



**IJITCE**

**ISSN 2347- 3657**

# International Journal of Information Technology & Computer Engineering

[www.ijitce.com](http://www.ijitce.com)



**Email : [ijitce.editor@gmail.com](mailto:ijitce.editor@gmail.com) or [editor@ijitce.com](mailto:editor@ijitce.com)**

# Nearest neighbors search using point location in balls with applications to approximate Voronoi decompositions <sup>☆</sup>

Ms. R. Latha Priyadharshini, Mr. S. Satheesh, Mrs. M. Indra Priya, Ms. R. Roopa

Associate Professor <sup>3</sup> Assistant Professor <sup>1,2,4</sup>

[lathapriyadharshini.r@actechnology.in](mailto:lathapriyadharshini.r@actechnology.in), [ssatheesh@actechnology.in](mailto:ssatheesh@actechnology.in), [indrapriya.m@actechnology.in](mailto:indrapriya.m@actechnology.in),  
[roopa.r@actechnology.in](mailto:roopa.r@actechnology.in)

Department of CS & BS, Arjun College of Technology, Thamaraikulam, Coimbatore-Pollachi Highway,  
Coimbatore, Tamilnadu-642 120

## Abstract

Our proposed solutions to the nearest neighbour searching problem for Point Location in Balls are an alternative to Sariel Har-Peled's recent work on Approximate Voronoi Diagrams, which maintains the logarithmic search time, and reduces the space bound to linear. This work is published in the Proc. of IEEE FOCS in 2001 and can be accessed online at <http://www.uiuc.edu/~sariel/papers>. We achieve this by streamlining the building of the algorithm that reduces the number of balls generated by it to  $O(n \log n)$ , as described in [S. Har-Peled, A replacement for Voronoi diagrams of near linear size, in: Proc. of IEEE FOCS, 2001, pp. 94-103, full version available from <http://www.uiuc.edu/sariel/papers>]. In order to accomplish linear space decomposition for closest neighbour searches, we further decrease the ball count by introducing a novel hierarchical decomposition strategy and expanding upon PLEBs. Our data structures are constructed with a temporal complexity of  $O(n \log n)$ .

**Keywords:** Voronoi diagrams; Approximate nearest neighbor; Data structures

## 1. Introduction

I will use a metric space and a collection of  $n$  points as  $P$ . Constructing a  $P$ -based data structure that effectively permits closest neighbour searches is one of the most basic challenges with varied applications. This issue has workable answers in two-dimensional planar space, but it becomes non-trivial in three-dimensional and higher-dimensional Euclidean spaces [5,9]. An overarching strategy for nearest neighbour searching is to construct a Voronoi diagram of  $P$  (Hereinafter  $\text{Vor}(P)$ ), a space partition that includes all points that are closer to a given point of  $P$  than any other point in  $P$ . Each cell in the diagram contains one such point. Computer scientists in computational geometry rely on voronoi diagrams for a wide variety of tasks, such as surface reconstruction, learning, motion planning, clustering, and more. Voronoi diagram generalisations to objects other than point sets are also very helpful (for a comprehensive review, see Aurenhammer [6]). Despite having many uses, voronoi diagrams have a major flaw: they are very structurally complicated.

computational bottleneck and the resulting ity. The exponentially complex worst-case complexity in  $R_d$  is known to be  $\theta(n^{\lceil d/2 \rceil})$  with constants in  $d$ . As a result, it becomes computationally impossible to generate and store these diagrams as their dimensions increase.

Therefore, in recent years, researchers have focused on finding alternative solutions to closest neighbour searches in greater dimensions [1,4,8,12–14]. The related relaxed issue of approximate nearest neighbour search has gained prominence and has showed promise of practical solutions because the exact problem looks closely tied to the complexity of Voronoi diagrams. We want to get a data structure by preprocessing the dataset  $P$ , given a constant  $\epsilon < 1$ .  $D$ ,

that, given a query point  $q$ , efficiently returns a point  $P$  such that the distance between  $p$  and  $q$ , denoted as  $pr P$ , is less than or equal to  $(1 \mp \epsilon)$  times the distance between  $pr$  and  $q$ .  $P$  is an  $\epsilon$  closest neighbour ( $\epsilon$ -NN) of  $q$ , and dist might be any arbitrary measure in this case.

The work of Indyk and Motwani [12], where they simplified the issue of  $\varepsilon$ -NNS to Approximate Point Location in Equal Balls ( $\varepsilon$ -PLEB), is particularly relevant to this research.

**Definition 1.1.** When given a query point  $q$  and a parameter  $r$ , an  $\varepsilon$ -PLEB( $P, r$ ) data structure finds a point  $p$  in  $P$  such that  $\text{dist}(p, q) < r$ , if one exists, given  $P$  and  $r$ . Returns nil if there is no point  $p \in P$  such that  $\text{dist}(p, q) < r(1 + \varepsilon)$ . Anything else may result in a return. It is referred to as PLEB( $P, r$ ) when  $\varepsilon$  equals zero. A new ring-cover tree was used by Indyk and Motwani to decrease  $\varepsilon$ -NNS to  $\varepsilon$ -PLEB. A basic hash-based method is used to resolve the  $\varepsilon$ -PLEB. Nevertheless, their reduction was extremely complicated; however, Har-Peled [11] made a complete improvement by presenting an alternate space decomposition known as an Approximate Voronoi Diagram (AVD), which allows for approximate nearest neighbour searching and has a much lower space complexity. His findings are summarised below.

**Theorem 1.1 (Har-Peled).** *Given  $P$  of  $n$  points in  $\mathbb{R}^d$  and a parameter  $\varepsilon > 0$  one can compute a set  $\mathcal{C}(P)$  of  $O(n^{\frac{\log n}{\varepsilon^d}} \log n/\varepsilon)$  regions where each region is a cube or an annulus of cubes. The regions of  $\mathcal{C}(P)$  are disjoint and cover the space and each region has an associated point  $p \in P$ , so that for any point  $q \in \mathcal{C}$ , the point  $p$  is a  $\varepsilon$ -NN of  $q$  in  $P$  and it can be computed in  $O(\log(n/\varepsilon))$  steps. The time for constructing the data structure is  $O(n^{\frac{\log n}{\varepsilon^d}} \log^2 n/\varepsilon)$ .*

Har-Peled also used a more sophisticated data structure for solving  $\varepsilon$ -NNS based on the BBD data structure of Arya et al. [1].

### 1.1. Our results

If AVD could be made linear-sized instead of near linear, it would be a significant improvement above Har-Peled's previous work. In this study, we generalise the concept of PLEBs to Point Location in Smallest Ball (PLSB), where the balls may have variable radii, and we accomplish space reduction by extending some of the principles of Har-Peled. For more specific definitions, see Section 4, however we also provide an approximation of it. This issue has also been implicitly—and more ad hoc—addressed by Har Peled [11].

It should be noted that PLEB is an example of PLSB where the radii of each ball is equal.

**Lemma 1.1.** *An  $\varepsilon$ -PLSB on  $n$  balls can be solved for  $d$ -dimensional space endowed with metric  $l_p$ , in  $O(\log n)$  query time and  $O(n/\varepsilon^d)$  space.*

Appendix C provides a synopsis of the approach, which is comparable to Har-Peled [11]. We shall limit the number of balls in the data structure, which is directly proportional to the space complexity, in the remaining work.

In the part that follows, we generate a set of balls straight from the clustering and decrease their number by a factor of  $\log n$ , thereby eliminating the recursive nature of the Har-Peled's approach. There is no change to the logarithmic search time and an improvement of  $O(\log n)$  to the preprocessing time.

on Section 4, we use a novel hierarchical clustering approach, which is also grounded on the MST of  $P$ , to enhance the space bound to linearity. Using a quicker approximation MST approach that requires  $O(n \log n)$  time—the same as Har-Peled—we are able to minimise the preprocessing time, which is practically quadratic due to the precise MST construction. As a general rule, in Har-Peled's data structure, we combined two clusters whenever their MST edge was the smallest. A comparison of hierarchical

This clustering method differs from Har-Peled's in that, rather than attempting to combine the two closest clusters, we merge many clusters in stages such that the resulting new cluster has a small diameter (Section 4).

Using Callahan and Kosaraju's [7] well separated space decomposition, Arya and Malamatos [2] independently obtained a comparable result. Additionally, they achieve space-time trade-offs for generalising the Voronoi cells with multiple

representatives, an improvement that has been made by Arya et al. [3]. A potential

The number of balls produced by our technique is much lower— $O(n)$ —than that of

A other search structure, not the BBD tree, could be able to take use of this, and the time complexity would be  $O((\alpha/d)n)$ .

## 2. A brief review of Har Peled's construction

Har-Peled [11], proposes a solution based on the following observation:

**Observation 2.1.** *Let the distance between two points  $A$  and  $B$  be  $|AB|$ . If a query point  $Q$  lies further than  $k \cdot |AB|$  from  $A$  and  $B$ , where  $k$  is set to  $1/\epsilon$ , then both  $A$  and  $B$  are  $\epsilon$ -neighbors of  $Q$ , i.e.,  $\frac{|QA|}{|QB|} < 1 + \epsilon$ .*

In addition, a clustering approach for organising the points into a hierarchical tree is suggested by Har-Peled [11]. If you were to grow balls around all the spots, you would end up with linked components. Further, we use the linked parts to build a forest, also known as a cluster tree.

At the outset, we work with only the collection of points. Accordingly, we begin with  $n$  linked components that comprise the individual points. There are  $n$  trees in the matching forest, and each tree has one leaf that represents a point. The balls are then evenly grown around all the spots. At the spots where two components' balls touch, a new component is formed by merging the two components, and this new component takes the place of the two components at those locations. By randomly hanging one tree on top of the other, we combine the two trees in the matching forest that represent the merging components (see Fig. 1).

Now we additionally attach a value of  $r_{\text{loss}}(p)$  to the tree's root,  $p$ , from the other tree. When these parts come together, this number represents the ball's radius. It is half of the distance between the two components' nearest locations, in fact.

Thus  $r_{\text{loss}}(p)$  = radius of the ball at  $p$ , when  $p$  ceases to be root.

### 2.1. Properties of the cluster tree

- Values of  $r_{\text{loss}}$  increase along a path towards the root.
- If  $L_{\text{max}}$  is the longest edge of the MST, the tree edge corresponding to this is the last edge to be added.

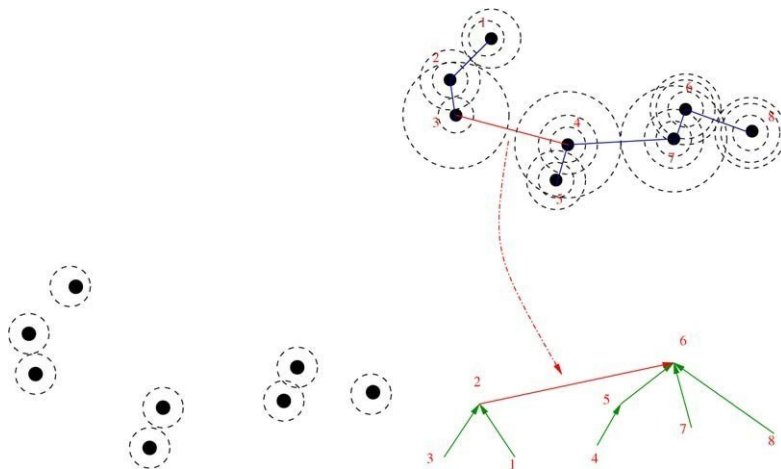


Fig. 1. Constructing the cluster tree.

- If any query point  $q$  that is outside of the union of balls of radii  $L_{\max} \cdot |P| \cdot 1/\epsilon$  centered at  $p \in P$ , then any point  $p$  is an  $\epsilon$ -nearest neighbor of  $q$ , i.e.,  $q$  is too far from  $P$ . Note that the diameter of  $P$  is bounded by  $L_{\max} \cdot |P|$ .

Any subtree (cluster) that forms during building also has the aforementioned qualities.

It requires  $O(n^2)$  time to calculate an accurate  $d$ -dimensional MST for high  $d$ , but Har-Peled settles for a  $nd$ -factor approximation that can be done in  $O(n \log n)$  time.

In order to answer the estimated closest neighbour queries in  $O(\log(n))$  time, he uses the  $nd$ -MST to build a  $nd$ -stretch Hierarchical Clustering and then provides a technique for building a family of PLEBs (Point Location in Equal Balls). The search continues by iteratively scanning the remaining points after trimming certain clusters based on some PLEBs. We need to build  $O(\log n)$  PLEBs for every recursive call, and there are  $O(\log n)$  levels total. At last, a compressed quadtree [1] is created by merging all of the  $O(n \log^2 n)$  balls. The space limit is about  $O(n \log^2 n)$  and is dependent on the number of balls created by the PLEB family.

## Improving the space bound

- 2.2. Here we provide an alternative viewpoint, explaining how we might reduce the amount of balls needed to answer approximate closest neighbour queries by removing overlapping and unnecessary balls from the collection of balls centred around a point  $p \in P$ . The approximate closest neighbour question may be answered by providing the centre of the smallest ball that includes a query point, and we provide an independent approach for building balls around  $P$  points.

### *Building the $nd$ -hierarchical clustering*

The following definitions are similar to Har-Peled that have been restated for convenience of the reader.

**Definition 3.1 ( $\lambda$ -MST).** A tree  $T$  having the points of  $P$  in its nodes, is a  $\lambda$ -MST of  $P$ , if it is a spanning tree of  $P$ , having a value  $\text{len}(\cdot)$  associated with each edge of  $T$ , such that for any edge  $e = uv$  of  $T$ ,  $\text{len}(e) \leq d(P_1, P_2) \leq \lambda \text{len}(e)$ , where  $P_1, P_2$  is the partition of  $P$  into two sets as induced by the two connected components of  $T - e$ , and  $d(P_1, P_2) = \min_{x \in P_1, y \in P_2} \|xy\|$  is the distance between  $P_1$  and  $P_2$ .

One can compute a  $nd$ -MST of  $P$  in  $O(n \log n)$  time. Har-Peled [11] shows how such an  $nd$ -MST can be constructed, and is similar to the fair-split tree construction of Callahan and Kosaraju [7].

**Definition 3.2 ( $\lambda$ -Stretch hierarchical clustering)** For any point  $p$  in a directed tree  $T$ , there exists an out-edge from  $p$  to its parent with a value  $\text{rloss}(p)$  that satisfies the following property:

If we have  $C_1$  and  $C_2$  as connected components obtained from building balls of radius  $r$  around all points and  $F$  as the forest obtained by removing all edges larger than  $r$  from  $T$  and  $X$  as the partition of  $P$  into subsets induced by the connected components of the forest, then  $C_1 \times C_2$ , where  $X \times Y$ , and  $A$  and  $B$  are in the same component in  $X$  if and only if they are in the same component in  $Y$ .

The following is the process for constructing the  $nd$ -stretch hierarchical clustering using the  $nd$ -MST:

We arrange the  $nd$ -edges MST's in ascending order of their  $\text{len}(\cdot)$  values. Afterwards, we go down the list. By inserting an edge between the roots of the two subtrees, we are able to link the two subtrees to which  $x$  and  $y$  belong in each iteration. We do this by taking the next edge,  $e = xy$ , from the sorted list and hanging the smaller tree onto the bigger one. With this edge  $e$ , we connect the value  $\text{rloss}(e) = \text{len}(xy)/2$ .

At the very beginning of the tree, at the root  $t$ , we set  $\text{rloss}(t) = \max(\text{rloss}(p_i))$ , where  $p_i$  is a child of  $t$ .

**Observation 3.1.** The height of the tree formed in the  $\lambda$ -stretch hierarchical clustering of  $P$  is at most  $\log n$ . (Hanging the smaller tree below the larger one insures this property.)

We will use  $b(p, r)$  to denote the ball of radius  $r$  centered at  $p$ .



**Definition 3.3** ( $r_{\text{death}}, r_{\text{low}}, r_{\text{high}}$ ). For an approximation factor  $\gamma$ , define  $r_{\text{death}}(p) = 6\lambda r_{\text{loss}}(p)n \log n/\gamma$  and define  $r_{\text{low}}(p) = (1/(1 + \gamma/3))r_{\text{loss}}(p)$ .  $r_{\text{low}}(p)$  gives us a value just smaller than  $r_{\text{loss}}(p)$ , s.t., if a query point  $q \in b(p, r_{\text{low}}(p))$  but  $q \notin b(p, r_{\text{loss}}(p))$ , then  $p$  is a  $(1 + \gamma/3)$ -approximate NN of the query point. Also, define  $r_{\text{high}}(p) = (36\lambda r_{\text{loss}}(p)n \log n/\gamma)$ . Note that for any point  $p$ ,  $r_{\text{high}}(p) > r_{\text{death}}(p)$ .

**Definition 3.4** ( $\text{parent}, \text{parent}_i$ ). In the tree that is created by the  $\lambda$ -stretch hierarchical clustering of  $P$ ,  $\text{parent}(p)$  is defined as the parent of node  $p$ . The  $i$ th parent of node  $p$  in the tree that is generated by the  $\lambda$ -stretch hierarchical clustering of  $P$  is defined as  $\text{parent}_i(p)$ , where  $\text{parent}_0(p)$  is  $p$  and  $\text{parent}_i(p)$  is  $\text{parent}(\text{parent}_{i-1}(p))$ . The function  $\text{parent}_i(p)$  is defined up to the point when it becomes the root of the tree that is generated by the  $\lambda$ -stretch hierarchical clustering of  $P$ , but it is not defined after the root.

We will use  $\text{Subtree}(p)$  to denote the subtree rooted at  $p$ .

### 2.3. Construction of balls (algorithm $\text{ConstructBalls}(P, \gamma)$ )

In this section  $\lambda = nd$ . Given a  $\lambda$ -stretch hierarchical clustering of  $P$ , for each point  $p$  in  $P$ , let  $r_0$  be the  $r_{\text{loss}}$  value for  $p$  and let  $p_1, p_2, \dots, p_m$  be the children of  $p$  in sorted order of their  $r_{\text{loss}}$  values, i.e.,  $r_{\text{loss}}(p_1) \leq r_{\text{loss}}(p_2) \leq \dots \leq r_{\text{loss}}(p_m)$ . Also, let  $x$  be the parent of  $p$  in the tree formed by the  $\lambda$ -stretch hierarchical clustering of  $P$ .

We construct the following ball sets around  $p$ :

- (1) Balls with radius  $r_i = r_{\text{loss}}(p)(1 + \gamma/3)^{j-1}$  for  $j = 0, \dots, M-1$ , where  $M = \lceil \log_{1+\gamma/3} (r_{\text{high}}(p)/r_{\text{loss}}(p)) \rceil$ .

This defines a ball set in the range  $[r_{\text{low}}(p), r_{\text{high}}(p)]$ .

- (2) Balls with radius  $r_i = r_{\text{loss}}(p_i)(1 + \gamma/3)^{j-1}$  for  $j = 0, \dots, M-1$ , where  $M = \lceil \log_{1+\gamma/3} (r_{\text{high}}(p_i)/r_{\text{loss}}(p_i)) \rceil$ .

This defines a ball set in the range  $[r_{\text{low}}(p_i), r_{\text{high}}(p_i)] \forall 1 \leq i \leq m$ .

We also construct a universal ball (of infinite radius) centered at the point that is the root of the tree formed by the  $\lambda$ -stretch hierarchical clustering of  $P$ , so that any point that is not in any of the balls constructed by the above rules lies in this ball.

### 2.4. Correctly answering $(1 + \epsilon)$ -approximate NN queries

In this subsection we prove that reporting the center of the smallest ball that contains a query point from the set of balls constructed by the algorithm  $\text{ConstructBalls}$  answers the approximate nearest neighbor query correctly.

**Observation 3.2.**  $\|a\|b \geq 2r_{\text{loss}}(p) \forall a \in \text{Subtree}(p) \text{ and } b \notin \text{Subtree}(p) \text{ as } 2r_{\text{loss}}(p) \text{ is the minimum separation between any point in the cluster from any point outside it by definition of } r_{\text{loss}}.$

Find the biggest ball such that for every query point beyond it we may claim another point exists which is a close neighbour of the query point with a limited cumulative approximation error; this will allow us to restrict the number of balls we build around a particular point.

If the query point is not within the biggest ball that can be created around the point  $x$ , then  $\text{parent}(x)$  is a neighbour of the query point with some accumulated approximation error, as shown in the following lemma, which also states that  $r_{\text{high}}(x)$  is a limit on the radius of the largest ball. Therefore, we may disregard the point  $x$ , but this will result in an approximation mistake that will add up over time.

The cumulative mistake is a factor of  $1/\gamma(3 \log n)$ , according to the  $r_{\text{high}}(x)$  formulation. As we shall see in a bit, this number is selected in such a way that the cumulative approximation mistakes can only grow up to a reasonable limit.

**Lemma 3.1.** For any query point  $q$ , if  $x$  is a  $(1 + \alpha)$ -NN of  $q$  in  $P$ , and if  $\|qx\| > r_{\text{high}}(x)$ , then  $\text{parent}(x)$  is a  $((1 + \gamma/(3 \log n))(1 + \alpha))$ -NN of  $q$  in  $P$ .

**Proof.** The proof of this lemma is similar to the proof of Lemma 2.13 in [11].

Let  $z$  be parent of  $x$  in the  $\lambda$ -stretch hierarchical clustering of the set of points  $P$ . Note that  $\|qx\| > r_{\text{high}}(x) > r_{\text{death}}(x)$ . By the same argument as in Lemma 2.10 of [11],  $\|zx\| \leq 2n\lambda r_{\text{loss}}(x) = (\gamma/(3 \log n))r_{\text{death}}(x) < (\gamma/(3 \log n))\|qx\|$ . Therefore,  $\|zq\| \leq \|qx\| + \|zx\| \leq (1 + \gamma/(3 \log n))\|qx\| \leq (1 + \gamma/(3 \log n))(1 + \alpha)d_P(q)$ .

It follows that  $\text{parent}(x)$  is a  $((1 + \gamma/(3 \log n))(1 + \alpha))$ -NN of  $q$  in  $P$ .  $\square$

If the query point is not in the biggest ball surrounding the current candidate, we may recursively consider the parent of the current candidate as the next closest neighbour candidate. This is shown in the following lemma. Any time we move the candidate closer to the parent, our guess becomes more inaccurate.

**Lemma 3.2.** Let  $p$  be the NN of a query point  $q$  in  $P$ . Let  $r$  be the radius of the smallest ball containing  $q$  and let it be centered at  $x$ . If  $r > r_{\text{high}}(\text{parent}^i(p)) \forall 0 \leq i \leq j - 1$ , then  $\text{parent}^j(p)$  is a  $((1 + \gamma/(3 \log n))^j)$ -approximate NN of  $q$  in  $P$ .

**Proof.** The proof is by induction on  $j$  using Lemma 3.1.  $\square$

Since we are making balls such that each one is larger by a factor of  $1/3$  compared to the one before it, we argue that reporting the nearest neighbour might result in an additional approximation error factor of  $1/3$  for the following lemma. What this means is that even though  $p$  is closer, we may end up reporting a different point  $x$  instead of  $p$  if  $q$  is in one of the balls in the ballset of  $p$  and  $p$  is the closest neighbour of  $q$ . This is due to the fact that the ball with  $q$  centred at  $x$  is smaller than the ball with  $q$  centred at  $p$ , even though  $p$  is closer. We can do this by using a discrete collection of balls, but there is an extra approximation error factor of  $1 + \gamma/3$ , which means that we will report  $x$  instead of  $p$ .

**Lemma 3.3.** For any query point  $q$ , if  $p$  is a  $(1 + \alpha)$ -NN of  $q$  in  $P$ , and if there exists a ballset, centered at  $p$ , in  $[r_{\text{low}}(t), r_{\text{high}}(t)]$  for some point  $t$ , and if the smallest ball containing  $q$  has radius  $r$ , such that  $r_{\text{loss}}(t) \leq r \leq r_{\text{high}}(t)$  and is centered at  $x$ , then  $x$  is a  $(1 + \gamma/3)(1 + \alpha)$ -approximate NN of  $q$  in  $P$ .

**Proof.** If  $x = p$ , then we are done.

Suppose  $x \neq p$ . Then two cases arise, either  $q \in B(p, r_{\text{high}}(t))$  or  $q \in B(p, r_{\text{high}}(t))$ . We will show in either case that  $d_P(q) \leq \|xq\| \leq (1 + \gamma/3)(1 + \alpha)d_P(q)$ .

**Case I.**  $q \in B(p, r_{\text{high}}(t))$ .

In this case,  $\|pq\| > r_{\text{high}}(t) \geq r > \|xq\|$ . This implies that  $x$  is closer to  $q$  than  $p$  and therefore trivially  $\|xq\| < (1 + \gamma/3)(1 + \alpha)d_P(q)$ .

**Case II.**  $q \in B(p, r_{\text{high}}(t))$ .

Clearly,  $q \notin B(p, r_{\text{low}}(t))$ , because otherwise the smallest ball containing  $q$  would have radius  $\leq r_{\text{low}}(t)$  and  $r \geq r_{\text{loss}}(t) > r_{\text{low}}(t)$  leading to a contradiction that the smallest ball containing  $q$  has radius  $r$ . Hence,  $\exists j$ , such that  $q \in B(p, r_j)$  and  $q \notin B(p, r_{j+1})$ , where  $r_0 = r_{\text{low}}(t)$ , implying that  $r_j < \|pq\| \leq r_{j+1} = (1 + \gamma/3)r_j$ .

Therefore,  $\|xq\| \leq r \leq r_{j+1} = (1 + \gamma/3)r_j < (1 + \gamma/3)\|pq\| \leq (1 + \gamma/3)(1 + \alpha)d_P(q)$ . This means that  $x$  is a near-NN of  $q$  in  $P$  that is  $(1 + \gamma/3)(1 + \alpha)$ -approximate.

Assume that the query point  $q$  is included in the ball with radius  $r_{\text{loss}}(p)$  for every given location  $p$ . Then, it's clear that the query point  $q$ 's nearest neighbours can only be points in the cluster formed by the subtree of the hierarchical clustering tree rooted at  $p$ . Even then, we have no idea which of these stances is nearest;  $p$  may be one of them. This can only be accomplished by creating additional balls around  $p$  with radii that are compatible with the ballsets.

comprised of nodes that are children of  $p$ . This method has the potential to reveal which of  $p$ 's progeny is geographically nearest to the point  $q$  that we're trying to identify. There is no difference between this and the set of balls generated by algorithm ConstructBalls in accordance with rule 2. It is worth noting that the child nodes of  $p$  would be more appropriate representations of the children farther down the hierarchy than the balls that may be built around  $p$ .

If we use method ConstructBalls to produce balls, the previous lemma shows that in that case, we can only answer the queries with a further compounded approximation inaccuracy of a factor of  $1 + \gamma/3$ .

**Lemma 3.4.** *For any query point  $q$ , if  $p$  is a  $(1 - \alpha)$ -NN of  $q$  in  $P$ , and we have balls constructed around the points of  $P$  as defined by ConstructBalls, and if the smallest ball containing  $q$  has radius  $r$ , such that  $r < r_{\text{loss}}(p)$  and is centered at  $x$ , then  $x$  is a  $(1 + \gamma/3)(1 + \alpha)$ -approximate NN of  $q$  in  $P$ .*

**Proof.** If  $q \in B(p, r_{\text{loss}}(p))$ , then  $\|pq\| > r_{\text{loss}}(p) > r \geq \|xq\|$  implying that  $x$  is closer to  $q$  than  $p$  and therefore trivially  $xq \in (1 + \gamma/3)(1 + \alpha)d_p(q)$ .  $\parallel \parallel$

Suppose  $q \notin B(p, r_{\text{loss}}(p))$ . Since  $q \in B(x, r)$  and  $q \notin B(p, r_{\text{loss}}(p))$ , these balls intersect. Then  $x \in \text{Subtree}(p)$  because all balls of radius  $< r_{\text{loss}}(p)$  centered around points  $\in \text{Subtree}(p)$  cannot intersect since  $\|a\|b\| \geq 2r_{\text{loss}}(p) \forall a \in \text{Subtree}(p)$  and  $b \in \text{Subtree}(p)$ . Let  $p_1$  be the child of  $p$ , such that  $x \in \text{Subtree}(p_1)$ . We consider three cases based on the value of  $r$ .

**Case I.**  $r < r_{\text{loss}}(p_1)$ .

We know that  $x \in \text{Subtree}(p_1)$  and  $p \notin \text{Subtree}(p_1)$ . Then  $\|pq\| > r$ , because otherwise  $\|px\| \leq \|pq\| + \|qx\| < r + r = 2r \leq 2r_{\text{loss}}(p_1)$ , and this would contradict Observation 3.2.

This implies that  $x$  is closer to  $q$  than  $p$  and therefore trivially  $\|xq\| \leq (1 + \gamma/3)(1 + \alpha)d_p(q)$ .

**Case II.**  $r > r_{\text{high}}(p_1)$ .

This case is not possible as we do not construct balls larger than  $r_{\text{high}}(t)$  around any point  $t$  of  $P$ , and  $r_{\text{high}}(t) \leq r_{\text{high}}(p_1) \forall t \in \text{Subtree}(p_1)$ , but we know that there is a ball centered at  $x \in \text{Subtree}(p_1)$  that contains  $q$ .

**Case III.**  $r_{\text{loss}}(p_1) \leq r \leq r_{\text{high}}(p_1)$ .

Since there is a ballset around  $p$  in the interval  $[r_{\text{low}}(p_1), r_{\text{high}}(p_1)]$ ,  $\exists j$ , such that  $q \in B(p, r_j)$  and  $q \notin B(p, r_{j+1})$ , where  $r_0 = r_{\text{low}}(p_1)$ . This implies that  $\|pq\| > r_j$ . It must be that  $r \leq r_{j+1}$ , because otherwise  $B(p, r_{j+1})$  would be a ball of radius smaller than  $r$  containing the query point  $q$ .

Therefore,  $\|xq\| \leq r \leq r_{j+1} = (1 + \gamma/3)r_j < (1 + \gamma/3)\|pq\| \leq (1 + \gamma/3)(1 + \alpha)d_p(q)$ .

Thus in all valid cases,  $x$  is a  $(1 + \gamma/3)(1 + \alpha)$ -approximate NN of  $q$  in  $P$ .  $\square$

In the following theorem, we combine the results of all the previous lemmas to show that the construction of balls in algorithm ConstructBalls does answer the approximate nearest-neighbor queries correctly.

**Theorem 3.1.** *For a point set  $P$  and an approximation factor  $\epsilon$ , Let the smallest ball containing  $q$ , in the balls formed by algorithm ConstructBalls( $P, \gamma$ ), where  $\gamma = \epsilon/2$ , be of radius  $r$ , centered at  $x$ , then  $x$  is a  $(1 + \epsilon)$ -approximate NN of  $q$  in  $P$ .*

**Proof.** Let  $p$  be the NN of  $q$  in  $P$ .

**Case I.**  $r < r_{\text{high}}(p)$ .

$- r_{\text{loss}}(p) \leq r \leq r_{\text{high}}(p)$ :

If  $r_{\text{loss}}(p) \leq r \leq r_{\text{high}}(p)$ , then by the construction in algorithm ConstructBalls, since there exists a ballset in  $[r_{\text{low}}(p), r_{\text{high}}(p)]$ , using Lemma 3.3,  $x$  is a  $(1 + \gamma/3)$ -NN of  $q$  in  $P$ .



–  $r < r_{\text{loss}}(p)$ :

If  $r < r_{\text{loss}}(p)$ , then using Lemma 3.4,  $x$  is a  $(1 + \gamma/3)$ -NN of  $q$  in  $P$ .

**Case II.**  $r > r_{\text{high}}(\text{parent}^i(p))$  and  $r < r_{\text{high}}(\text{parent}^{i+1}(p))$  for some  $i$ .

By Lemma 3.2,  $\text{parent}^{i+1}(p)$  is a  $((1 + \gamma/(3 \log n))^{i+1})$ -approximate NN of  $q$  in  $P$ .

- If  $r_{\text{loss}}(\text{parent}^{i+1}(p)) \leq r \leq r_{\text{high}}(\text{parent}^{i+1}(p))$ , then by the construction in algorithm ConstructBalls, since there exists a ballset in  $[r_{\text{low}}(\text{parent}^{i+1}(p)), r_{\text{high}}(\text{parent}^{i+1}(p))]$ , centered at  $\text{parent}^{i+1}(p)$ . Using Lemma 3.3,  $x$  is a  $((1 + \gamma/(3 \log n))^{i+1}(1 + \gamma/3))$ -approximate NN of  $q$  in  $P$ .
- If  $r < r_{\text{loss}}(\text{parent}^{i+1}(p))$ , then using Lemma 3.4,  $x$  is a  $((1 + \gamma/(3 \log n))^{i+1}(1 + \gamma/3))$ -NN of  $q$  in  $P$ .  
Also  $i + 1 \leq \log n$ , since height of the tree formed by the  $\lambda$ -stretch hierarchical clustering of  $P$  is bound by  $\log n$ .  
Therefore  $x$  is a  $((1 + \gamma/(3 \log n))^{\log n}(1 + \gamma/3))$ -approximate NN of  $q$  in  $P$ .

**Case III.**  $r > r_{\text{high}}(t)$ , where  $t$  is the root of the tree formed by the  $\lambda$ -stretch hierarchical clustering of  $P$ .

There is no ball in  $P$  of radius  $> r_{\text{high}}(p)$ . Thus  $t$  is reported as the approximate-NN of  $q$  in  $P$  as  $q$  lies in the universal ball centered at  $t$ .

By Lemma 3.2,  $t$  is  $((1 + \gamma/(3 \log n))^i)$ -approximate NN of  $q$  in  $P$ , where  $i \leq \log n$ , since height of the tree formed by the  $\lambda$ -stretch hierarchical clustering of  $P$  is bound by  $\log n$ .

Therefore  $x$  is a  $((1 + \gamma/(3 \log n))^{\log n}(1 + \gamma/3))$ -approximate NN of  $q$  in  $P$ .

And as  $(1 + \gamma/3)(1 + \gamma/(3 \log n))^{\log n} \leq 1 + 2\gamma$ ,  $x$  is a  $(1 + 2\gamma)$ -approximate NN of  $q$  in  $P$ .

Hence  $x$  is a  $(1 + \epsilon)$ -approximate NN of  $q$  in  $P$ .  $\square$

## 2.5. A bound on the total number of balls required

In the following theorem we analyze the number of balls that we construct in the algorithm ConstructBalls.

**Theorem 3.2.** The total number of balls constructed by the algorithm ConstructBalls( $P, \epsilon/2$ ) is  $O((1/\epsilon^2)n \log(n/\epsilon))$ .

**Proof.** The number of balls in a ball set is  $O((1/\epsilon) \log_{1-\epsilon}(n/\epsilon)) = O((1/\epsilon^2) \log(n/\epsilon))$ .

For each point, one ballset is constructed for the point by rule 1 of algorithm ConstructBalls. Additionally, one ballset is constructed for each child of the point by rule 2 of algorithm ConstructBalls. By charging the balls constructed around the child nodes (using rule 2) by algorithm ConstructBalls to the respective child nodes, the charge incurred at each point is at most that of 2 ball sets, i.e.,  $O((1/\epsilon^2) \log(n/\epsilon))$ .

Therefore, the total charge incurred over the  $n$  points is  $O((1/\epsilon^2)n \log(n/\epsilon))$ .  $\square$

This improves the bound on the number of balls constructed in [11] by a factor of  $\log n$ . These balls can be used to construct cells, i.e., quadtree boxes, and then store them in a BBD-tree as described in [11] resulting in an improvement of a factor of  $\log n$  in space complexity and preprocessing time while keeping the query time logarithmic.

## 3. A linear space solution based on $\epsilon$ -PLSB and a new hierarchical clustering

As we showed in the last section, the recursive nature of the data structure proposed by Har-Peled [11] may be eliminated, resulting in an  $O(n \log(n))$  space. We generate balls of radius between  $r_{\text{loss}}(p)$ ,  $r_{\text{loss}}(p) (1 - \epsilon)$ ,  $r_{\text{loss}}(p) (1 - \epsilon)^2, \dots, r_{\text{death}}(p)$  around each point in the database, which causes the additional  $\log(n)$  factor to emerge. Every point in the database contributes  $\Omega(\log(n))$  since  $r_{\text{death}}(p)$  is defined as  $\Omega(n/\epsilon)$  times  $r_{\text{loss}}(p)$ . We really do not need  $r_{\text{death}}(p)$  to be so high, however. It is necessary for  $r_{\text{death}}$  to be the distance beyond which no point in the subtree( $p$ ) is any closer to the query point  $q$  than any other point in the tree. Assuming that  $\text{diam}$  is the diameter of a collection of points, an acceptable distance for that may be  $O(\text{diam}(\text{subtree}(\text{parent}(p)))/\epsilon)$ .

Despite using this  $r_{\text{death}}$  definition, the number of balls that might be created could still be  $O(n \log(n))$ . Take the following example to illustrate the point: there is a line that connects  $n$  locations (let's call them  $p_1, p_2, \dots, p_n$ ) such that

the distance between points  $p_i$  and  $p_{i+1}$  is  $1 + ih$ , where  $h$  is

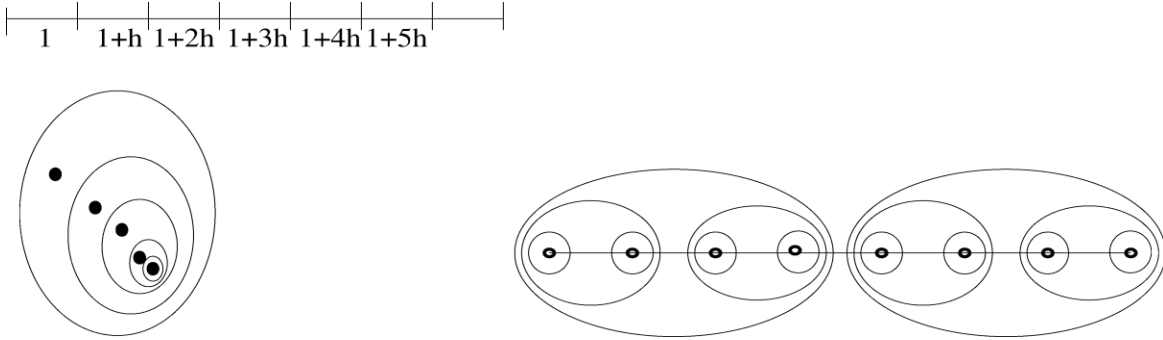


Fig. 2. Two different clustering for almost equidistant points in a line.

- 4 so little that it defies calculation. The MST is obviously only one chain. Even if we rethink  $r_{death}$  as proposed, the quantity of balls produced by the method in [11] (and its adaptation here) remains  $\Omega(n \log(n))$ . The issue is that while building the hierarchical clustering, the method in [11] doesn't take the edge-lengths of the approximation MST into account, other from their relative ordering. We rethink the way hierarchical clustering is built in order to address that issue. We also care about the absolute values of the MST's edge lengths, not only their relative ordering; in other words, we don't differentiate between edges with almost similar lengths. Hence, the hierarchical clustering in the first case may seem different, as shown in Figure 2. Obviously, the second one is better; we'll go back to this specific instance later.

#### 4.1. $\epsilon$ -NNS and $\epsilon$ -PLSB

The problem that we address here is a related problem called *Approximate Point Location in Smallest Ball* ( $\epsilon$ -PLSB) which is an approximate version of *Point Location in Smallest Ball* (PLSB) problem. In defining  $\epsilon$ -PLSB, we take care of the following.

- (1) The approximate version of PLSB can be solved by searching for cubes in a hierarchical grid as [11] (see Appendix C).
- (2) The approximate nearest neighbor searching problem is reducible to the  $\epsilon$ -PLSB problem.

For this, we make use of another point  $q^r$  which will ensure that the approximation factor  $1 + \epsilon$  is achieved from the nearest neighbor of  $p$ . A point  $p$  is an  $\epsilon$ -PLSB solution for a query point  $q$ , if there exists  $q^r$  such that (i)  $p$  is a PLSB solution for  $q^r$  and (ii)  $q^r$  is near  $q$ . The formal definitions follow.

##### Definition 4.1.

- (1) **PLSB.** Given a set of points  $P$ , and a finite set  $Q \subset P \times \mathbb{R}$ , we want to build a data structure  $D$ , such that given any query point  $q$ , we are able to return a pair  $(p, r) \in Q$ , such that  $b(p, r)$  is the smallest ball containing  $q$ . If there is no such ball, then it should return NIL.
- (2)  **$\epsilon$ -PLSB.** Given a set of points  $P$ , and a finite set  $Q \subset P \times \mathbb{R}$ , we want to build a data structure  $D$ , such that given any query point  $q$ :
  - (i) If  $q \in b(p, r)$  for some  $(p, r) \in Q$ , then we must return a pair  $(p^r, r^r) \in Q$ , such that  $r^r \leq r$  and there exists a point  $q^r \in b(p^r, r^r)$  such that  $\text{dist}(q, q^r) \leq r^r \epsilon$  and  $b(p^r, r^r)$  is the smallest ball containing  $q^r$  among balls from  $\{b(p, r) \mid (p, r) \in Q\}$ 's solution to PLSB query for  $q^r$ .
  - (ii) If  $q$  is not contained in any ball from the set  $\{b(p, r(1 + \epsilon)) \mid (p, r) \in Q\}$ , then we should return NIL.
  - (iii) Otherwise, we can return anything (even NIL).

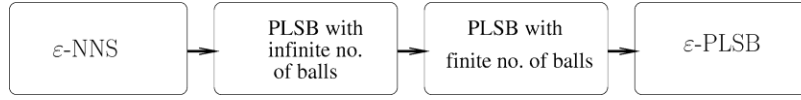


Fig. 3. The reduction of  $\varepsilon$ -NNS to  $\varepsilon$ -PLSB—the intermediate steps illustrate the proof strategy.

If we have a collection of points  $P$ , we can build a special set of balls  $B$  around those points (see to Definition 4.5) in such a way that the value of  $\varepsilon$ -PLSB in  $B$  gives us the approximate neighbour in  $P$ . To do this, we use the method shown in Figure 3.

A meticulously built hierarchical clustering of the point set  $P$  accomplishes the first reduction phase. Quite broadly, our hierarchical clustering is defined. The leaves of this binary tree are subsets of the input point set  $P$  that contain exactly one element each. As a subtree's leaves converge at a certain node, we say that node is internal. We build an instance  $I$  of  $\varepsilon$ -PLSB for any given tree in such a way that performing the query  $q$  for  $I$  can solve any nearest neighbour query for the point set  $P$ . To represent the minimum distance between two sets  $X$  and  $Y$ , we use the notation  $\text{dist}(X, Y)$ . Here we define hierarchical clustering using this notation.

**Definition 4.2 (Hierarchical clustering).** Given a point set  $P$ , its hierarchical clustering is a *binary* tree with singleton subsections of  $P$  serving as leaves, with each subsection of  $P$  appearing a unique number of times. The junction of the sub-tree leaves at each internal node indicates that node. Thus, every branch in the tree represents a cluster, which is a subset of  $P$ . Do not forget that cluster  $P$  is the root. Two numbers,  $r_{\min}(v)$  and  $r_{\max}(v)$ , are associated with the following attributes for each internal node  $v$  in the tree.

- (1)  $r_{\min}$  values increase from leaves to root.
- (2)  $r_{\min}(v) \leq \frac{1}{2} \text{dist}(\text{left}(v), \text{right}(v))$ .
- (3)  $r_{\max}(v) \geq 8 \frac{\text{diam}(v)}{\varepsilon}$ .

**Definition 4.3.** For every cluster  $v$  of  $T$ , we pick an arbitrary point  $p \in v$  and call it  $\text{leader}(v)$ .

**Definition 4.4 (BallRange( $p, r_1, r_2$ )).** It is defined to be the set of all balls centered at  $p$  of radius between  $r_1$  and  $r_2$ .  $\text{BallRange}(p, r_1, r_2) = \{b(p, r) \mid r_1 \leq r \leq r_2\}$ . Note that there are infinite number of balls in  $\text{BallRange}(p, r_1, r_2)$ .

**Definition 4.5 (BallSet( $v$ )).** For a cluster  $v$  of a hierarchical clustering  $T$ , we define  $\text{BallSet}(v)$  to be

- If  $v$  is leaf, then  $\text{BallSet}(v) = \varnothing$ .
- Otherwise,  $\text{BallSet}(v) = \text{BallSet}(\text{left}(v)) \cup \text{BallSet}(\text{right}(v)) \cup B_l \cup B_r$ , where  $B_l = \text{BallRange}(\text{leader}(\text{left}(v)), r_{\min}(v), r_{\max}(v))$ ,  $B_r = \text{BallRange}(\text{leader}(\text{right}(v)), r_{\min}(v), r_{\max}(v))$ .

Note that  $\text{BallSet}(v)$  is actually a union of  $2|v| - 2$  BallRanges.

**Definition 4.6 (Inner ball versus outer ball).** A ball is called an inner ball, if it is the smallest ball for one of the BallRanges of  $\text{BallSet}(v)$ . Otherwise, it is called an outer ball.

**Observation 4.1.** If  $b(p, r) \in \text{BallSet}(v)$  is an inner ball of  $\text{BallSet}(v)$ , then  $r = r_{\min}(u)$ , for  $u \in \text{Subtree}(v)$  in the hierarchical clustering of  $v$ .

**Observation 4.2.** Suppose, the ball  $b(p, r)$  is the smallest ball in  $\text{BallSet}(v)$  containing  $q$ . If  $b(p, r)$  is an outer ball, then  $\text{dist}(q, p) = r$ .

#### 4.2. Details of the reduction

We provide the reduction approach from  $\varepsilon$ -NNS to  $\varepsilon$ -PLSB for an lp metric on Rd, given a hierarchical clustering. Three stages might be seen as dividing this procedure. The first thing we do is make an endless supply of balls P1 such that

The answer to the PLSB inquiry q to P1 is a spherical shape, with its centre located around q's  $\varepsilon/7$ -nearest neighbour. Moving on to the next phase, we create a limited collection of balls P2 from P1 in such a way that the ball centred around the  $2\varepsilon/3$ -nearest neighbour of q is the result of the PLSB query q to P2. The outcome of the  $\varepsilon/7$ -PLSB query q to P2 is a ball that revolves around the  $\varepsilon$ -nearest neighbour of q in P, as we have shown in our proof.

#### 4.2.1. Generating $P_1$ (reduction to PLSB with infinite No. of balls)

Here we have to reduce  $\varepsilon/7$ -NNS to PLSB with infinite number of balls. For that, we prove that for any cluster  $v$  of  $T$ ,  $\text{BallSet}(v)$  is such a point set. The exact statement is the Theorem 4.1.

**Theorem 4.1.** For a cluster  $v$  of  $T$ , given a query point  $q$ , a point  $p \in v$  is  $\varepsilon/7$  nearest neighbor of  $q$  among points of  $v$ , if one of the following is true.

- If  $p$  is the center of the smallest ball in  $\text{BallSet}(v)$  containing  $q$ .
- There is no ball in  $\text{BallSet}(v)$  containing  $q$ .

Note that the second condition depends only on  $q$ , and so any arbitrary database point  $p$  is  $\varepsilon/7$ -nearest in that case.

**Proof.** We prove it by induction on the height of  $v$  in  $T$ . In the base case,  $v$  is a leaf, and  $\text{BallSet}(v)$  is empty. So, it is trivially proved.

Now, in the induction step, we assume the result for  $\text{left}(v)$  and  $\text{right}(v)$ , and prove it for  $v$ . Let  $p_1 = \text{leader}(\text{left}(v))$  and  $p_2 = \text{leader}(\text{right}(v))$ . By definition,  $\text{BallSet}(v) = (\text{BallSet}(\text{left}(v)) \cup B_l) \cup (\text{BallSet}(\text{right}(v)) \cup B_r)$ , where  $B_l$  and  $B_r$  are a set of balls around  $p_1$  and  $p_2$ , respectively. Note that  $\text{BallSet}(\text{left}(v)) \cup B_l = B_1$  (say) is a set of balls with centers among points of  $\text{left}(v)$  and  $\text{BallSet}(\text{right}(v)) \cup B_r = B_2$  (say) is a set of balls with centers among points of  $\text{right}(v)$ . We consider three cases (exhaustive but not necessarily disjoint).

**Case 1.** Since,  $b(p_1, r_{\max}(v)) \in B_1$ , so  $\text{dist}(q, p_1) \geq r_{\max}(v) \geq 8 \text{diam}(v)/\varepsilon$ . So, distance between  $q$  and  $p_1$  is large compared to the diameter of  $v$ . It is easy to see that for any point  $p \in v$ ,  $\frac{\text{dist}(q, p)}{\text{dist}(q, v)} \leq \frac{r_{\max}(v)}{r_{\max}(v) - \text{diam}(v)} \leq \frac{8}{8 - \varepsilon} \leq 1 + \varepsilon/7$ .

**Case 2.** There exists no ball in  $B_2$  containing  $q$ : Similar to Case 1.

**Case 3.** There exists a ball each in  $B_1$  and  $B_2$  containing  $q$ : Let, the smallest ball in  $B_1$  and  $B_2$  containing  $q$  be  $b(p_1^*, r_1)$  and  $b(p_2^*, r_2)$ , respectively. To handle this case, we use the Lemma 4.1 proved later.

**Lemma 4.1.**  $p_1^*$  is an  $\varepsilon/7$ -nearest neighbor of  $q$ , among points of  $\text{left}(v)$ , and  $p_2^*$  is an  $\varepsilon/7$ -nearest neighbor of  $q$ , among points of  $\text{right}(v)$ .

Now we consider the three sub-cases arising in Case 3. We have  $\text{dist}(q, p_1^*) \leq r_1$  and  $\text{dist}(q, p_2^*) \leq r_2$ .

**Case 3a.** The ball  $b(p_1^*, r_1)$  is an inner ball: In this case, from Observation 4.1, we have  $r_1 = r_{\min}(v')$  for some  $v' \subseteq v$ . Using this with the definition of  $r_{\min}$  we get:  $\text{dist}(q, p_1^*) \leq r_1 = r_{\min}(v') \leq r_{\min}(v) \leq 0.5 \text{dist}(\text{left}(v), \text{right}(v))$ . This implies that  $\text{dist}(q, \text{left}(v)) \leq 0.5 \text{dist}(\text{left}(v), \text{right}(v))$ . But,  $\text{dist}(q, \text{right}(v)) \geq \text{dist}(\text{left}(v), \text{right}(v)) - \text{dist}(q, \text{left}(v))$ . Therefore,  $\text{dist}(q, \text{right}(v)) \geq 0.5 \text{dist}(\text{left}(v), \text{right}(v)) \geq \text{dist}(q, \text{left}(v))$ . So,  $\text{dist}(q, v) = \text{dist}(q, \text{left}(v))$ . Now,  $\frac{\text{dist}(q, p_1^*)}{\text{dist}(q, v)} = \frac{\text{dist}(q, p_1^*)}{\text{dist}(q, \text{left}(v))} \leq 1 + \varepsilon/7$ . So,  $p_1^*$  is an  $\varepsilon/7$  nearest neighbor of  $q$  in  $v$ .

Also, any ball in  $B_2$  containing  $q$  must have a radius larger than  $\text{dist}(q, \text{right}(v))$ . But,  $\text{dist}(q, \text{right}(v)) \geq 0.5 \text{dist}(\text{left}(v), \text{right}(v)) \geq r_{\min}(v') = r_1$ . So, any ball in  $B_2$  containing  $q$  is larger than  $b(p_1^*, r_1)$ . So, we see that the smallest ball containing  $q$  in  $\text{BallSet}(v)$  is centered around  $\varepsilon/7$  nearest neighbor of  $q$  in  $v$ .

**Case 3b.** The ball  $b(p_2^*, r_2)$  is an inner ball: Similar to Case 3a.

**Case 3c.** Both the balls  $b(p_1^*, r_1)$  and  $b(p_2^*, r_2)$  are outer balls: From Observation 4.2, we have  $r_1 = \text{dist}(q, p_1^*)$  and  $r_2 = \text{dist}(q, p_2^*)$ . Without loss of generality,  $r_1 \leq r_2$ . So,  $b(p_1^*, r_1)$  is the smallest ball containing  $q$ , in  $B_1 \cup B_2 = \text{BallSet}(v)$ . Also,  $\frac{\text{dist}(q, p_1^*)}{\text{dist}(q, \text{left}(v))} \leq 1 + \epsilon/7$  and  $\frac{\text{dist}(q, p_2^*)}{\text{dist}(q, \text{right}(v))} \leq 1 + \epsilon/7$  from the Lemma 4.1. Since,  $p_1^*$  is nearer to  $q$  than  $p_2^*$ , therefore  $\frac{\text{dist}(q, \text{left}(v))}{\text{dist}(q, p_1^*)} \leq 1 + \epsilon/7$ . Therefore  $\frac{\text{dist}(q, p_1^*)}{\min(\text{dist}(q, \text{left}(v)), \text{dist}(q, \text{right}(v)))} \leq 1 + \epsilon/7$ . This implies that  $p_1^*$  is the  $\epsilon/7$  nearest neighbor of  $q$  in  $v$ .  $\square$

Now we give the proof of Lemma 4.1.

**Proof.** Lemma 4.1 will only be proven for  $B_1$  here. The same holds true for  $B_2$ . We may express it as  $B_1 \cup \text{BallSet}(\text{left}(v))$ . Assuming that no ball in  $\text{BallSet}(\text{left}(v))$  includes  $q$ , all the points in  $\text{left}(v)$  are  $\epsilon/7$ -nearest neighbours of  $q$  according to the induction hypothesis. The smallest ball in  $\text{BallSet}(\text{left}(v))$  that contains  $q$  is centred around an  $\epsilon/7$ -nearest neighbour of  $q$  in  $\text{left}(v)$  if  $q$  is included in a ball of  $\text{BallSet}(\text{left}(v))$ . Assume that the smallest ball in  $\text{BallSet}(\text{left}(v))$  that contains  $q$  is  $b(p, r)$ . Our task is complete if this is the tiniest ball in  $B_1$ . Forget about it. From  $B_1$  comes the tiniest ball in  $B_1$ . Take  $B_1$  (and by extension,  $B_1$ ) as a whole, and consider the smallest ball  $b(p_1, r_1)$ .

**Case (i).**  $b(p^r, r^r)$  is an outer ball: In this case,  $\text{dist}(q, p^r) = r^r$ . So, if the smaller ball is  $b(p_1, r_1)$ , then  $\text{dist}(q, p_1) \leq r^r \leq r^r = \text{dist}(q, p^r)$ . So, we choose a ball nearer than the  $\epsilon/7$ -nearest neighbor. Clearly, that is also  $\epsilon/7$ -nearest neighbor.

**Case (ii).**  $b(p^r, r^r)$  is an inner ball: In this case,  $r^r$  is  $r_{\min}(v')$  for some  $v' \subset v$ . The smallest ball in  $B_1$  is of radius  $r_{\min}(v)$ . So  $r^r \geq r_{\min}(v) \geq r_{\min}(v') = r^r$ . So, the  $b(p, r)$  is a smaller ball, which is a contradiction. So, this proves Lemma 4.1.  $\square$

#### 4.2.2. Generating $P_2$ from $P_1$ (reduction to PLSB with finite No. of balls)

In this section, we show that we need to retain only a finite number of balls of  $P_1$ . For that we consider each ball range separately.

**Theorem 4.2.** Given a hierarchical clustering  $T = (V, E)$  for a point set  $P$ , we can construct a set of

$$2 \max_{v \in V} \frac{1}{1 + \epsilon} \log \frac{r_{\max}(v)}{r_{\min}(v)}$$

balls such that following is true for a PLSB constructed for them. Given a query point  $q$ , the point returned by PLSB data structure is  $1 + 2\epsilon/5$  nearest neighbor of  $q$  among points of  $P$ .

**Proof.** We start with  $\text{BallSet}(P)$ . Recall that there are infinite number of balls. We shall retain a carefully selected finite set from these. Recall that balls in the  $\text{BallSet}(P)$  can be expressed as union of  $2|P| - 2$  BallRanges. We replace each maximal  $\text{BallRange}(p, r_1, r_2) \subset \text{BallSet}(P)$  with

$$b(p, r) : r = r_1 + \frac{\epsilon}{7}^k \wedge r \leq r_2(1 + \epsilon/7) \wedge k \in \{0, 1, 2, \dots\}.$$

Let  $B$  be this finite set of balls. We will show that, a PLSB constructed for  $B$  returns a ball which is centered around a  $(1 + \epsilon/7)(1 + \epsilon/7) (\leq 1 + 2\epsilon/5)$  nearest neighbor. This completes the proof of the theorem.

Suppose,  $q$  is a query point, for  $(\epsilon/7)$ -PLSB of  $B$ . Suppose,  $b(p, r)$  is the smallest ball containing  $q$  among balls of  $\text{BallSet}(P)$ .

**Case 1** ( $b(p, r)$  is an inner ball of  $\text{BallSet}(P)$ ). From Theorem 4.1, we know, that  $p$  is an  $\epsilon/7$ -nearest neighbor of  $q$ . Also, we know that  $b(p, r)$  is contained in  $B$  ( $B$  contains all the inner balls of  $\text{BallSet}(P)$ ). So,  $b(p, r)$  is the smallest ball in  $B$  containing  $q$ , and  $p$  is an  $\epsilon/7$ -nearest neighbor of  $q$ . Hence proved.



**Case 2** ( $b(p, r)$  is an outer ball of  $\text{BallSet}(P)$ ). In this case, we observe that there is a ball  $b(p, r^*) \in B$ , with  $r^*$  at most  $r(1 + \epsilon/7)$  which contains  $q$ . So, the smallest ball in  $B$  containing  $q$  is of radius utmost  $r(1 + \epsilon/7)$ . So, the center of the smallest ball in  $B$  containing  $q$  is at a distance of at most  $r(1 + \epsilon/7)$  from  $q$ , which is  $(1 + \epsilon/7)^2 \text{dist}(q, P)$ .  $\square$

#### 4.2.3. Reduction to $\epsilon/7$ -PLSB with finite No. of balls

**Theorem 4.3.** Given a hierarchical clustering  $T = (V, E)$  for a point set  $P$ , we can construct a set of

$$2 \max_{v \in V} 1, \log_{1+\epsilon} \frac{r_{\max}(v)}{r_{\min}(v)}$$

balls such that following is true for an  $\epsilon/7$ -PLSB constructed for it. Given a query point  $q$ , the point returned by  $\epsilon/7$ -PLSB data structure is  $1 + \epsilon$  nearest neighbor of  $q$  among points of  $P$ .

**Proof.** We use the same finite set of balls  $B$  as used in Theorem 4.2, and show that an  $\epsilon/7$ -PLSB made for it satisfies the above condition. Suppose,  $b(p, r)$  is the ball returned by the  $\epsilon/7$ -PLSB data structure for  $B$ .

**Case 1** ( $q \in b(p, r)$ ). From definition of  $\epsilon/7$ -PLSB, we get that  $b(p, r)$  is the smallest ball in  $B$  containing  $q$ . Using Theorem 4.2, we get that  $p$  is  $(1 + \epsilon/7)^2$  nearest neighbor of  $q$ .

**Case 2** ( $q \notin b(p, r)$ ). In this case, we observe from the definition of  $\epsilon/7$ -PLSB, that there exists a query point  $q^*$ , such that  $\text{dist}(q, q^*) \leq \epsilon r/7$  and  $b(p, r)$  is the smallest ball in  $B$  containing  $q^*$ . Using this with Theorem 4.2,  $\text{dist}(q, P) + \text{dist}(q, q^*) \geq \text{dist}(q^*, p)/(1 + \epsilon/7)^2$ . This implies  $\text{dist}(q, P) \geq \text{dist}(q^*, p)/(1 + \epsilon/7)^2 - \epsilon r/7$ . Simplifying, we get  $\text{dist}(q, P) \geq \text{dist}(q^*, p)/(1 + \epsilon)$ . So,  $p$  is a  $\epsilon$  nearest neighbor of  $q$ .  $\square$

$B$  is composed of ball ranges, two for each cluster. Number of balls in ballranges corresponding to cluster  $v$  is  $2 \max(1, \log_{1+\epsilon} \frac{r_{\max}(v)}{r_{\min}(v)})$ . So, we need at most  $2 \max_{v \in V} (1, \log_{1+\epsilon} \frac{r_{\max}(v)}{r_{\min}(v)})$  balls.

#### 4.3. Some observations on MST based hierarchical clustering

Given a point set  $P$ , we want to construct a hierarchical clustering such that

$$2 \max_{v \in V} 1, \log_{1+\epsilon} \frac{r_{\max}(v)}{r_{\min}(v)} = O(|P| \log_{1+\epsilon} (1/\epsilon)).$$

The procedure to construct such a hierarchical clustering is quite involved. We shall first describe a simpler hierarchical clustering which achieves  $O(|P| \log_{1+\epsilon} (|P|/\epsilon))$  bound. Then we describe methods for improving it to obtain  $O(|P| \log_{1+\epsilon} (1/\epsilon))$  bound.

The first hierarchical clustering that we describe is directly based on the MST of the point set. Suppose, the MST of the point set  $P$  is  $M = (P, E_M)$ . Then we construct the hierarchical clustering as follows: The longest edge in  $E_M$  divides  $P$  into two parts, say  $P_1$  and  $P_2$ . We find the hierarchical clusterings of  $P_1$  and  $P_2$  recursively, say  $C_1$  and  $C_2$ . The MST hierarchical clustering of  $P$ , then, has  $P$  as root and  $C_1$  and  $C_2$  as two children of the root.

Suppose,  $v$  is a cluster in the MST hierarchical clustering. We use

$$\begin{aligned} r_{\min}(v) &= \frac{1}{2} \text{dist}(\text{left}(v), \text{right}(v)) \quad \text{and} \\ r_{\max}(v) &= 8 \text{diam}(v)/\epsilon. \end{aligned}$$

Consider any two points  $p_1, p_2 \in v$ . Consider the path connecting  $p_1$  and  $p_2$  in the MST. Suppose,  $e_1, e_2, \dots, e_l$  are the edges in the path. Then by repeated application of triangle inequality, we get  $\text{dist}(p_1, p_2) \leq \text{len}(e_1) + \text{len}(e_2) + \dots + \text{len}(e_l)$ . Since,  $l$  is at most  $n$  and  $\text{len}(e_i) \leq \text{dist}(\text{left}(v), \text{right}(v))$  (because  $\text{dist}(\text{left}(v), \text{right}(v))$  is the length of the largest edge in MST over  $v$ ). Therefore, for any  $p_1, p_2 \in v$  we have  $\text{dist}(p_1, p_2) \leq n \text{dist}(\text{left}(v), \text{right}(v))$ . So, diameter of  $v$  is at most  $n \text{dist}(\text{left}(v), \text{right}(v))$ . Thus,

$$r_{\max}(v) \leq 8n \text{dist}(\text{left}(v), \text{right}(v)) / \epsilon.$$

So, 
$$\frac{r_{\max}(v)}{r_{\min}(v)} = O(n/\epsilon).$$

Thus 
$$2 \max_{v \in V} \log \frac{r_{\max}(v)}{r_{\min}(v)} = O(P \log \frac{P}{\epsilon}).$$

Here are some issues with MST hierarchical clustering:

- (1) Building the MST is a part of the preprocessing that takes almost quadratic time.
- (2) Finding the diameters of point sets is computationally intensive and is part of the preprocessing.
- (3) Instead of being linear, the space need is  $n \log n$ .

Approximate MSTs are used to address the initial issue. We use different interpretations of  $r_{\max}(v)$  to deal with the second issue. Before we provide our answer in the following part, let's examine some cases pertaining to the third difficulty.

$\Omega(P \log(P))$  is the worst-case space requirement for the current MST hierarchical clustering. A collection of  $n$  points  $p_1, p_2, \dots, p_n$  on a straight line with  $\text{dist}(p_i, p_{i+1}) = ih$  is one example. In this case,  $h$  is a very little number. A line graph (Fig. 2) represents the MST in this instance, and

$$\log \frac{r_{\max}(p)}{r_{\min}(p)} = O(\log 1 + \log 2 + \dots + \log n) = \Omega(n \log n).$$

*For the case where the edge lengths are growing geometrically, say,  $2^i$ , for  $i = 1, 2, \dots$ , the MST hierarchical clustering is actually linear.*

#### 4.4. An alternate hierarchical clustering

We provide a one-dimensional solution for points in a line case in linear space for the sake of exposition. Ignoring the exact ordering, let's use the prior example based on the fact that all MST edges are almost equal. Building

on this concept, we create the completely balanced tree-shaped hierarchical clustering shown in Figure 2 (the second clustering). For this instance,

$$\log \frac{r_{\max}(v)}{r_{\min}(v)} = O\left(\frac{n}{2} \log 2 + \frac{n}{4} \log 4 + \dots + \log n\right) = O(n).$$

Partitioning the edges of this MST into sets where all the edges inside a set have lengths within a factor of two of each other is the essence of our method. All edges within a single partition are treated as if they were of the same length when building the hierarchical clustering from the approx MST. Our standard MST hierarchical clustering procedure involves sorting the edges in ascending order of length and continuously merging clusters linked by the edges. The sequence of merging will be rearranged here. A lot more work goes into this reordering for the general  $d$ -dimensional clustering. A more straightforward approach is shown, nevertheless, for a single dimension.

Let  $T = (P, E_M)$  be the MST for  $P$ . We partition  $E_M$  into  $\{E_k \mid k \in \mathbb{Z}\}$ , where  $E_k = \{e \in E_M \mid 2^k \leq \text{len}(e) < 2^{k+1}\}$ . Suppose  $E_{k_1}, E_{k_2}, \dots, E_{k_l}$  are the set of all non-empty  $E_k$ 's with  $k_{i+1} < k_i$  (Fig. 4). We consider all edges in each  $E_{k_i}$  as being of effectively the same length and try to keep the clustering more balanced. Suppose we need to combine the clusters  $c_0, c_1, \dots, c_{k-1}$  into one because of the edges in  $E_{k_i}$ . This is because, as a line graph, the MST naturally forms contiguous intervals, and these clusters would fit exactly into that description. We will assume, without limiting

ourselves, that the sequence  $c_0, c_1, c_2, \dots, c_{k-1}$  is the same as it appears on the actual line. Next, we arrange the nodes  $c_0, c_1, \dots, c_{k-1}$  in a height-balanced binary tree. This tree suggests a natural sequence for the merging, therefore that's how it happens. Any cluster  $v$  that is generated in phase  $i$  will have an  $r_{\min}(v)$  equal to  $2^{k_i} - 1$  for this clustering.

*Proof for points in a line case*

The result that the generated clustering is linear follows from following two observations.

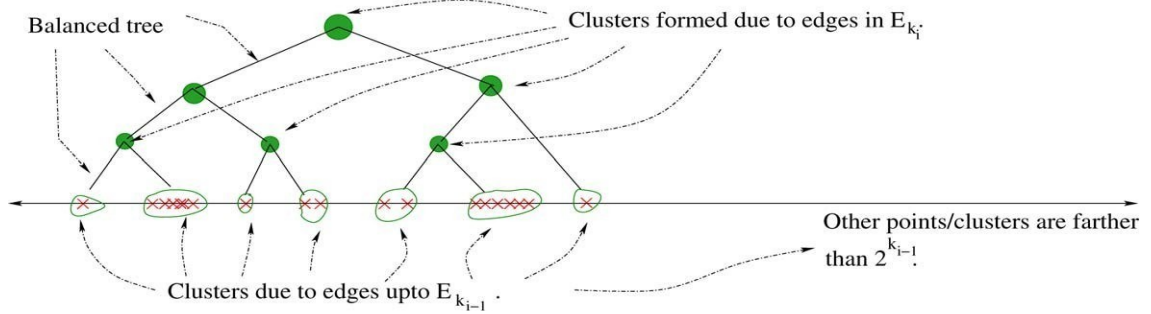


Fig. 4. Linear Hierarchical clustering for points in a line.

- (1) If  $W_i$  is the set of clusters that we add when we consider edges in  $E_{k_i}$  (so-called phase  $i$ ), then

$$\log_{\substack{v \in W_i \\ r_{\min}(v)}} \frac{r_{\max}(v)}{r_{\min}(v)} \quad O \quad = \quad \frac{\sum_{e \in E_{k_1} \cup E_{k_2} \cup \dots \cup E_{k_i}} \text{len}(e)}{2^{k_i}}.$$

- (2) This result requires some technical details that we provide for the general  $d$ -dimensional case. The size of hierarchical clustering is linear. This follows by observing that in the previous result, if we add the contribution of each edge  $e \in E_{k_i}$  to the total sum across all phases, we obtain  $O(\text{len}(e) \cdot \sum_{j=1}^i \frac{1}{2^{k_j}}) = O(1)$  as

$\text{len}(e) = O(2^{k_i})$ . Since there are  $O(|P|)$  edges in the MST, the total sum is  $O(|P|)$ .

#### 4.5. Constructing the general hierarchical clustering

As in the [11], we use a  $\lambda$ -MST instead of an exact MST.

The central idea in our solution is to partition the edges of a  $\lambda$ -MST into sets, such that all the edges within a set have lengths within a factor of two of each other (actually any constant factor would work). While constructing the hierarchical clustering from the approximate MST, all the edges within a single set are considered as if of same length.

Let  $T = (P, E_M)$  be a  $\lambda$ -MST for  $P$ . We partition  $E_M$  into  $\{E_k \mid k \in \mathbb{Z}\}$ , where  $E_k = \{e \in E_M \mid 2^k \leq \text{len}(e) < 2^{k+1}\}$ . Suppose  $E_1, E_2, \dots, E_l$  be the set of all non-empty  $E_k$ 's with  $k < l$ .

Recall that a hierarchical clustering is a labeled tree with the vertex set as a subset of power set of  $P$ . We construct the hierarchical clustering bottom up in  $l$  phases ( $l$  is the number of non-empty  $E_i$ 's). In any phase  $i$ , we start with a forest  $F_i$ . We initialize with the singletons i.e.  $F_1 = (\{p \mid p \in P\}, \emptyset)$ . Then iteratively merge the trees of  $F$  by creating new clusters in the forest, until we finally get a hierarchical clustering for  $P$ . By merging of two trees with roots  $v_1$  and  $v_2$ , we mean adding a new cluster  $v_1 v_2$ , with  $v_1$  and  $v_2$  as its two children. For merging more than two trees, we mean to keep on doing the binary merge until they become part of a single tree.  $F_i$  is the forest at the start of phase  $i$ . So,  $F_1 = (\{p \mid p \in P\}, \emptyset)$ . We maintain the invariant that,  $C$  is a root of a tree in  $F_i$  iff  $C$  is a maximal connected component of  $(P, E_{k_1} \cup E_{k_2} \cup \dots \cup E_{k_{i-1}})$ .

For any cluster  $v$  added in phase  $i$ , we assign

$$r_{\min}(v) = \frac{2^{k_i}}{2\lambda}.$$

#### Observation 4.3.

- (1)  $r_{\min}$  values increase from leaves to root. This is because, vertices added in later phases (which have larger  $r_{\min}$  values) are always "above" the vertices in earlier phases, in the hierarchical clustering.
- (2)  $r_{\min}(v) \leq \frac{1}{2} \text{dist}(\text{left}(v), \text{right}(v))$ . The proof is as follows. First observe that, the  $\lambda$ -MST edge going across any two maximal connected components of  $(P, E_{k_1} \cup E_{k_2} \cup \dots \cup E_{k_{i-1}})$  is of length at least  $2^{k_i}$ . So, the distance

<sup>1</sup> A forest is a graph, which can be expressed as a union of vertex-disjoint trees.

between those components is at least  $\frac{2^{k_i}}{\lambda}$ . Since, children of any cluster added in phase  $i$  is also a union of maximal connected components of  $(P, E_{k_1} \cup E_{k_2} \cup \dots \cup E_{k_{i-1}})$ , the distance between them is also lower bounded by  $\frac{2^{k_i}}{\lambda}$ .

**Definition 4.7** ( $\text{Path}_T(p_1, p_2)$ ). Given a  $\lambda$ -MST  $T$ , we define  $\text{Path}_T(p_1, p_2)$  for two points  $p_1$  and  $p_2$  as the set of edges on the unique path on  $T$  connecting  $p_1$  and  $p_2$ .

**Definition 4.8** ( $S(v)$ ). Given a  $\lambda$ -MST  $T$ , for a  $v \subseteq P$ , we define

$$S(v) = \bigcup_{p_1, p_2 \in v} \text{Path}_T(p_1, p_2).$$

Note that if  $v$  spans<sup>2</sup> a connected component in  $T$ , then  $S(v)$  is simply the set of edges in  $T$ , which go across vertices of  $v$ .

**Definition 4.9** ( $S(v)$ ). We define  $S(v)$  to be sum of lengths of all the edges in  $S(v)$ .

In our hierarchical clustering, we shall use

$$r_{\max}(v) = 8 \frac{S(v)}{\epsilon}.$$

**Observation 4.4.**

$$\text{diam}(v) \leq S(v).$$

The proof is a simple application of triangle inequality along the path joining the two most distant points in  $v$ .

From Observations 4.4 and 4.3, we conclude that the  $r_{\min}$  and  $r_{\max}$  values are well defined (consistent with the definition of the hierarchical clustering).

In what follows, we let  $\alpha = \log_{1+\epsilon} \left( \frac{\lambda}{\epsilon} \right)$ .

**Theorem 4.4.** Given a point set  $P$  and a  $\lambda$ -MST (say  $T$ ) of  $P$ , there exists a hierarchical clustering of  $P$ , with set of clusters  $V$  such that  $2^{-\max_{v \in V} \max(1, \log_{1+\epsilon} \frac{r_{\max}(v)}{r_{\min}(v)})} = O(\alpha |P|)$ . Further, this clustering can be done in  $O(n \log(n))$  time.

**Proof.** The proof of this theorem is based on the following lemma.

**Lemma 4.2.** If  $W_i$  is the set of clusters that we add in phase  $i$ , then

$$\log_{1+\epsilon} \frac{r_{\max}(v)}{r_{\min}(v)} = O \left( \alpha \times \frac{\text{dist}(p_1, p_2)}{2^{k_i}} \right)_{(p_1, p_2) \in E_{k_1} \cup E_{k_2} \cup \dots \cup E_{k_i}}$$

and the time taken in phase  $i$  is  $O(|W_i| \log(n))$ .

Using Lemma 4.2, we count the contribution of each edge of the  $\lambda$ -MST. Suppose,  $e \in E_{k_i}$ . Then, it does not contribute anything up to phase  $i-1$ . In any phase  $j \geq i$ , it contributes  $\alpha \text{len}(e)/2^{k_j}$ . So, total contribution is at most  $\alpha \text{len}(e) \sum_{j=i}^{\infty} (1/2^{k_j}) \leq 2\alpha \text{len}(e)/2^{k_i}$ . Also, since  $e \in E_{k_i}$ , so  $\text{len}(e) \leq 2^{k_i+1}$ . Therefore, the contribution of each edge is at most  $4\alpha$  and the total is  $O(\alpha(\text{number of edges in } \lambda\text{-MST}))$ , which is clearly  $O(\alpha |P|)$ .

As for the total time taken, time taken in phase  $i$  is upper bounded by  $|W_i| \log(n)$ . Since, we add a total of  $|W_i|$  clusters, so total time taken is  $O(n \log(n))$ .

This completes the proof of Theorem 4.4.  $\square$

<sup>2</sup> Spans as in a spanning tree. A vertex set  $S$  spans a connected component in a graph  $G = (V, E)$  if the graph  $G^S = (S, E \cap \binom{S}{2})$  is a connected graph.

#### 4.6. Merging within phase $i$ —Lemma 4.2

In this section we prove Lemma 4.2, using a series of intermediate results. We start by giving a graph theoretic result.

**Lemma 4.3.** *Given a tree  $T$  and  $k$  of its vertices  $\{v_1, v_2, \dots, v_k\}$  (possibly repeated), we can pair them up (if odd, 1 vertex would be left out) as  $\{(v_{i1}, v_{j1}), \dots, (v_{ik/2}, v_{jk/2})\}$ , so that the set of paths from  $v_{i1}$  to  $v_{j1}$  are all edge disjoint in  $T$ .*

*Furthermore, the tree  $T$  can be preprocessed in time  $O(n \log n)$ , such that once that is done, the pairing up of any given  $k$  vertices can be done in  $O(k \log n)$  time.*

**Proof.** We prove it by induction on the number of nodes in the tree. The base case is simple. In induction step, pick any leaf, say  $v$  of the tree  $T$ .

**Case 1** ( $v$  is not among the vertices to be paired). Remove it from the tree and use the induction hypothesis to pair all the vertices.

**Case 2** ( $v$  is among the vertices to be paired). Then  $v$  must be connected to some other vertex in  $T$ . Let it be  $u$ .

After deleting  $v$  and reversing the selection of  $u$  in  $T$ , couple  $u$  and  $v$ , and then pair the remainder of the vertices using the induction hypothesis. This process may be repeated if  $u$  is also to be paired.

To avoid pairing  $u$  with  $v$ , delete  $v$  and choose  $u$  using the induction hypothesis. The next step is to pair each vertex with an adjacent vertex in  $T$ . Let us pretend that  $u$  and  $w$  are coupled. After that, keep all of the pairings except for  $v$  and  $w$ . Every one of the pathways is clearly edge-joint.

We will create an algorithm that pairs  $k$  vertices in  $O(n)$  time using the aforementioned technique. With a preprocessing time of  $O(n \log(n))$ , we want an algorithm that can pair the inputs in  $O(k \log(n))$  time. You can figure out the features from this sketch of the building. First things first: build an approximate MST balanced edge separator tree. The temporal complexity for this task is  $O(n \log n)$ . Then, we divide the  $k$  input points in half along the separator tree sides in order to solve a pairing issue for  $k$  points. The next step is to use the corresponding subtree to recursively couple the two components. Lastly, if there are any unpaired vertices in either side, they are joined up with each other to provide a partial solution. The obtained query time is  $O(k \log(n))$ .  $\square$

We give some additional definitions.

#### Definition 4.10.

- (1)  $G_i$  is defined to be a sub-graph of the  $\lambda$ -MST of the point set  $P$  with vertex set  $P$  and edge set  $E_{k_1} \cup E_{k_2} \cup \dots \cup E_{k_{j-1}}$ . Note that it is a forest.
- (2)  $C_i$  is defined to be the set of connected components of  $G_i$ .
- (3)  $H_i$  is a forest over the connected components of  $G_i$  formed by edges in  $E_{k_i}$ . Its vertex set is  $C_i$ . The edge set is  $\{(c_s, c_t) \mid \exists (p_s, p_t) \in E_{k_i}, p_s \in c_s, p_t \in c_t\}$ .

#### Observation 4.5.

$C_i$  = The set of all the roots of the rooted forest  $F_i$ .

This follows from the discussion in Section 4.5.

We shall use the Lemma 4.3 to merge the clusters in the form of what we call a *star*. The definition of a star follows.

**Definition 4.11** (*Star*). A set of clusters  $c_0, c_1, \dots, c_k$  is called a star of the phase  $i$  with  $c_0$  as center, if they satisfy the following conditions.



- $c_0$  spans a connected component in  $G_{i+1}$ .
- For each  $j \in \{1, 2, \dots, k\}$ ,  $c_0 \cup c_j$  spans a connected component in  $G_{i+1}$ .

$c_0$  is called the center of the star and  $c_1, c_2, \dots, c_k$  are called the leaves of the star.

#### 4.6.1. Merging of a star

For simplicity, we assume that number of leaves of the star is a power of two. First we make the following claim.

**Lemma 4.4.** *Given a star with center  $c_0$  and leaves  $c_1, c_2, \dots, c_k$  we can pair each of  $c_1, c_2, \dots, c_k$  into  $k/2$  pairs, say,  $(c_1, c_2), \dots, (c_{k-1}, c_k)$ , such that*

$$S(c_{2j-1} \cup c_{2j}) \leq S(c_j) + k2^{ki+1}.$$

$j=1 \qquad \qquad \qquad j=0$

Time taken for the pairing is  $O(k \log(n))$ .

**Proof.** From the definition of the star,  $c_0 \cup c_j$  spans a connected component in  $G_{i+1}$ . So, there exists an edge in  $E_{G_{i+1}}^{c_0 \cup c_j}$  which connects  $c_0$  to  $c_j$ . Suppose, for each  $c_i$ ,  $(v_i, v_i')$  be the edge connecting  $c_0$  to  $c_i$  with  $v_i \in c_0$  and  $v_i' \in c_i$ .

#### Observation 4.6.

$$S(c_m \cup c_n) \leq S(c_m) + S(c_n) + \text{PathLength}(v_m', v_n'),$$

where  $\text{PathLength}(v_m, v_n)$  is the length of the path between  $v_m$  and  $v_n$  in the  $\lambda$ -MST.

#### Observation 4.7.

$$\text{PathLength}(v_m', v_n') \leq \text{PathLength}(v_m, v_n) + 2 \times 2^{ki}.$$

This follows from the fact that both the edges  $(v_m', v_m)$  and  $(v_n', v_n)$ , which appear at the two ends of the path from  $v_m'$  to  $v_n'$ , are of length at most  $2^{ki}$ .

Since,  $c_0$  spans a connected component in  $G_{i+1}$ , which is a sub-graph of the  $\lambda$ -MST, so  $c_0$  spans a tree in  $G_{i+1}$ . Each of  $v_i$  is contained in  $c_0$ . Use Lemma 4.3 to pair each  $v_i$ , using the tree spanned by  $c_0$  in  $G_{i+1}$ . Without losing generality, we assume that the pairing is  $(v_1, v_2), (v_3, v_4), \dots, (v_{k-1}, v_k)$ . Then, since the paths are edge-disjoint, we get the following result:

$$\text{PathLength}(v_{2j-1}, v_{2j}) \leq S(c_0).$$

$j=1$

Using Observations 4.6 and 4.7, we obtain

$$S(c_{2j-1} \cup c_{2j}) \leq S(c_j) + S(c_0) + k2^{ki+1}.$$

$j=1 \qquad \qquad \qquad j=1$

Also, from Lemma 4.3, time taken in pairing  $v_i$ 's is  $O(k \log n)$ . This is the time taken in pairing  $c_i$ 's as well. This proves Lemma 4.4.  $\square$

Armed with Lemma 4.4, we now describe the procedure to merge a star. We use the following procedure to merge the leaves of the star.

**Merge-Star**( $\{c_0, c_1, c_2, \dots, c_k\}$ ). The procedure works in iterations. In each iteration, we try to merge a star. At the end of each iteration, we are left with a star with lesser (roughly half) number of leaves than we had at the beginning of the iteration. The center of the stars in each iteration remains  $c_0$ .

Let  $s_1, s_2, \dots, s_{k^r}$  be the leaves of the star to be merged in the current iteration. Here  $k^r$  is initialized with  $k$ , and  $s_i$ 's are initialized with  $c_i$ 's.

- If  $(k^r = 1)$ , then merge the only leaf with the center of the star and exit.
- Use Lemma 4.4 to pair up the leaves of the star. Without loss of generality, the pairing is  $(s_1, s_2), \dots, (s_{k^r-1}, s_{k^r})$ .
- Merge the clusters in a pair.
- Repeat the procedure with

$$k^r \leftarrow k^r/2$$

and

$$\{s_1, s_2, \dots\} \leftarrow \{s_1 \cup s_2, s_3 \cup s_4, \dots\}.$$

(Note that this is again a star.)

**Lemma 4.5.** *A star of phase  $i$  with  $c_0$  as center and  $c_1, c_2, \dots, c_k$  as leaves can be merged by adding a set of clusters  $W$ , such that*

$$\log_{2^{k_i}} \frac{S(v)}{2^{k_i}} = O(k + k \log \frac{S(c_0)}{2^{k_i}} + \sum_{j=1}^k \frac{S(c_j)}{2^{k_i}}).$$

Also, the time taken in the merging of the star is  $O(|W| \log n)$ .

The proof is described in Appendix A.

#### 4.6.2. Merging of a connected component of $H_i$

Now that we know how to merge a star, we describe how to merge a connected component of  $H_i$ . Let that component be  $C$ . Let  $D = \{c \in C \mid c \text{ is a star}\}$ . In the following procedure to merge the components,  $A$  denotes the set of clusters yet to be merged. It is initially assigned the value  $C$ .

#### Merge-Component.

- If  $|A| = 1$ , then there is nothing to be done.
- Define  $T^r$  to be a tree, with vertex set  $A$  and edge set as  $\{(c_1, c_2) \mid c_1, c_2 \in A \wedge \exists (u_1, u_2) \in G_{i+1}, u_1 \in c_1, u_2 \in c_2\}$ . (Note that this graph is always a tree.)
- Find a minimum vertex cover of  $T^r$ . Note that the smallest vertex cover of a tree has size at most half the size of the tree  $|V| \leq \frac{1}{2}|A|$ .
- For each cluster  $c_0 \in V$ , we have a set of clusters  $\{c_1, c_2, \dots, c_k\} \subset A \setminus V$ , directly connected to it (in other words covered by it).

**Observation 4.8.**  $c_0, c_1, \dots, c_k$  form a star with  $c_0$  as center and  $c_1, c_2, \dots, c_k$  as leaves.

Use the vertex cover  $V$  to partition  $A$  into stars. Since, each star is centered around a vertex of the vertex cover, the number of stars is  $|V|$ .

- Merge the stars using the procedure referred to in Lemma 4.5. After the merging, number of components left is  $|V| \leq |A|/2$ . So, we have reduced the number of components from  $|A|$  to at most  $|A|/2$ .
- Repeat the procedure with  $A$  replaced by set of merged stars.

**Lemma 4.6.** *It is possible to merge a connected component  $C$  of  $H_i$  by adding a set of clusters  $W$ , such that*

$$\log_{1+\epsilon} \frac{r_{\max}(v)}{r_{\min}(v)} = O(\alpha \times \sum_{c \in C} \frac{S(c)}{2^{k_i}}).$$

Time taken is  $O(|W| \log(n))$ .

The proof is described in Appendix A.

The proof of Lemma 4.2 follows by using the bound of Lemma 4.6 on each of the connected components of  $H_i$  and then adding the contributions. This completes the proof of Lemma 4.2.  $\square$

## Appendix A

**Proof of Lemma 4.5.** Since in each iteration of the *merge-star* procedure, the number of leafs decreases by a factor of two, so in the iteration  $t$ ,  $k^r = \frac{k}{2^t}$ .

**Observation A.1.** In the  $t$ th iteration,

$$k^r = \frac{k}{2^t}.$$

Also from the Lemma 4.4, in each iteration, the sum of the  $r_{\max}$  values of clusters increase by at most  $S(c_0) + k2^{k_i+1}$ . So, we get the following observation.

**Observation A.2.** In the  $t$ th iteration,

$$S(s_j) \leq S(c_j) + t S(c_0) + k2^{k_i+1}.$$

**Observation A.3.** Time taken in the  $t$ th iteration is  $O(\frac{k}{2^t} \log(n))$ . This follows, from the fact that, in the merging algorithm, pairing up of the clusters takes  $O(k^r \log(n))$  time and merging takes  $O(k^r)$  time. Since, in the  $t$ th iteration,  $k^r \leq 2k/2^t$ , so the time taken is  $O(\frac{k}{2^t} \log(n))$ .

Now, the contribution of set of clusters that are added in the  $t$ th iteration is given by

$$k^r/2 \log \frac{S(s_{2j-1} \cup s_{2j})}{2^{k_i}}.$$

Using weighted Geometric mean  $\leq$  weighted arithmetic mean,

$$k^r/2 \log \frac{S(s_{2j-1} \cup s_{2j})}{2^{k_i}} \leq \frac{1}{2} k^r \log \frac{S(s_{2j-1}) + S(s_{2j})}{k^r 2^{k_i}}.$$

Using Observations A.1 and A.2, we get for the  $t$ th iteration

$$k^r/2 \log \frac{S(s_{2j-1} \cup s_{2j})}{2^{k_i}} \leq \frac{k}{2^{t-1}} \log 2^t \frac{S(c_j) + t S(c_0) + t2^{k_i+1}}{2^{k_i} k}.$$

So the total contribution is

$$\frac{k}{2^{t-1}} \log 2^t \frac{S(c_j) + t S(c_0) + t2^{k_i+1}}{2^{k_i} k}.$$

It is easily shown that this is

$$O \left( k + k \log \frac{S(c_0)}{k2^{k_i}} + \sum_{j=1}^k \frac{S(c_j)}{k2^{k_i}} \right).$$

Also, adding time taken in each iteration from Observation A.3, we find the total time taken to be  $O(k \log n)$ . This proves the Lemma 4.5.  $\square$

**Proof of Lemma 4.6.** Consider any iteration of the procedure. Suppose,  $c_{rs}$  is the  $s$ th cluster of the  $r$ th star of  $A$ . Suppose,  $h_r$  is the number of clusters in the  $r$ th star.

Then, from Observations 4.4 and 4.3 and Lemma 4.5, we get—we can merge all the stars created in this iteration by adding a set of clusters  $W$ , with

$$\log_{1+\epsilon} \frac{r_{\max}(v)}{r_{\min}(v)} = O \alpha \log_{2^{k_i}} \frac{S(v)}{2^{k_i}} = O \alpha \sum_{r=1}^{h_r} h_r + h_r \log_{h_r 2^{k_i}} \frac{S(c_{rs})}{h_r 2^{k_i}}.$$

From (geometric mean  $\leq$  arithmetic mean), we get

$$\log_{1+\epsilon} \frac{r_{\max}(v)}{r_{\min}(v)} = O \alpha \sum_{r=1}^{h_r} h_r + \sum_{r=1}^{h_r} h_r \log_{h_r 2^{k_i}} \frac{S(c_{rs})}{h_r 2^{k_i}}.$$

Clearly,  $\sum_{rs} S(c_{rs}) \leq S(D)$ , and  $\sum_r h_r = |A|$ . So, this gives

$$\log_{1+\epsilon} \frac{r_{\max}(v)}{r_{\min}(v)} = O \alpha |A| + |A| \log_{2^{k_i} |A|} \frac{S(D)}{2^{k_i} |A|}.$$

Now, in each iteration  $|A|$  reduces by at least a factor of two. Initially, it is  $|C|$ . So, the total contribution from all the clusters added in all the iterations is

$$O \alpha |C| \left( 1 + \frac{1}{2} + \frac{1}{4} + \dots \right) + \sum_{k_i} \frac{1}{2^{k_i} |C|} \left( 1 + \log_{2^{k_i} |C|} \frac{S(D)}{2^{k_i} |C|} + \alpha |C| \left( 1 + \frac{\log 2}{2} + \frac{\log 4}{4} + \dots \right) \right) \\ = O \alpha |C| + |C| \log_{2^{k_i} |C|} \frac{S(D)}{2^{k_i} |C|}.$$

Now, in the sub-graph spanned by  $D$ , there are  $|C| - 1$  edges of length between  $2^{k_i}$  and  $2^{k_i+1}$ . So, from definition of  $S(D)$ ,  $|C| \leq \frac{S(D)}{2^{k_i}}$ . Also,  $|C| \log_{2^{k_i} |C|} \frac{S(D)}{2^{k_i} |C|}$  is  $O(\frac{S(D)}{2^{k_i}})$ . So, total contribution is  $O(\alpha S(D)/2^{k_i})$ .

The time taken in  $k$ th iteration is easily seen to be  $|A| \log(n)$  (using Lemma 4.5). Since, in the  $k$ th iteration,  $|A| \leq |C|/2^k$ , so total time taken in all the iterations is  $O(|C| \log(n) + |C| \log(n)/2 + |C| \log(n)/4 + \dots) = O(|C| \log(n))$ . This proves Lemma 4.6.  $\square$

## Appendix B. Constructing an approximate MST

Two methods were outlined by Har-Peled for estimating MSTs. Building a  $1/\delta$ -MST requires an amount of time equal to  $O(n \log(n))$ . Well Separated Pair Decomposition by Callahan Kossaraju is employed. In addition, he provides a method for building a  $nd$ -MST that uses  $O(n \log(n))$  time. The first method may be used to build a 2-MST that will serve our needs.

Instead of using WSPD, we may build a  $\lambda$ -MST by utilising the methods outlined in this study. This is the way it is done.

- (1) Apply the Har-Peled method to construct  $nd$ -MST  $T_1$  of  $P$ . Just so you know, WSPD is not being used here.
- (2) Create an estimate of  $\epsilon \mu$  for Voronoi Diagram  $V$  using the methods outlined in this article and  $T_1$ . It would consume space of order  $O(n \log n)$ .
- (3) The Delaunay graph of a Voronoi diagram (or approximation Voronoi diagram) is defined as a graph over the diagram's sites, where edges between points  $p_1$  and  $p_2$  are defined as cells of  $p_1$  and  $p_2$  sharing a border, if and only if this formula is true. Define  $V$  and build its Delaunay graph  $F$ .

$G$ 's MST is equal to the  $1 + 3\mu$ -MST of the set  $P$  of points.

The evidence rests on the fact that an approximation MST is included in an approximate Delaunay graph. When dealing with accurate MST and Delaunay graphs, this may be translated into a comparable outcome. To show this, we will show that for any cut  $(P_1, P_2)$  in  $P$ , there is an edge  $(p_1, p_2)$  in the Delaunay graph where the distance between the two edges is at most  $(1 + 3\mu)$ . The distance between points  $P_1$  and  $P_2$  is  $(1 + 3\mu)$ .

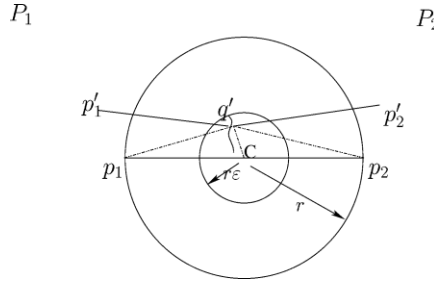


Fig. 5. A cut in an approximate MST.

**Proof.** Suppose,  $p_1, p_2$  be the points such that  $p \in P_i$  and  $\text{dist}(p_1, p_2) = \text{dist}(P_1, P_2)$ . Consider the ball  $b(C, r)$  with  $\overline{p_1, p_2}$  as diameter (Fig. 5). It is easy to show that no point of  $P$  is contained in  $b(C, r)$ . Consider all the Voronoi cells that intersect the ball  $b(C, \epsilon r)$ . Observe that those cells cannot all correspond to one of  $P_i$ . So, there is a point (say  $q'$ ) in  $b(C, \epsilon r)$  which lies at the boundary of two Voronoi cells, one of them corresponding to a point in  $P_1$  and other to a point in  $P_2$ . Let these points be  $p_1'$  and  $p_2'$ . So, both  $p_1'$  and  $p_2'$  are  $\epsilon$ -nearest neighbors of  $q'$ . So,  $\text{dist}(q', p_1') \leq (1 + \epsilon) \text{dist}(q', p_1)$ . Using this and from geometry, we get  $\text{dist}(p_1', p_2') \leq \text{dist}(p_1', q') + \text{dist}(p_2', q') \leq (1 + \epsilon)(\text{dist}(p_1, q') + \text{dist}(p_2, q')) \leq (1 + \epsilon)(r + r\epsilon + r + r\epsilon) \leq (1 + \epsilon)^2 \text{dist}(P_1, P_2) \leq (1 + 3\epsilon) \text{dist}(P_1, P_2)$ . Since there is an edge in the Delaunay graph between  $p_1'$  and  $p_2'$ , the result follows.  $\square$

### Appendix C. Hierarchical clustering for construction of PLSBs

Similar to Har-Peled's [11] method, our PLSB creation strategies rely on a hierarchical clustering of the point set. The main difference between both approaches is in the construction of the hierarchical clustering, but otherwise, the concepts are same. To begin, we have the discrete points that make up the hierarchical clustering's leaves. Each stage involves merging two clusters according to a different metric for "proximity" between them. Each of these groups has a delegate who will lead the construction of a set of balls. You may utilise these balls to find the cluster that has a point that is almost closer to your query point. The representation that the ball holding the query point is built around determines this. You don't need to build extra balls to differentiate between the closeness of the query point to points in the two clusters when it's far enough away from both to report a point from either cluster. The procedure is carried out until every point has been combined into one cluster.

The balls in the PLSB issue are the set of all the balls made in this way.

#### C.1. Answering PLSB using compressed quadtree

- D An effective data structure for solving a PLSB issue was outlined by Har-Peled [11] as follows: after estimating the balls using axis-parallel cubes taken from a meticulously built hierarchical grid, a quadtree search tree was built over these cubes. An overview of the data structure is provided below.

Here we have a PLSB issue using  $m$  balls. Our goal is to build a data structure that efficiently finds the smallest ball that contains a query point  $q$  and reports it. A hierarchical grid of cubes is used to approximate the balls (Fig. 6).

The equation (u) represents the division of  $d$  into a uniform axis parallel grid with the origin as its centre and side length equal to  $2 \log u$  for a real value  $u$ . Changing  $u$ 's value yields a space-covering multi-resolution grid.

Now think of an  $r$ -radius ball with a centre at point  $p$ , denoted as  $b = B(p, r)$ . A collection of cubes of the grid  $G(r/\epsilon d)$  is denoted as  $GC(b, \epsilon)$ .

- $b_\epsilon$  is an  $\epsilon$ -approximation to  $b$ ; and
- $|GC(b, \epsilon)| = O(1/\epsilon^d)$ .

Now consider  $B$ . Using the technique described above, one can derive a set  $GC(B, \epsilon) = \bigcup_{b \in B} GC(b, \epsilon)$  of  $O(m/\epsilon^d)$  cubes from the axis-parallel hierarchical grid that  $\epsilon$ -approximates  $B$ . There is also one "infinite" cube,



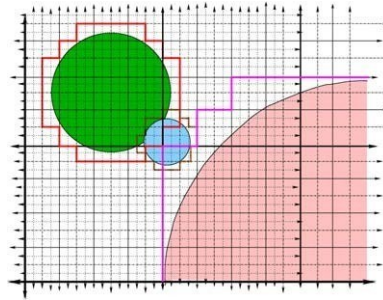


Fig. 6.  $\epsilon$ -PLSB data structure.

it encompasses the whole area. When the query point is distant enough from  $P$  that any point  $x \in P$  is an  $\epsilon$ -closest neighbour to the query point, this backdrop cube is used to answer  $\epsilon$  estimated nearest neighbour queries. Keep in mind that the smaller cube, or the cube that corresponds to the smaller ball, always takes precedence whenever two or more cubes from different balls intersect. This collection contains all the cubes that were created. The cubes may be kept in a quadtree and used to answer point-location queries among the prioritised cubes since they are all obtained from a hierarchical grid representation of  $d$ . Areas that make up the compressed quadtree may be either cubes with the same axis or the annulus, or difference, of two cubes that are within each other. The areas encompass space and are not connected to one another. In addition, there is a ball that corresponds to each zone. If there is any overlap, the area is linked to the smallest ball. The conventional algorithms allow for the computation of in  $O(\log n)$  time. The data structure from [10] may be used to preprocess the leaves of for point-location. The time required to answer a point location query for a point  $q$  is thus  $O(\log |Q|) = O(\log (m/\epsilon))$ .

## References

- [1] S. Arya, D.M. Mount, N.S. Netanyahu, R. Silverman, A.Y. Wu, An optimal algorithm for approximate nearest neighbor searching in fixed dimensions, *J. ACM* 45 (6) (1998) 891–923.
- [2] Sunil Arya, Theodoros Malamatos, Linear-size approximate Voronoi diagrams, in: *Proc. of SODA*, 2002.
- [3] Sunil Arya, Theodoros Malamatos, David Mount, Space-efficient approximate Voronoi diagrams, in: *Proc. of ACM STOC*, 2002.
- [4] Sunil Arya, David M. Mount, Approximate nearest neighbor queries in fixed dimension, in: *Proc. of SIAM–ACM SODA*, 1993, pp. 271–280.
- [5] D. Attali, J.D. Boissonnat, Complexity of the Delaunay triangulation of points on a smooth surface, <http://www-sop.inria.fr/prisme/personnel/boissonnat/papers.html>, 2001.
- [6] F. Aurenhammer, Voronoi diagrams: A survey of a fundamental geometric data structure, *ACM Comput. Surv.* 23 (1991) 345–405.
- [7] P.B. Callahan, S.R. Kosaraju, A decomposition of multidimensional point sets with applications to  $k$ -nearest-neighbors and  $n$ -body potential fields, *J. ACM* 42 (1995) 67–90.
- [8] Kenneth L. Clarkson, An algorithm for approximate closest-point queries, in: *Proc. 10th Annu. ACM Sympos. Comput. Geom.*, 1994, pp. 160–164.
- [9] J. Erickson, Nice points sets can have nasty Delaunay triangulations, in: *Proc. 17th Annu. ACM Sympos. Comput. Geom.*, 2001.
- [10] G.N. Frederickson, Data structures for on-line updating of minimum spanning trees, with applications, *SIAM J. Comput.* 14 (4) (1985) 781–798.
- [11] S. Har-Peled, A replacement for Voronoi diagrams of near linear size, in: *Proc. of IEEE FOCS*, 2001, pp. 94–103, full version available from <http://www.uiuc.edu/~sariel/papers>.
- [12] P. Indyk, R. Motwani, Approximate nearest neighbors: Towards removing the curse of dimensionality, in: *Proc. 30th Annu. ACM Sympos. Theory Comput.*, 1998, pp. 604–613.
- [13] Piotr Indyk, PhD thesis, Stanford University, 1999.
- [14] Eyal Kushilevitz, Rafail Ostrovsky, Yuval Rabani, Efficient search for approximate nearest neighbour search in high dimensional spaces, in: *ACM Symposium on Theory of Computing*, 1998, pp. 614–623.