \Delta/2^{d-1}$. Thus $\pi'(p)$ is qualified (distance-wise) to be the parent of $p$ according to the insertion rule.

To show that insertion rule is satisfied we need to show that when $p$ first arrives, there was no other point at a depth higher than $d-1$ which was close enough to $p$ to be its parent. Let $q'$ be any point at depth $j > d-1$, which

arrived before $p$. Note that $q' \in \{1, \ldots, p-1\}$ since by assumption points arrive in the order they are labelled by the farthest algorithm. We need to show that $d(p, q') > \Delta/2^j$. Note that by definition of $d_p$, $d(p, q') \geq \min_{j \in [1, p-1]} d(p, j) = d_p$. Using the fact that $\text{level}(p) = d$ and equation 1 with $d_2 = \Delta$ we get

$$d_p > \Delta/2^{\text{level}(p)} = \Delta/2^d \ .$$

Combining the two statements above we have that

$$d(p, q') \geq d_p > \Delta/2^d \geq \Delta/2^j \ ,$$

where the last inequality follows since $j > d - 1 \Rightarrow j \geq d$. Since $d(p, q') > \Delta/2^j$ point $q'$ cannot be parent of $p$. □

We have shown that $T$ satisfies the insertion rule and thus it can be constructed by the tree-doubling algorithm when the assumptions of Theorem 2 hold. Thus for the rest of the proof assume that the tree doubling algorithm constructs $T$.

Given $k$, the farthest algorithm removes exactly $k - 1$ edges from the $\Pi'$ tree and returns a clustering $F(k)$ with exactly $k$ clusters. The tree-doubling algorithm looks for the deepest level of the $T$ tree with at most $k$ nodes and thus returns a clustering $D(k)$ with $\leq k$ clusters. We first show the two clusterings $D(k)$ and $F(k)$ are equivalent when they both have exactly $k$ clusters. The refinement property then follows easily.

**Lemma 3.** *Let $k$ be such that the tree doubling tree $T$ has a depth $d$ with exactly $k$ vertices, then the clusterings $F(k)$ and $D(k)$ are the same.*

*Proof.* Let $F_1, \ldots F_k$ be the clusters returned by the farthest algorithm, where $F_i$ contains point $i$. Let $D_1, \ldots D_k$ be the clusters returned by the tree-doubling algorithm, where the cluster are defined by the $k$ points at depth $d$ in $T$. Since depth $d$ contains exactly $k$ vertices, the monotonicity of $(\text{level}(i))_i$ implies that these points must be exactly the points $1, \ldots, k$. We will show that for any $1 \leq i \leq k$ if a point, $x$, is in $D_i$ then $x \in F_i$. Since the $k$-clustering is a partition of the points, this immediately implies that $D_i = F_i$ for all $1 \leq i \leq k$.

Let $x$ be a point in $D_i$. Since $D_i$ contains the points in the subtree under $i$, there is a $i$-to-$x$ path $P = (i = p_1, p_2, \ldots p_l = x)$ in $T$. Let $S = (i = s_1, s_2, \ldots s_m = x)$ be the sequence of points obtained by deleting all repetitions of points from $P$. By the construction of $T$ we have that $s_j = \pi'(s_{j+1})$ for all $1 \leq j \leq m$ which implies that $S$ is a valid $i$-to-$x$ path in the $\Pi'$-tree. Since depth $d$ of $T$ contains all points in $\{1, \ldots, k\}$, only point $i$ can appear in sequence $S$. Thus none of the points $\{1, \ldots, k\}$ except $i$ are in the $i$-to-$x$ path in $\Pi'$ tree. This implies that $x \in F_i$. □

**Corollary 1.** *Let $k_1$ be an input such that $D(k_1)$ has strictly less than $k_1$ clusters. Let $k_2$ be the minimum input such that $k_2 > k_1$ and $D(k_2)$ has exactly $k_2$ clusters. Then $D(k_1) \preceq F(k_1) \preceq D(k_2)$, where $A \preceq B$ stands for "B is a refinement of A".*

*Proof.* Let $k < k_1$ be the number of clusters in $D(k_1)$. Thus $D(k_1) = D(k)$ and $T$ has a level with exactly $k$ vertices. By Lemma 3, $F(k) = D(k)$. By definition of a hierarchical clustering, $F(k) \preceq F(k_1)$ as $k < k_1$. Thus we have $D(k_1) = D(k) = F(k) \preceq F(k_1)$.

Similarly, on input $k_2$, the tree-doubling algorithm produces a clustering with exactly $k_2$ clusters which implies that $T$ has a level with exactly $k_2$ vertices. By Lemma 3, $F(k_2) = D(k_2)$. By definition of a hierarchical clustering, $F(k_1) \preceq F(k_2)$ as $k_1 < k_2$. Thus we have $F(k_1) \preceq F(k_2) = D(k_2)$. □

### 3.2 Proof of the Refinement Theorem, Randomized Version

Suppose the random parameter $r$ in the randomized versions of the farthest and the tree-doubling algorithms is chosen to have the same value. Then Lemma 3 and Corollary 1 also apply to the randomized algorithms. The only change to the analysis is to use inequalities (3) instead of inequalities (2) and (1) in the proof of correctness for algorithm 1.

### 3.3 Nondeterministic Choices

To prove the refinement theorem, we made some assumptions about the nondeterministic choices of the two algorithms. But how much do these choices affect the performance of the algorithms?

The first point chosen by the farthest algorithm determines the value of $d_2$ and this in turn determines the level threshold of the $\Pi'$ tree, i.e level one contains the points which are at distance $[d_2, d_2/2)$ from previously labelled points and level two contains points which are at distance $[d_2/2, d_2/2^2)$ from previously labelled points and so on. The initial point can affect the performance of the farthest algorithm by a factor up to 8 as demonstrated on the example we present in Section 4, Fig. 2. On this example when the farthest algorithm chooses initial point $p_1$ it outputs a 5-clustering which has cost arbitrarily close to 8OPT. However the optimal 5-clustering can be obtained if $p_4$ is chosen as the initial point.

Points arrive to the tree doubling algorithm in an arbitrary order. How much can the ordering of points affect the performance of the tree doubling algorithm? By the refinement theorem, if points arrive in the order labelled by the farthest algorithm, there is always a way to break ties so that the tree doubling clustering is no better than the farthest clustering. However the arrival order of points can help the tree doubling algorithm perform better than the farthest algorithm. We demonstrate this on the tight example presented in 4, Fig. 2. If the points arrive as labelled by the farthest algorithm, the tree doubling and the farthest 5-clustering have cost 8OPT, while if the order starts with $p_2, p_5, p_5'$, then tree doubling can construct the cost 2OPT, 5-clustering:

$$\left\{ (p_2), (p_3), (p_3')(p_1, p_4, p_5, q_1 \ldots q_n), (p_1', p_4', p_5', q_1' \ldots q_n') \right\} .$$

Combining these observations, we see that the farthest algorithm can produce clusterings which are 8 times better than the tree doubling algorithm clusterings if the farthest algorithm starts with the best possible initial point and the tree doubling is given its points in the worst possible ordering. On the other hand the tree doubling clusterings can be 4 times better than the farthest clusterings when its points are ordered favorably and the farthest algorithm starts at the worst possible initial point.

## 4  Proof of the Tightness Theorem

We will prove that, for any $\epsilon > 0$, there exists an input on which the farthest algorithm produces a hierarchical clustering where the $k = 5$ clustering is worse than the optimal 5-clustering by a factor of at least $8 - 4\epsilon$.

Choose any $\epsilon > 0$ and let $n = \lceil 2\log(1/\epsilon) \rceil$. The input set $S$ will have $2n + 9$ points; nine standard points, $p_1, p_1', p_2, p_3, p_3', p_4, p_4', p_5, p_5'$, and $2n$ additional points $q_1, q_1', q_2, q_2', \ldots q_n, q_n'$ with distance as shown in Fig. 2. Note that the distance from $q_i$ to $q_{i+1}$ and the distance from $p_1$ to $q_i$ for $i \in [1, n-1]$ is $1/2^i$ and the same holds for the distance from $q_i'$ to $q_{i+1}'$ and the distance from $p_1'$ to $q_i'$.



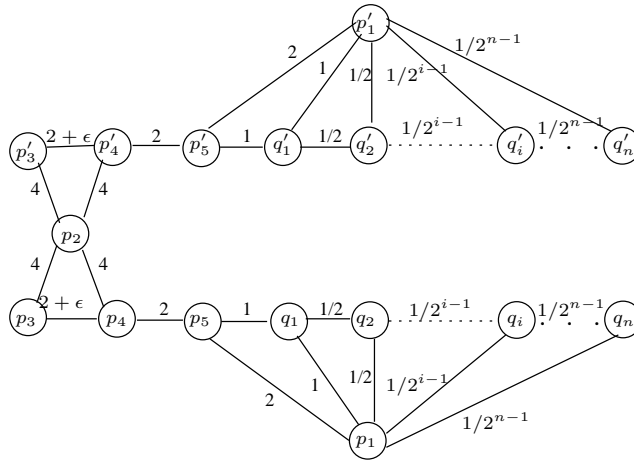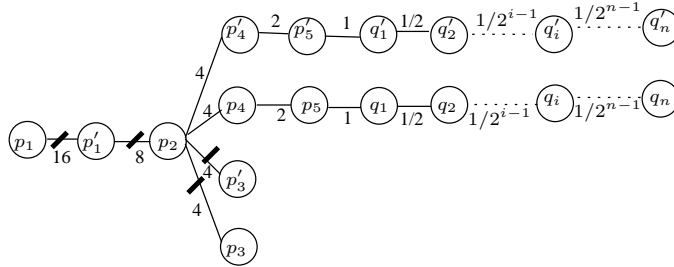**Fig. 2.** Graph for Tight Example

It is easy to verify that the optimal 5-clustering of $S$ is:

$$\left\{ (p_1, p_5, q_1, q_2 \ldots q_n), (p_1', p_5', q_1', q_2', \ldots, q_n'), (p_2), (p_3, p_4), (p_3', p_4') \right\}$$

where clusters $(p_3, p_4)$ and $(p_3', p_4')$ have the largest diameter of $2 + \epsilon = \text{OPT}(5)$.

We carry out the steps of farthest algorithm and show that its 5-clustering can have a cluster of diameter $16 - (2/2^{n-1})$. The algorithm starts with point $p_1$ and obtains the ordering: $p_1, p_1', p_2, p_3, p_3', p_4, p_4', p_5, p_5', q_1, q_1', \ldots q_n, q_n'$. Thus $d(p_1, p_1') = 16 = \Delta$ is used to define the levels for the points. The algorithm connects point $p \neq p_1$ to its parent $\pi'(p)$, the closest point to $p$ at a strictly lower level. The resulting $\Pi'$-tree is shown in Fig. 3.



**Fig. 3.** $\Pi$-Tree for Tight Example

To obtain a 5-clustering, the algorithm removes edges $(p_i, \pi'(p_i))$ for $p_i \in \{p_1', p_2, p_3, p_3'\}$ which yields the clustering:

$$\left\{ (p_1), (p_1'), (p_3)(p_3'), (p_4, p_5, q_1 \ldots q_n, \ p_2, \ p_4', p_5', q_1 \ldots q_n', ) \right\}$$

The diameter of the last cluster is the distance from $q_n$ to $q_n'$ which is $16 - 2/2^{n-1} = 16 - 4\epsilon^2 = (8 - 4\epsilon)\mathrm{OPT}(5)$. This proves the Tightness theorem.
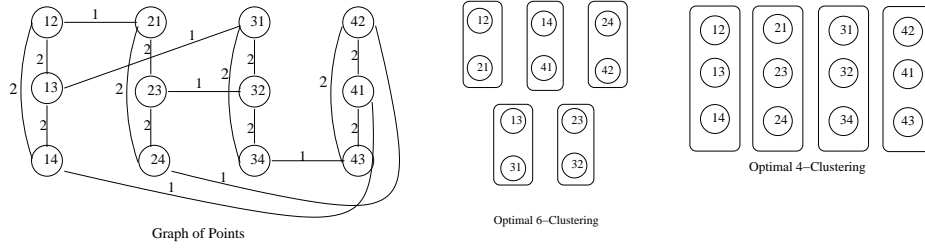
## 5 Proof of the Hierarchical Lower Bound Theorem

We demonstrate an input set $S$ on which every deterministic algorithm obtains a competitive ratio at least 2 and every randomized algorithm obtains a competitive ratio at least 1.5. $S$ has points $p_{ij}$ for $i, j \in [1, 4]$ and $i \neq j$ with distances, $d(p_{ij}, p_{ji}) = 1$ and $d(p_{ij}, p_{ik}) = 2$ as shown in Fig. 4. [4] It is easy to verify that the optimal 6-clustering consists of the six pairs $p_{ij}p_{ji}$ each of diameter 1. Let $B_i = \{p_{ij} | j \in [1, 4], i \neq j\}$. Observe that $B_i$ for $i \in [1, 4]$ is the optimal 4-clustering with each cluster having diameter 2.

### 5.1 The deterministic Lower Bound

Let $A$ be any deterministic hierarchical clustering algorithm.

---

[4] The input set $S$ is also used in [2] to derive a lower bound for the online setting, but the edge weights (distances) are different. The lower bound from [2] is an adversarial argument where the choice of the next arriving points depends on the current clustering and so it does not extend to the offline setting.

**Fig. 4.** Lower Bound Example

**Case 1:** Suppose $A$ produces the optimal 6-clustering. Then $A$'s clusters must be the 6 pairs $p_{ij}p_{ji}$. Since $A$ is a hierarchical clustering algorithm, it must merge some of these pairs to obtain the 4-clustering. Merging any two of these pairs results in a cluster of diameter 4, giving us a competitive ratio of at least $A(4)/\text{OPT}(4) = 4/2 = 2$.

**Case 2:** Suppose $A$ does not produce the optimal 6-clustering. Then some cluster in $A$'s 6-clustering consist of points other than some pair $p_{ij}p_{ji}$. This cluster must have diameter $\geq 2$. Thus the competitive ratio for $A$ is at least $A(6)/OPT(6) \geq 2/1 = 2$.

### 5.2 The randomized Lower Bound

Let $B$ be any randomized hierarchical clustering algorithm. Let $p$ be the probability that $B$ outputs the optimal 6-clustering. Thus the maximum diameter is 1 with probability $p$, and at least 2 with probability $1 - p$. (See analysis for the deterministic scenario). We compute the expected competitive ratio of $B$ for $k = 4$ and $k = 6$, and by definition, the expected competitive ratio of $B$ over all values of $k$ is at least the maximum of these two values.

For $k = 4$ the competitive ratio is

$$E(B_4(S))/\text{OPT}_4(S) \geq (4p + 2(1 - p))/2 = 1 + p \ .$$

where the first inequality follows from the fact that when $B$ chooses the optimal 6-clustering, its 4-clustering will have a cluster of diameter $\geq 4$.

For $k = 6$ the competitive ratio is

$$E(B_6(S))/\text{OPT}_6(S) \geq (p + 2(1 - p))/1 = 2 - p \ .$$

The expected competitive ratio is $\max(1 + p, 2 - p) \geq 1.5$.

## 6  Proof of the Online Hierarchical Supplier Theorem

Our online algorithm for $k$-supplier will use the (online) tree-doubling algorithm as a subroutine. Note that for the off-line hierarchical supplier problem we can

use either the tree-doubling or the farthest algorithm and achieve the same performance guarantees described below. In fact, we conjecture that in the offline case, a similar result may also be obtainable using methods from [3, 11].

## 6.1 The Algorithm

We denote a supplier as *active* if it is the closest supplier to one of the current customers. Throughout the algorithm, we will maintain a hierarchical clustering of the active suppliers by inserting them into the (deterministic or randomized) tree-doubling algorithm tree $T^+$.

When a new customer $c$ arrives, we find the supplier $s$ who is closest to $c$. If $s$ is not yet in $T^+$, we mark $s$ as an active supplier and add $s$ to $T^+$ (using the deterministic or randomized tree-doubling algorithm).

To obtain a hierarchical $k$-supplier solution, find the largest depth $d$ in $T^+$ which contains $k' \leq k$ active suppliers $s_1, s_2, \ldots s_{k'}$ and output these suppliers. For each customer $c$ with closest supplier $s_0$, assign $c$ to $s_i$ for $i \in [1, k']$, if $s_0 = s_i$ or if $s_0$ is in the subtree below $s_i$ in depth $d$ of $T^+$.

## 6.2 The Deterministic Analysis

Suppose $d$ is the largest depth containing at most $k$ active suppliers. Let $s$ (at depth $d$) be the supplier that customer $c$ was assigned to and $s_0$ be the active supplier that $c$ is closest to. Then there is a $s_0$-to-$s$ path in $T^+$. Let $s_0, s_1, \ldots s_p$ be the sequence of the suppliers on the $s_0$-to-$s$ path, where $s_p = s$. By the triangular inequality, the distance from $c$ to $s$ can be bounded as:

$$d(c, s) \leq d(c, s_0) + \sum_{i=0}^{p-1} d(s_i, s_{i+1}) \ . \tag{4}$$

Let $\Delta$ be the maximum distance between any two suppliers. By the close-parent property of the tree-doubling algorithm, the distance from $s_i$ to $s_{i+1}$ for $i \in [0, p-1]$ is at most $\Delta/2^{depth(s_i)-1}$. Since the depths of suppliers on the $s_0$-to-$s$ path are strictly decreasing, and $s_{p-1}$ is on level $d+1$, we have that,

$$\sum_{i=0}^{p-1} d(s_i, s_{i+1}) \leq \frac{\Delta}{2^{depth(s_{p-1})-1}}(1 + 1/2 + 1/4 + \ldots) \leq 2\frac{\Delta}{2^d} \ . \tag{5}$$

Now we derive two lower bounds for $\mathrm{OPT}_k$. First, since $s_0$ is the closest supplier to $c$, we have that $\mathrm{OPT}_k \geq d(c, s_0)$. Next, since $d$ is the largest depth in $T^+$ with at most $k$ active suppliers, depth $d+1$ contains at least $k+1$ active suppliers, $s_1, s_2, \ldots, s_{k+1}$. Using Lemma 4, we have $\mathrm{OPT}_k \geq \delta/4$ where $\delta = \min_{1 \leq i < j \leq k+1} d(s_i, s_j)$ . By the Far-Cousins property of $T^+$, $\delta$ is at least $\Delta/2^{d+1}$. Applying these bounds we obtain

$$d(c, s) \leq d(c, s_0) + \frac{2\Delta}{2^d} \leq \mathrm{OPT}_k + 4\delta \ .$$

Lemma 4 below shows that $\delta \leq 4\mathrm{OPT}_k$. Thus the final result that $d(c, s) \leq 17\mathrm{OPT}_k$ follows as a corollary to lemma 4.

**Lemma 4.** *Let $d$ be the largest depth in $T^+$ with at most $k$ active suppliers and let $s_1, s_2, \ldots, s_{k+1}$ be active suppliers at depth $d+1$. Let $\delta = \min_{1 \leq i < j \leq k+1} d(s_i, s_j)$, and $OPT_k$ be the maximum distance from a customer to a supplier in the optimal $k$-supplier solution. Then $\delta \leq 4OPT_k$.*

*Proof.* Since suppliers $s_1, s_2, \ldots, s_{k+1}$ are active, each of them is the closest supplier to some customer $c_i$. The solution $\mathrm{OPT}_k$ uses at most $k$ suppliers, so it will have to assign two of those customers, $c_i$ and $c_j$, to the same supplier $s^*$. Thus,

$$\mathrm{OPT}_k \geq \max(d(c_i, s^*), d(c_j, s^*)) \geq (d(c_i, s^*) + d(c_j, s^*))/2 \ .$$

Applying the triangle inequality on $d(s_i, s_j)$ we have that:

$$d(s_i, s_j) \leq d(s_i, c_i) + d(c_i, s^*) + d(s^*, c_j) + d(c_j, s_j)$$

Using the fact that $s_i$ is the closest supplier to $c_i$ and $s_j$ is closest for $c_j$, we obtain

$$\delta \leq d(s_i, s_j) \leq 2(d(c_i, s^*) + d(c_j, s^*)) \leq 4\mathrm{OPT}_k \ .$$

$\square$

### 6.3 The Randomized Analysis

Equation 4 still holds. Instead of Equation 5 we now have:

$$\sum_{i=0}^{p-1} d(s_i, s_{i+1}) \leq \frac{e^r \Delta}{e^{depth(s_{p-1})-1}}(1 + 1/e + 1/e^2 + \ldots) \leq \frac{e}{e-1}\frac{e^r \Delta}{e^d} \ .$$

Now, by Property 3 the minimum distance $\delta$ between $s_1, \ldots, s_{k+1}$ satisfies $e^r \Delta/e^{d+1} < \delta \leq e^r \Delta/e^d$. Write $\delta = e^\epsilon e^r \Delta/e^{d+1}$, where $\epsilon$ is distributed uniformly in $[0, 1)$. In expectation we have

$$E(e^r \Delta/e^{d+1}) = \delta \int_0^1 e^{-\epsilon} d\epsilon = \delta \frac{e-1}{e} \ .$$

Lemma 4 still holds, so we finally get:

$$E(d(c, s)) \leq d(c, s_0) + \frac{e}{e-1}E(\frac{e^r \Delta}{e^d}) \leq \mathrm{OPT}_k + \frac{e}{e-1}e\delta\frac{e-1}{e} \leq (1+4e)\ \mathrm{OPT}_k \ .$$

# 7 Conclusions and open questions

Hierarchical clustering provides a useful way to view large amounts of data in an organized manner and is popular among statisticians, biologists and social scientists [10]. In this work we have studied this problem when the objective is to minimize the maximum cluster diameter and showed that two previously known algorithms, tree doubling [2] and farthest algorithm [4] produce essentially the same output. Both algorithms also work for the closely related objective of minimizing the maximum cluster radius, where the radius of a cluster is defined as the maximum distance from a designated center point to other points in the cluster. Our result that the farthest clustering is a refinement of the tree doubling clustering extend to radius objective. We also showed that the analyses of both algorithms are tight under the diameter objective. However the example we used to show tightness does not extend to the radius objective and finding a tight example for this remains an open question.

We exhibited new lower bounds for hierarchical clustering for the diameter objective. Do similar lower bounds exist for the radius objective? The derivation of our bounds assumed no computational restrictions on the algorithm, however it might be possible to get stronger lower bounds by placing such restrictions.

# References

1. P. Arabie, L. J. Hubert and G. De Soete, editors. *Clustering and Classification.* World Scientific, River Edge, NJ, 1998.
2. M. Charikar, C. Chekuri, T. Feder and R. Motwani. Incremental clustering and dynamic information retrieval. *SIAM Journal on Computing*, 33(6):1417–1440, 2004.
3. Marek Chrobak, Claire Kenyon, John Noga and Neal E. Young. Online Medians via Online Bribery. *Lecture Notes in Computer Science* 3887:311-322 (2006); Latin American Theoretical Informatics, 2006.
4. S. Dasgupta, P. Long. Performance guarantees for hierarchical clustering. *Journal of Computer and System Sciences,* 70(4):555-569, 2005.
5. R. O. Duda, P. E. Hart, and D. G. Sork. *Pattern Classification.* Wiley and Sons, 2001.
6. T. F. Gonzalez. Clustering to Minimize the Maximum Intercluster Distance. In *Proceedings of the 17th Annual ACM Symposium on the Theory of Computing*, 38:293-306, 1985.
7. D.S Hochbaum. Various Notions of Approximations: Good, Better, Best and More. In D.S Hochbaum, editor, *Approximation Algorithms for NP-Hard Problems.* PWS Publishing Company. 1996.
8. Dorit S. Hochbaum and David B. Shmoys. A best possible heuristic for the k-center problem. *Mathematics of Operations Research*, 10:180–184, 1985.
9. A. K. Jain and R. C. Dubes. *Algorithms for Clustering Data.* Prentice Hall, Englewood Cliffs, NJ, 1988.
10. L. Kaufman and Peter J. Rousseeuw, *Finding Groups in Data: An Introduction to Cluster Analysis*, Wiley, NY, 1990.

11. Guolong Lin, Chandrashekhar Nagarajan, Rajmohan Rajamaran, and David P. Williamson. A general approach for incremental approximation and hierarchical clustering. In *Proceedings of the seventeenth annual ACM-SIAM Symposium on Discrete algorithm (SODA)*, pages 1147-1156, 2006.

12. NSF Workshop Report on *Emerging Issues in Aerosol Particle Science and Technology* (NAST), UCLA, 2003, Chapter 1, Section 18, "Improved and rapid data analysis tools (Chemical Characterization)". Available at `http://www.nano.gov/html/res/NSFAerosolParteport.pdf`.