

# Faster DBScan and HDBScan in Low-Dimensional Euclidean Spaces\*

Mark de Berg<sup>†1</sup>, Ade Gunawan<sup>2</sup>, and Marcel Roeloffzen<sup>‡3</sup>

1 Department of Computing Science, TU Eindhoven, Eindhoven, The Netherlands

[mberg@win.tue.nl](mailto:mberg@win.tue.nl)

2 Department of Computing Science, TU Eindhoven, Eindhoven, The Netherlands

3 National institute of informatics, Tokyo and JST ERATO, Kawarabayashi Large Graph Project, Japan

[marcel@nii.ac.jp](mailto:marcel@nii.ac.jp)

---

## Abstract

We present a new algorithm for the widely used density-based clustering method DBSCAN. Our algorithm computes the DBSCAN-clustering in  $O(n \log n)$  time in  $\mathbb{R}^2$ , irrespective of the scale parameter  $\varepsilon$ , but assuming the second parameter MINPTS is set to a fixed constant, as is the case in practice. We also present an  $O(n \log n)$  randomized algorithm for HDBSCAN in the plane – HDBSCAN is a hierarchical version of DBSCAN introduced recently – and we show how to compute an approximate version of HDBSCAN in near-linear time in any fixed dimension.

**1998 ACM Subject Classification** F.2.2 Nonnumerical Algorithms and Problems, H.3.3 Information Search and Retrieval

**Keywords and phrases** Density-based clustering, hierarchical clustering

**Digital Object Identifier** 10.4230/LIPIcs.ISAAC.2017.25

## 1 Introduction

Clustering is one of the most fundamental tasks in data mining. Due to the wide variety of applications where clustering is important, the clustering problem comes in many variants. These variants differ for example in the dimensionality of the data set  $D$  and in the underlying metric, but also in the objective of the clustering. Thus a multitude of clustering algorithms has been developed [21], each with their own strengths and weaknesses. We are interested in *density-based clustering*, where clusters are defined by areas in which the density of the data points is high and clusters are separated from each other by areas of low density.

One of the most popular density-based clustering methods is DBSCAN; the paper by Ester *et al.* [12] on DBSCAN has been cited over 8,800 times, and in 2014 DBSCAN received the test-of-time award from KDD, a leading data-mining conference. DBSCAN has two parameters,  $\varepsilon$  and MINPTS, that together determine when the density around a point  $p \in D$  is high enough for  $p$  to be part of a cluster as opposed to being noise; see Section 2 for a precise definition of the DBSCAN clustering. Typically MINPTS is a constant – in the original article [12] it is concluded that  $\text{MINPTS} = 4$  works well – but finding the right value for  $\varepsilon$

---

\* A full version of the paper is available at [7], <https://arxiv.org/abs/1702.08607>.

<sup>†</sup> MdB is supported by the Netherlands Organization for Scientific Research under grant 024.002.003.

<sup>‡</sup> MR is supported by JST ERATO Grant Number JPMJER1201, Japan.



© Mark de Berg, Ade Gunawan, and Marcel Roeloffzen;  
licensed under Creative Commons License CC-BY

28th International Symposium on Algorithms and Computation (ISAAC 2017).

Editors: Yoshio Okamoto and Takeshi Tokuyama; Article No. 25; pp. 25:1–25:13

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

is more difficult. The worst-case running time of the original DBSCAN algorithm is  $\Theta(n^2)$ . It is often stated that the running time is  $O(n \log n)$  for Euclidean spaces when a suitable indexing structure such as an R-tree is used to support the DBSCAN algorithm. While this may be true in certain practical cases, it is not true from a theoretical point of view.

Several variants of DBSCAN algorithm have been proposed, often with the goal to speed up the computation. Some (IDBSCAN [5] and FDBSCAN [16]) do so at the expense of computing a slightly different, and not clearly defined, clustering. Others (GridBSCAN [17]) compute the same clustering as DBSCAN, but without speeding up the worst-case running time.

A fundamental bottleneck of the original DBSCAN algorithm is that it performs a query with each point  $p \in D$  to find  $N_\varepsilon(p, D)$ , the set of points within distance  $\varepsilon$  of  $p$ . The algorithm uses these points to continue expanding the cluster, hence, range counting would not be sufficient. It follows that  $\sum_{p \in D} |N_\varepsilon(p, D)|$  is a lower bound on the running time of the DBSCAN algorithm. In the worst case  $\sum_{p \in D} |N_\varepsilon(p, D)| = \Theta(n^2)$ , so even with a fast indexing structure the worst-case running time of the original DBSCAN algorithm is  $\Omega(n^2)$ . (Apart from this, the worst-case query time of R-trees and other standard indexing structures is not logarithmic even if we disregard the time to report points.) In most practical instances the DBSCAN algorithm is much faster than quadratic. The reason is that  $\varepsilon$  is typically small so that the sets  $N_\varepsilon(p, D)$  do not contain many points and the range queries can be answered quickly. However, the fact that the algorithm always explicitly reports the sets  $N_\varepsilon(p, D)$  makes the running time sensitive to the choice of  $\varepsilon$  and the density of the point set  $D$ . For example, suppose we have a disk-shaped cluster with a Gaussian distribution around the disk center. Then a suitable value of  $\varepsilon$  will lead to large sets  $N_\varepsilon(p, D)$  for points  $p$  near the center of the cluster.

Chen *et al.* [9] overcame the quadratic bottleneck of the standard approach, and designed an algorithm<sup>1</sup> with  $O(n^{2-\frac{2}{d+2}} \text{polylog } n)$  worst-case running time. They also present an approximate algorithm. Note that for  $d = 2$  the running time of the exact algorithm is  $O(n^{1.5} \text{polylog } n)$ . Chen *et al.* remark that their exact algorithm is mainly of theoretical interest. The natural question is then whether or not it is possible to compute the DBSCAN clustering in subquadratic time in the worst case, irrespective of the value of  $\varepsilon$ , with a simple and practical algorithm.

Although DBSCAN is used extensively and performs well in many situations, it has its drawbacks. One is that it produces a flat, non-hierarchical clustering which heavily depends on the choice of the scale parameter  $\varepsilon$ . Ankerst *et al.* [3] therefore introduced OPTICS, which can be seen as a hierarchical version of DBSCAN. Recently Campello *et al.* [8] proposed an improved density-based hierarchical clustering method – similar to OPTICS but cleaner – together with a cluster-stability measure that can be used to automatically extract relevant clusters. The new method, called HDBSCAN, only needs the parameter MINPTS, which is much easier to choose than  $\varepsilon$ . Campello *et al.* used MINPTS=4 in all their experiments. While HDBSCAN is very powerful, the algorithm to compute the HDBSCAN hierarchy runs in quadratic time; not only in the worst-case, but actually also in the best-case. There have been only few papers dealing with speeding up HDBSCAN or its predecessor OPTICS. A notable recent exception is POPTICS [20], a parallel algorithm that computes a similar, but not the same, hierarchy as OPTICS. We do not know of any algorithm that computes the HDBSCAN or OPTICS hierarchy in subquadratic time. Thus the second question we study is: is it possible to compute the HDBSCAN hierarchy in subquadratic time.

---

<sup>1</sup> As described, the algorithm actually computes a variation of the DBSCAN clustering, but it is easily adapted to compute the true DBSCAN clustering.

**Our results.** We present an  $O(n \log n)$  algorithm to compute the DBSCAN clustering for a set  $D$  of  $n$  points in the plane, irrespective of the setting of the parameter  $\varepsilon$  used to define the DBSCAN clustering. Here, and in our other results, we assume that the parameter  $\text{MINPTS}$  is a fixed constant. As mentioned this is the case in practice, where one typically uses  $\text{MINPTS} = 4$ . We remark that our algorithm is not only fast in theory, but a slightly simplified version is also competitive in practice and much less sensitive to the choice of  $\varepsilon$  than the original DBSCAN algorithm. In this submission we focus on our theoretical contributions. Experimental results can be found in the full version [7].

We also present a new algorithm for planar HDBSCAN: we show how to compute the HDBSCAN hierarchy in  $\mathbb{R}^2$  in  $O(n \log n)$  expected time, thus obtaining the first subquadratic algorithm for the problem.

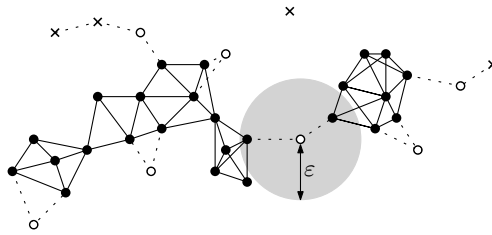
In higher dimensions exact algorithms are much slower and so we consider approximations instead. We extend the concept of an approximate DBSCAN clustering as defined by Chen *et al.* [9] and by Gan and Tao [13] (see below) to the hierarchical version. We thus obtain  $\delta$ -approximate HDBSCAN, an approximate version of the HDBSCAN hierarchy of Campello *et al.* [8], where the parameter  $\delta$  specifies the accuracy of the approximation. Intuitively, a  $\delta$ -approximate HDBSCAN hierarchy has the same clusters as the standard HDBSCAN hierarchy at any level  $\varepsilon$ , except that clusters at distance  $(1 - \delta) \cdot \varepsilon$  from each other may be merged, see Section 5 for a precise definition. We show that a  $\delta$ -approximate HDBSCAN hierarchy in  $\mathbb{R}^d$  can be computed in  $O((n/\delta^{(d-1)/2}) \log^{d-1} n)$  time.

**Further related work.** Our paper is the conference paper corresponding to the so far unpublished master’s thesis of the second author [15], which contained the results on DBSCAN, extended with results on HDBSCAN. In the meantime, Gan and Tao [13] published a paper in which they extend the work from the master’s thesis to  $\mathbb{R}^d$ , resulting in an algorithm for DBSCAN with a running time of  $O(n^{2 - \frac{2}{d+1} + \gamma})$ ; we briefly comment on how this is done at the end of Section 3. Gan and Tao also prove that computing the DBSCAN clustering in  $\mathbb{R}^d$  for  $d \geq 3$  is at least as hard as the so-called unit-spherical emptiness problem, which is believed to require  $\Omega(n^{4/3})$  time [11]. Finally, Gan and Tao show that a  $\delta$ -approximate DBSCAN clustering can be computed in  $O(n/\delta^{d-1})$  expected time, using a modified version of the exact algorithm. Their approximate clustering is the same as the approximate clustering defined by Chen *et al.* [9], who already showed how to compute it in  $O(n \log n + n/\delta^{d-1})$  time deterministically. (Gan and Tao were unaware of the paper by Chen *et al.*.) As we remark in Section 5 our algorithm can also be used to obtain a deterministic algorithm with  $O(n \log n + n/\delta^{d/3+c})$  running time for some constant  $c$ .

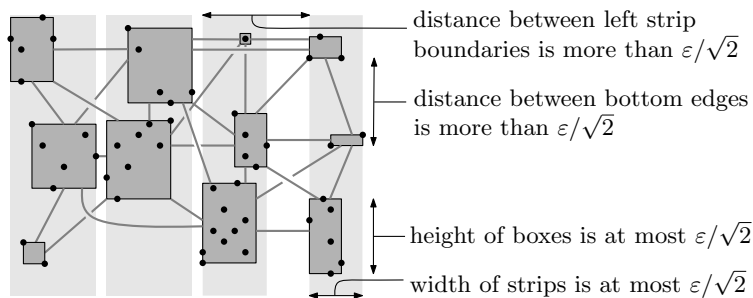
## 2 Preliminaries on DBScan and DBScan\*

Let  $D$  be a set of points in  $\mathbb{R}^d$ . DBSCAN distinguishes three types of points: *core points* in the “interior” of a cluster, *border points* on the boundary of a cluster, and *noise points* not in any cluster. The distinction is based on two global parameters,  $\varepsilon$  and  $\text{MINPTS}$ . Define  $N_\varepsilon(p, D) := \{q \in D : |pq| \leq \varepsilon\}$  to be the *neighborhood* of a point  $p$ , where  $|pq|$  denotes the (Euclidean) distance between  $p$  and  $q$ ; the point  $p$  itself is included in  $N_\varepsilon(p, D)$ . A point  $p \in D$  is a *core point* if  $|N_\varepsilon(p, D)| \geq \text{MINPTS}$ , and a non-core point  $q$  in the neighborhood of a core point is a *border point*. We denote the set of core points by  $D_{\text{core}}$ , and the set of border points by  $D_{\text{border}}$ . The remaining points are *noise*. In DBSCAN\* [8] border points are not part of a cluster but are considered noise.

Ester *et al.* [12] define the DBSCAN clusters based on the concept of density-reachability.



■ **Figure 1** A neighborhood graph with  $\text{MINPTS} = 4$  and  $\varepsilon$  as indicated. Solid disks are core points, open circles are border points, and crosses are noise. Edges between core points are solid, other edges are dotted. The solid disks and edges form the core graph.



■ **Figure 2** Example of a box graph.

Equivalently, we can define the clusters as the connected components of a certain graph.

To this end, define the *neighborhood graph*  $\mathcal{G}(D, E)$  as the undirected graph with node set  $D$  and edges connecting pairs of points within distance  $\varepsilon$ ; see Fig. 1. Note that a point  $p \in D$  is a core point if and only if its degree in  $\mathcal{G}$  is at least  $\text{MINPTS} - 1$ , since then its neighborhood contains at least  $\text{MINPTS}$  points, including  $p$  itself. Now consider the subgraph  $\mathcal{G}_{\text{core}}(D_{\text{core}}, E_{\text{core}})$  induced by the core points, that is,  $\mathcal{G}_{\text{core}}$  is the graph whose nodes are the core points and whose edges connect two core points when they are within distance  $\varepsilon$  from each other. We call  $\mathcal{G}_{\text{core}}$  the *core graph*. The connected components of  $\mathcal{G}_{\text{core}}$  are the clusters in DBSCAN\*. The clusters in DBSCAN are the same, except that they also contain border points. Formally, a border point  $q$  belongs to a cluster  $C$  if  $q$  has an edge in  $\mathcal{G}$  to a core point  $p \in C$ . Thus a border point can belong to multiple clusters. The original DBSCAN algorithm constructs clusters one by one and assigns a border point  $p$  to the first cluster that finds  $p$ ; we assign border points to the cluster of their nearest core point.

### 3 A fast algorithm for DBScan

The original DBSCAN algorithm reports, while generating and exploring the clusters, for each point  $p \in D$  all its neighbors. In other words, it spends time on every edge in the neighborhood graph. Our new algorithm avoids this by working with a smaller graph, the *box graph*  $\mathcal{G}_{\text{box}}$ . Its nodes are disjoint rectangular boxes with a diameter of at most  $\varepsilon$  that together contain all the points in  $D$ , and its edges connect pairs of boxes within distance  $\varepsilon$ ; see Fig. 2.

The boxes are generated such that (i) any two points in the same box are in each other's neighborhood, and (ii) the degree of any node in the box graph is  $O(1)$ . Property (i) allows us to immediately classify all points in a box as core points when it contains at least  $\text{MINPTS}$  points, and property (ii) allows us to quickly retrieve the neighbors of any given point in a box. Next we describe the algorithm, which consists of four easy steps, in detail.

**Step 1: Compute the box graph  $\mathcal{G}_{\text{box}}$ .** To compute  $\mathcal{G}_{\text{box}}$ , we first construct a collection of vertical *strips* that together cover all the points. Let  $p_1, \dots, p_n$  be the points in  $D$  sorted by  $x$ -coordinate, with ties broken arbitrarily. The first strip has  $p_1$  on its left boundary. We continue from left to right, adding points to the first strip as we go, until we encounter a point  $p_i$  whose distance to the left strip boundary is more than  $\varepsilon/\sqrt{2}$ . We then start a new strip with  $p_i$  on its left boundary, and we add points to that strip until we encounter a point whose distance to the left strip boundary is more than  $\varepsilon/\sqrt{2}$ , and so on, until we handled all the points. Constructing the strips takes  $O(n)$  time, after sorting the points by  $x$ -coordinate.

Within each strip we perform a similar procedure, going over the points within the strip in order of increasing  $y$ -coordinate and creating boxes instead of strips. Thus the first box in the strip has the lowest point on its bottom edge, and we keep adding points to this box (enlarging it so that the new point fits, ensuring a tight bounding box) until we encounter a point whose vertical distance to the bottom edge is more than  $\varepsilon/\sqrt{2}$ . We then start a new box, and so on, until we handled all points in the strip. If the number of points in the  $j$ -th strip is  $n_j$ , then the time needed to handle all the strips is  $\sum_j O(n_j \log n_j) = O(n \log n)$ .

Let  $m$  be the number of strips and  $\mathcal{B}_j$  the set of boxes in the  $j$ -th strip. We sometimes refer to a set  $\mathcal{B}_j$  as a strip, even though formally  $\mathcal{B}_j$  is a set of boxes. Let  $\mathcal{B} := \mathcal{B}_1 \cup \dots \cup \mathcal{B}_m$ . The nodes of the box graph  $\mathcal{G}_{\text{box}}$  are the boxes in  $\mathcal{B}$  and there is an edge  $(b, b')$  when  $\text{dist}(b, b') \leq \varepsilon$ , where  $\text{dist}(b, b')$  denote the minimum distance between  $b$  and  $b'$ . Two boxes  $b, b'$  are *neighbors* when they are connected by an edge. Let  $\mathcal{N}_\varepsilon(b, \mathcal{B})$  be the neighbors of  $b$ .

► **Lemma 1.**  $\mathcal{G}_{\text{box}}$  has at most  $n$  nodes, each having  $O(1)$  neighbors.

The lemma above follows from the fact that any box  $b \in \mathcal{B}_j$  can have neighbors only in  $\mathcal{B}_{j-2}, \mathcal{B}_{j-1}, \mathcal{B}_j, \mathcal{B}_{j+1},$  or  $\mathcal{B}_{j+2}$ , and within any of these five strips,  $b$  can have at most five neighbors. (A more precise proof giving a bound of 22 neighbors can be found in the full version [7].) This also gives us an easy way to compute the edge set  $E_{\text{box}}$  of the box graph, because the edges between boxes in strips  $\mathcal{B}_j$  and  $\mathcal{B}_{j'}$  with  $|j - j'| \leq 2$  can be computed in  $O(|\mathcal{B}_j| + |\mathcal{B}_{j'}|)$  time in total by scanning the boxes in  $\mathcal{B}_j$  and  $\mathcal{B}_{j'}$  in a coordinated manner. The total time to compute all edges of the box graph is thus

$$O\left(\sum_{j=1}^m \sum_{j'=\max(j-2,1)}^{\min(j+2,m)} (|\mathcal{B}_j| + |\mathcal{B}_{j'}|)\right) = O\left(\sum_{j=1}^m |\mathcal{B}_j|\right) = O(n).$$

Adding the time to construct the strips and boxes, we see that Step 1 takes  $O(n \log n)$  time and we obtain the following lemma.

► **Lemma 2.** The box graph  $\mathcal{G}_{\text{box}}(\mathcal{B}, E_{\text{box}})$  can be computed in  $O(n \log n)$  time.

**An alternative for Step 1.** An alternative approach is to define the boxes as the non-empty cells in a grid whose cells have height and width  $\varepsilon/\sqrt{2}$ . If we store the boxes in a hash-table based on the coordinates of their lower left corners, then finding the neighbors of a box  $b$  can be done by checking each potential neighbor cell for existence in the hash-table – we do not need to store the box graph explicitly. Creating the boxes (with their corresponding point sets) can be done in  $O(n)$  time if the floor function can be computed in  $O(1)$  time.

**Step 2: Find the core points.** The graph  $\mathcal{G}_{\text{box}}$  allows us to determine the core points in a simple and efficient manner. The key observation is that the maximum distance between any two points in the same box is at most  $\varepsilon$ . Hence, if a box contains more than  $\text{MINPTS}$  points, then all of them are core points. The following algorithm suffices to find the core points.

For a box  $b \in \mathcal{B}$ , let  $D(b) := D \cap b$  be the set of point inside  $b$ , and let  $n_b := |D(b)|$ . If  $n_b \geq \text{MINPTS}$  then label all points in  $b$  as core points. Otherwise, for each point  $p \in D(b)$ , count the number of points  $q$  in neighboring boxes of  $b$  for which  $|pq| \leq \varepsilon$ . If this number is at least  $\text{MINPTS} - n_b$ , then label  $p$  as core point. The counting is done brute-force, by checking all points in neighboring boxes. Hence, this takes  $O(\sum_{b' \in \mathcal{N}_\varepsilon(b, \mathcal{B})} n_{b'})$  time for each point  $p \in b$ .

► **Lemma 3.** *Given  $\mathcal{G}_{\text{box}}$ , we can find all core points in  $D$  in  $O(n)$  time.*

**Proof.** The total time spent to handle boxes  $b$  with  $n_b \geq \text{MINPTS}$  is clearly  $O(n)$ . The time needed to handle a box  $b$  with  $n_b < \text{MINPTS}$  is

$$O\left(n_b \cdot \sum_{b' \in \mathcal{N}_\varepsilon(b, \mathcal{B})} n_{b'}\right) = O\left(\text{MINPTS} \cdot \sum_{b' \in \mathcal{N}_\varepsilon(b, \mathcal{B})} n_{b'}\right).$$

Now charge  $O(\text{MINPTS}) = O(1)$  time to each point in every  $b' \in \mathcal{N}_\varepsilon(b, \mathcal{B})$ . Because any box  $b'$  is the neighbor of  $O(1)$  other boxes by Lemma 1, each point is charged  $O(1)$  times. ◀

**Step 3: Compute the cluster cores.** The *core of a cluster* is the set of core points in that cluster. In Step 3 we assign to each core point a *cluster-id* so that core points in the same cluster have the same cluster-id. Again, this can be done in an efficient manner using  $\mathcal{G}_{\text{box}}$ . To this end, we first remove certain boxes and edges from  $\mathcal{G}_{\text{box}}$  to obtain a reduced box graph  $\mathcal{G}_{\text{box}}^*$ . More precisely, we keep only the boxes with at least one core point, and we keep only the edges  $(b, b')$  for which there are core points  $p \in b, p' \in b'$  with  $|pp'| \leq \varepsilon$ . Because any two core points in a given box  $b$  are connected in  $\mathcal{G}_{\text{core}}$ , we have the following lemma.

► **Lemma 4.** *The connected components in  $\mathcal{G}_{\text{box}}^*$  correspond one-to-one to the connected components in the core graph  $\mathcal{G}_{\text{core}}$  and, hence, to the DBSCAN\* clusters.*

Thus the cluster cores can be computed by computing the connected components in  $\mathcal{G}_{\text{box}}^*$ . The latter can be done in  $O(n)$  time using DFS [10]. We then give every core point  $p$  a cluster-id corresponding to the connected component of the box  $b$  that contains  $p$ .

To construct  $\mathcal{G}_{\text{box}}^*$ , we need to decide for two given boxes  $b, b'$  whether there are core points  $p \in D(b), p' \in D(b')$  with  $|pp'| \leq \varepsilon$ . For ease of discussion we call the points in  $D(b)$  blue and those in  $D(b')$  red. It is well known [2] that the bichromatic closest pair defines an edge of the Delaunay triangulation of the points, so it suffices to compute the Delaunay triangulation of  $D(b) \cup D(b')$  and find the shortest red-blue edge. If it is at most  $\varepsilon$  we connect  $b$  and  $b'$  in  $\mathcal{G}_{\text{box}}^*$  and otherwise we do not. This leads to the following lemma.

► **Lemma 5.** *Computing the cluster cores can be done in  $O(n \log n)$  time.*

**Proof.** The most time consuming part of the construction of  $\mathcal{G}_{\text{box}}^*$  is to determine for each pair of neighboring boxes in  $\mathcal{B}$  whether there are core points  $p \in b, p' \in b'$  with  $|pp'| \leq \varepsilon$ . Let  $\mathcal{B}^*$  be the set of boxes containing at least  $\text{MINPTS}$  points. Then the total time spent on the pairs of boxes from  $\mathcal{B}^*$  is

$$\sum_{b \in \mathcal{B}^*} \sum_{b' \in \mathcal{N}_\varepsilon(b, \mathcal{B}^*)} O((n_b + n_{b'}) \log(n_b + n_{b'})),$$

which is  $O(n \log n)$  because  $|\mathcal{N}_\varepsilon(b, \mathcal{B}^*)| = O(1)$  for any box  $b$  and  $\sum_{b \in \mathcal{B}^*} n_b \leq n$ . ◀



► **Remark.** In practice, computing the Delaunay triangulation is not necessary. Instead we can use a brute-force algorithm that checks every pair of points in  $b$  and  $b'$  and stops when a sufficiently close pair is found. The number of points in each box is expected to be small and if it is large one may expect many pairs to have a short distance, hence, testing pairs in random order should find such a pair fairly quickly.

**Step 4: Assigning border points to clusters.** It remains to decide for non-core points  $p$  whether  $p$  is a border point or noise. If  $p$  is a border point, it has to be assigned to the nearest cluster. Again, a brute-force method suffices: for each box  $b \in B$  and each non-core point  $p \in b$ , we check all points in  $b$  and its neighboring boxes to find  $p$ 's nearest core point,  $p'$ . If  $|pp'| \leq \varepsilon$ , then  $p$  is a border point in the same cluster as  $p'$ , otherwise  $p$  is noise. We only need to consider boxes  $b$  with  $n_b < \text{MINPTS}$  – otherwise all points in  $b$  are core points – so the argument from the proof of Lemma 3 shows that this takes  $O(n)$  time.

**Putting it all together.** Steps 1 and 3 take  $O(n \log n)$  time and Steps 2 and 4 take  $O(n)$  time. We thus obtain the following theorem.

► **Theorem 6.** *Let  $D$  be a set of  $n$  points in  $\mathbb{R}^2$ , and  $\varepsilon$  and  $\text{MINPTS}$  be given constants. Then we can compute a DBSCAN clustering on  $D$  according to  $\varepsilon$  and  $\text{MINPTS}$  for the Euclidean metric in  $O(n \log n)$  time.*

► **Remark (Extension to higher dimensions.)** The algorithm just described can easily be extended to  $\mathbb{R}^d$  for  $d > 2$ , as already observed by Gan and Tao [13]. The resulting running time is  $O(n^{2 - \frac{2}{\lceil d/2 \rceil + 1} + \gamma})$ .

## 4 A fast algorithm for HDBScan in the plane

Campello *et al.* [8] introduced HDBSCAN, a hierarchical version of DBSCAN\* similar to OPTICS [3]. The algorithm described by Campello *et al.* to compute the HDBSCAN hierarchy runs in quadratic time. We show that in  $\mathbb{R}^2$  and under the Euclidean metric, the HDBSCAN hierarchy can be computed in  $O(n \log n)$  time.

**Preliminaries on HDBScan.** Recall that DBSCAN\* is the version of DBSCAN in which border points are considered noise. The HDBSCAN hierarchy is a tree structure encoding the clusterings of DBSCAN\* that arise as  $\varepsilon$  increases from  $\varepsilon = 0$  to  $\varepsilon = \infty$  for a fixed  $\text{MINPTS}$ . Initially, when  $\varepsilon = 0$ , all points are noise. As  $\varepsilon$  increases, three types of events can happen to the DBSCAN\* clustering:

- *Type (i): the status of a point changes.* In this event, a point changes from being noise to being a core point. The value of  $\varepsilon$  at which this happens for a point  $p$  is called the *core distance* of  $p$ ; we denote it by  $d_{\text{core}}(p)$ .
- *Type (ii): a new cluster starts.* This event is triggered by a type (i) event, when a point becoming a core point forms a new singleton cluster.
- *Type (iii): two clusters merge.* This event can be triggered by a type (i) event or it can happen when  $\varepsilon = |pq|$  for core points  $p, q$  from different clusters.

Note that all events happen at values of  $\varepsilon$  such that  $\varepsilon = |pq|$  for some pair of points  $p, q \in D$ . This process can be modeled as a *dendrogram*: a tree whose leaves correspond to the points in  $D$  and whose nodes correspond to clusters arising during the process. This dendrogram, where each node stores the value of  $\varepsilon$  at which the corresponding cluster was created, is the HDBSCAN hierarchy. Campello *et al.* compute the HDBSCAN hierarchy as follows.

For two points  $p, q \in D$ , define  $d_{\text{mr}}(p, q) := \max(d_{\text{core}}(p), d_{\text{core}}(q), |pq|)$  to be the *mutual reachability distance* of  $p$  and  $q$ . The *mutual reachability graph*  $\mathcal{G}_{\text{mr}}$  is defined as the complete graph with node set  $D$  in which each edge  $(p, q)$  has weight  $d_{\text{mr}}(p, q)$ . Campello *et al.* observe that HDBSCAN hierarchy can easily be computed from a minimum spanning tree (MST) on  $\mathcal{G}_{\text{mr}}$ . (Indeed, the cluster-growing process corresponds to the computation of an MST on  $\mathcal{G}_{\text{mr}}$  using Kruskal's algorithm [10].) Hence, they compute the HDBSCAN hierarchy as follows.

1. Compute the core distances  $d_{\text{core}}(p)$  for all points  $p \in D$ .
2. Compute an MST  $\mathcal{T}$  of the mutual reachability graph  $\mathcal{G}_{\text{mr}}$ .
3. Convert  $\mathcal{T}$  into a dendrogram where each internal node stores the value of  $\varepsilon$  at which the corresponding cluster is formed.

**Our planar algorithm.** The most time-consuming parts in the algorithm above are Steps 1 and 2; Step 3 takes  $O(n)$  time after sorting the edges of  $\mathcal{T}$  by weight.

For Step 1 we observe that  $d_{\text{core}}(p)$  is the distance of point  $p$  to its  $\ell$ -th nearest neighbor for  $\ell = \text{MINPTS} - 1$ . Hence, to compute all core distances it suffices to compute for each point its  $k$  nearest neighbors. This can be done in any fixed dimension in  $O(n\ell \log n)$  time [22]. Since  $\ell = \text{MINPTS} - 1 = O(1)$  this implies that Step 1 takes  $O(n \log n)$  time.

Step 2 is more difficult to do in subquadratic time. The main problem is that we cannot afford to look at all edges of  $\mathcal{G}_{\text{mr}}$  when computing  $\mathcal{T}$ . To overcome this problem we need the following generalization of Delaunay triangulations, introduced by Gudmundsson *et al.* [14]. Recall that a pair of points  $p, q \in D$  forms an edge in the Delaunay triangulation of  $D$  if and only if there is a circle with  $p$  and  $q$  on its boundary and no points from  $D$  in its interior [6]. We say that the pair  $p, q \in D$  forms a *k-th order Delaunay edge*, or *k-OD edge* for short, if and only if there exists a circle with  $p$  and  $q$  on its boundary and at most  $k$  points from  $D$  in its interior [14]. Thus the 0-OD edges are precisely the edges of the Delaunay triangulation. The  $k$ -OD edges are useful for us because of the following lemma.

► **Lemma 7.** *Let  $\bar{\mathcal{G}}_{\text{mr}}$  be the subgraph of  $\mathcal{G}_{\text{mr}}$  that contains only the  $k$ -OD edges, where  $k := \max(\text{MINPTS} - 3, 0)$ . Then an MST of  $\bar{\mathcal{G}}_{\text{mr}}$  is also an MST of  $\mathcal{G}_{\text{mr}}$ .*

**Proof.** Imagine computing an MST  $\mathcal{T}$  on  $\mathcal{G}_{\text{mr}}$  using Kruskal's algorithm [10]. This algorithm treats the edges  $(p, q)$  of  $\mathcal{G}_{\text{mr}}$  in order of increasing weight, that is, increasing values of  $d_{\text{mr}}(p, q)$ . When it processes  $(p, q)$  it checks if  $p$  and  $q$  are already in the same connected component – in our application this component corresponds to a cluster at the current value of  $\varepsilon$  – and, if not, merges these components. We will argue that whenever we process an edge  $(p, q)$  that is not in  $\bar{\mathcal{G}}_{\text{mr}}$ , that is, an edge that is not a  $k$ -OD edge, then  $p$  and  $q$  are already in the same connected component. Hence, there is no need to process  $(p, q)$ , which proves that an MST of  $\bar{\mathcal{G}}_{\text{mr}}$  is also an MST of  $\mathcal{G}_{\text{mr}}$ .

Let  $C_{pq}$  be the circle such that  $p$  and  $q$  form a diametrical pair of  $C$ , and let  $D(C_{pq}) \subset D$  be the set of points lying in the interior of  $C_{pq}$ . If  $|D(C_{pq})| \leq k$ , then  $(p, q)$  is a  $k$ -OD edge, so assume  $|D(C_{pq})| \geq k + 1$ . Note that  $d_{\text{core}}(r) < |pq|$  for all  $r \in D(C_{pq})$ . Indeed, since  $r$  is an interior point in a disk with diameter  $|pq|$ , the distance from  $r$  to any other point in  $C_{pq}$ , including  $p$  and  $q$ , is smaller than  $|pq|$ . Hence, for  $\varepsilon = |pq|$  we have  $|N_\varepsilon(r, D)| \geq |D(C_{pq})| + 2 = k + 3 \geq \text{MINPTS}$ . Thus all points  $r \in C_{pq}$  are core points when we process  $(p, q)$ . Moreover, for all edges  $(s, t)$  with  $s, t \in D(C_{pq}) \cup \{p, q\}$  we have  $d_{\text{mr}}(s, t) \leq |pq|$ . Hence, it suffices to prove the following.

► **Claim.** *Let  $C$  be a circle with two points  $p, q$  on its boundary and let  $D(C) \subset D$  be the set of points from  $D$  in the interior of  $C$ . Then there is a path from  $p$  to  $q$  in  $\bar{\mathcal{G}}_{\text{mr}}$  that uses only points in  $D(C) \cup \{p, q\}$ .*



We prove this claim by induction on  $|D(C)|$ . If  $|D(C)| \leq k$  then  $(p, q)$  is a  $k$ -OD edge itself and we are done. Otherwise, pick any point  $r \in D(C)$ . Now shrink  $C$ , while keeping  $p$  in its boundary, until we obtain a circle  $C_1$  that also has  $r$  on its boundary. By induction, there is a path from  $p$  to  $r$  in  $\overline{\mathcal{G}}_{\text{mr}}$  that uses only points in  $D(C_1) \cup \{p, r\} \subset D(C) \cup \{p, q\}$ . A similar argument shows that there is a path from  $r$  to  $q$  that uses only points in  $D(C) \cup \{p, q\}$ . This proves the claim and, hence, the lemma.  $\blacktriangleleft$

Gudmundsson *et al.* showed that the number of  $k$ -OD edges is  $O(n(k+1))$  and that the set of all  $k$ -OD edges can be computed in  $O(n(k+1) \log n)$  time with a randomized incremental algorithm. Lemma 7 implies that after computing the core distances and the  $k$ -OD edges in  $O(n \log n)$  time with  $k = \max(\text{MINPTS} - 3, 0) = O(1)$  we can compute the MST for  $\mathcal{G}_{\text{mr}}$  by considering only  $O(n)$  edges. Thus computing the MST can be done in  $O(n \log n)$  time [10]. Since the rest of the algorithm takes linear time, we obtain the following theorem.

► **Theorem 8.** *Let  $D$  be a set of  $n$  points in  $\mathbb{R}^2$  and  $\text{MINPTS}$  be a given constant. We can compute the HDBSCAN hierarchy on  $D$  for the Euclidean metric with a randomized algorithm in  $O(n \log n)$  expected time.*

## 5 Approximate HDBScan

In this section we introduce an approximate version of HDBSCAN which can be computed in near-linear time in any fixed dimension.

**Approximate DBSCAN\*.** Before we can define approximate HDBSCAN, we need to define approximate DBSCAN\*. Our definition of approximate DBSCAN\* is essentially the same as the definitions of Chen *et al.* [9] and Gan and Tao [13]. The main difference is that we base our definition on DBSCAN\* instead of DBSCAN, which avoids some technical difficulties in the definition.

Let  $\text{MINPTS}$  be a fixed constant. Let  $\mathcal{C}_\varepsilon(D)$  denote the set of clusters in the DBSCAN\* clustering for a given value of  $\varepsilon$ . We call a clustering  $\mathcal{C}_1$  a *refinement* of a clustering  $\mathcal{C}_2$ , denoted by  $\mathcal{C}_1 \prec \mathcal{C}_2$ , when for every cluster  $C_1 \in \mathcal{C}_1$  there is a cluster  $C_2 \in \mathcal{C}_2$  with  $C_1 \subseteq C_2$ . Recall that, as  $\varepsilon$  increases, the DBSCAN\* clusters merge or expand and new singleton clusters may appear, but clusters do not shrink or disappear. Hence, if  $\varepsilon < \varepsilon'$  then<sup>2</sup>  $\mathcal{C}_\varepsilon(D) \prec \mathcal{C}_{\varepsilon'}(D)$ . An approximate DBSCAN\* clustering is now defined as follows.

► **Definition 9.** A  $\delta$ -approximate DBSCAN\* clustering of a data set  $D$ , for given parameters  $\varepsilon$  and  $\text{MINPTS}$ , and a given error  $\delta > 0$ , is defined as a clustering  $\mathcal{C}^*$  of  $D$  into clusters and noise such that  $\mathcal{C}_{(1-\delta)\varepsilon}(D) \prec \mathcal{C}^* \prec \mathcal{C}_\varepsilon(D)$ .

Thus if we choose  $\delta$  sufficiently small, then a  $\delta$ -approximate DBSCAN\* clustering is very similar to the exact DBSCAN\* clustering for the given parameter values.

► **Remark.** An approximate DBSCAN\* clustering can be computed by using the *approximate bichromatic closest pair* algorithm by Arya and Chan [4] as a subroutine in our exact algorithm. The resulting algorithm finds a  $\delta$ -approximate DBSCAN\* clustering in  $\mathbb{R}^d$  in  $O(n \log n + n/\delta^{d/3+c})$  time. This is similar to the running time of Chen *et al.* [9] and the expected running time of Gan and Tao [13], but it has a better dependency on  $\delta$ . Note, however, that Gao and Tao are able to avoid the  $O(n \log n)$  term. The easy details can be found in the full version [7].

<sup>2</sup> Here it is important that we consider DBSCAN\* and not DBSCAN. Indeed, in DBSCAN border points can “flip” between clusters as  $\varepsilon$  increases, and so we do not necessarily have  $\mathcal{C}_\varepsilon(D) \prec \mathcal{C}_{\varepsilon'}(D)$ .

**Approximate HDBScan.** Our definition of an approximate HDBSCAN hierarchy is based on the definition of  $\delta$ -approximate DBSCAN\* clusterings: we say that a hierarchy is a  $\delta$ -approximate HDBSCAN hierarchy if, for any value of  $\varepsilon$ , the clustering extracted from the hierarchy is a  $\delta$ -approximate DBSCAN\* clustering for that value of  $\varepsilon$ . Next we show how to compute a  $\delta$ -approximate HDBSCAN hierarchy in  $O(n \log n)$  time, in any fixed dimension  $d$ .

As in Section 4 we follow the algorithm by Campello *et al.* [8]. Steps 1 and 3 can still be done in  $O(n \log n)$  and  $O(n)$  time, respectively. We speed up Step 2 of the algorithm by computing an MST on a subgraph of the mutual reachability graph  $\mathcal{G}_{\text{mr}}$  rather than on the whole graph. The difference with the exact algorithm of Section 4 is that we will select the edges of the subgraph in a different manner, using ideas from so-called  $\theta$ -graphs [19].

Let  $p \in D$  be a point. We partition  $\mathbb{R}^d$  into simplicial cones with apex  $p$  and whose angular diameter is  $\theta$ , where  $\theta$  will be specified later. (The angular diameter of a cone  $c$  with apex  $p$  is the maximum angle between any two vectors emanating from  $p$  and inside  $c$ .) Let  $\Gamma_p$  be the resulting collection of cones and consider a cone  $c \in \Gamma_p$ . Let  $D(c) \subseteq D$  denote the set of points inside  $c$ . If a point lies on the boundaries of several cones we can assign it to one of these cones arbitrarily. Pick a half-line  $\ell_c$  with endpoint  $p$  that lies inside  $c$ . A  $\theta$ -graph would now be obtained by projecting all points from  $D(c)$  orthogonally onto  $\ell_c$ , and adding an edge from  $p$  to the point closest to  $p$  in this projection, with ties broken arbitrarily. We do the same, except that we add edges to the  $k$  closest points for  $k := 2 \cdot \text{MINPTS} - 3$ . If  $c$  contains fewer than  $k$  points, we simply connect  $p$  to all points in  $D(c)$ . Doing this for all the cones  $c \in \Gamma_p$  gives us a set  $E_p$  of  $O(k/\theta^{d-1}) = O(1/\theta^{d-1})$  edges for point  $p$ . Let  $E(\theta) := \bigcup_{p \in D} E_p$ . The set  $E(\theta)$  can be computed by making a straightforward adaptation to the algorithm to compute a  $\theta$ -graph in  $\mathbb{R}^d$  [19, Chapter 5], leading to the following result.

► **Lemma 10.**  $E(\theta)$  has  $O(n/\theta^{d-1})$  edges and can be computed in  $O((n/\theta^{d-1}) \log^{d-1} n)$  time.

The set  $E(\theta)$ , where  $\theta$  is chosen such that  $\cos \theta \geq 1 - \delta$ , defines the subgraph  $\bar{\mathcal{G}}_{\text{mr}}(\delta)$  on which we compute the MST in Step 2. Since  $\cos \theta > 1 - \theta^2/2$ , we have  $\cos \theta \geq 1 - \delta$  when  $\theta := \sqrt{2\delta}$ . Next we show that an MST on  $\bar{\mathcal{G}}_{\text{mr}}(\delta)$  defines a  $\delta$ -approximate HDBSCAN clustering.

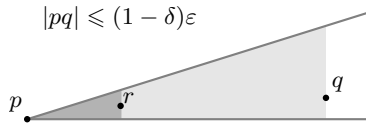
► **Lemma 11.** Let  $\mathcal{T}$  be an MST of  $\bar{\mathcal{G}}_{\text{mr}}(\delta)$  and let  $\varepsilon > 0$ . Let  $\mathcal{C}(\mathcal{T}, \varepsilon)$  be the clustering induced by  $\mathcal{T}$ . Then  $\mathcal{C}$  is a  $\delta$ -approximate DBSCAN\* clustering for the given  $\varepsilon$ .

**Proof.** For a weighted graph  $\mathcal{G}$  and threshold weight  $\tau$ , let  $\mathcal{G}[\tau]$  denote the subgraph obtained by removing all edges of weight greater than  $\tau$ . In order to show that  $\mathcal{C}(\mathcal{T}, \varepsilon) \prec \mathcal{C}_\varepsilon(D)$  we must show that any connected component of  $\mathcal{T}[\varepsilon]$  is contained in a connected component of  $\mathcal{G}_{\text{mr}}[\varepsilon]$ . Since  $\mathcal{T}$  is a subgraph of  $\mathcal{G}_{\text{mr}}$  this is obviously the case.

Next we prove that  $\mathcal{C}_{(1-\delta)\varepsilon}(D) \prec \mathcal{C}(\mathcal{T}, \varepsilon)$ . For this we must prove that any connected component of  $\mathcal{G}_{\text{mr}}[(1-\delta)\varepsilon]$  is contained in a connected component of  $\mathcal{T}[\varepsilon]$ . Since  $\mathcal{T}$  is an MST of  $\bar{\mathcal{G}}_{\text{mr}}(\delta)$ , the connected components of  $\mathcal{T}[\varepsilon]$  are the same as the connected components of  $\bar{\mathcal{G}}_{\text{mr}}(\delta)[\varepsilon]$ . It thus suffices to show the following: for any edge  $(p, q) \in \mathcal{G}_{\text{mr}}[(1-\delta)\varepsilon]$ , there is a path from  $p$  to  $q$  in  $\bar{\mathcal{G}}_{\text{mr}}(\delta)[\varepsilon]$ . We show this by induction on  $|pq|$ , similarly to the way in which it is shown that a  $\theta$ -graph has a small dilation.

Let  $(p, q)$  be an edge in  $\mathcal{G}_{\text{mr}}[(1-\delta)\varepsilon]$ . Consider the set  $\Gamma_p$  of cones with apex  $p$  that was used to define the edge set  $E_p$ , and let  $c \in \Gamma_p$  be the cone containing  $q$ . Recall that we added an edge from  $p$  to the  $k$  points in  $c$  that are closest to  $p$  when projected onto the half-line  $\ell_c$ , where  $k := 2 \cdot \text{MINPTS} - 3$ . Hence, when  $q$  is one of these  $k$  closest points we are done. Otherwise, let  $r \in D(c)$  be the  $(\text{MINPTS} - 1)$ -th closest point.

► **Claim.** (i)  $d_{\text{core}}(r) \leq (1 - \delta)\varepsilon$ , (ii)  $|pr| \leq \varepsilon$ , and (iii)  $|rq| < |pq|$ .



Both the dark and the light grey region contain at least  $\text{MINPTS} - 2$  points, not counting  $p, q, r$ . Depending on the position of  $r$ , all points in the light region or all points in the dark region have distance at most  $(1 - \delta)\varepsilon$  from  $r$ .

■ **Figure 3** Illustration for the proof of Lemma 11.

Before we prove this claim, we first we argue that the claim allows us to finish our inductive proof. Since  $(p, q)$  is an edge in  $\mathcal{G}_{\text{mr}}[(1 - \delta)\varepsilon]$  we have  $d_{\text{mr}}(p, q) \leq (1 - \delta)\varepsilon$ . Thus  $|pq| \leq (1 - \delta)\varepsilon$  and  $d_{\text{core}}(q) \leq (1 - \delta)\varepsilon$ . Together with parts (i) and (iii) of the claim this implies that  $(r, q)$  is an edge in  $\mathcal{G}_{\text{mr}}[(1 - \delta)\varepsilon]$  with  $|rq| < |pq|$ .

In the base case of our inductive proof, where  $(p, q)$  is the shortest edge in  $\mathcal{G}_{\text{mr}}[(1 - \delta)\varepsilon]$ , this cannot occur. Thus  $q$  must be one of the  $k$  closest points in the cone  $c$ , and we have an edge between  $p$  and  $q$  in  $\overline{\mathcal{G}}_{\text{mr}}(\delta)[\varepsilon]$  by construction.

If we are not in the base case, then we have a path from  $r$  to  $q$  in  $\overline{\mathcal{G}}_{\text{mr}}(\delta)[\varepsilon]$  by the induction hypothesis. Moreover,  $(p, r)$  is an edge in  $\overline{\mathcal{G}}_{\text{mr}}(\delta)$  by construction. Since  $|pr| \leq \varepsilon$  by part (ii) of the claim, we have a path from  $p$  to  $q$  in  $\overline{\mathcal{G}}_{\text{mr}}(\delta)[\varepsilon]$ .

It remains to prove the claim. For this we use the following fact [19, Lemma 4.1.4], which is also used to prove that a  $\theta$ -graph has small dilation. Note that although Lemma 4.1.4 in [19] is stated in 2 dimensions, but the proof never assumes that the line on which is projected lives in the same plane and clearly three points  $s, t, p$  live in a single plane.

► **Fact.** *Let  $s, t$  be any two points in a cone  $c \in \Gamma_p$  such that, when projected onto the half-line  $\ell_c$ , the distance from  $p$  to  $s$  is smaller than the distance from  $p$  to  $t$ . Then  $|ps| \leq |pt|/\cos\theta$  and  $|st| < |pt| - (\cos\theta - \sin\theta)|ps| \leq |pt|$ , since we can assume  $\theta$  is sufficiently small that  $\cos\theta - \sin\theta > 0$ .*

Part (iii) of the claim immediately follows from this fact by taking  $s := r$  and  $t := q$ . Part (ii) follows again by taking  $s := r$  and  $t := q$ , using that  $|pq| \leq (1 - \delta)\varepsilon$  and that we have chosen  $\delta$  such that  $\cos\theta = 1 - \delta$ . For part (i) we must prove that there are at least  $\text{MINPTS} - 1$  points within distance  $(1 - \delta)\varepsilon$  from  $r$ . Recall that  $r$  is the  $(\text{MINPTS} - 1)$ -th closest point to  $p$  in the cone  $c$ , measured in the projection onto the half-line  $\ell_c$ . Let  $r_1, \dots, r_k$  be the  $k$  closest points; thus  $r = r_i$  for  $i = \text{MINPTS} - 1$ . We distinguish two cases:  $|pr| \leq (1 - \delta)\varepsilon$  and  $|pr| > (1 - \delta)\varepsilon$ . See also Fig. 3.

In the former case we can conclude that  $|r_i r| \leq (1 - \delta)\varepsilon$  for all  $1 \leq i \leq \text{MINPTS} - 2$  by setting  $s := r_i$  and  $t := r$  and using  $|pr| \leq (1 - \delta)\varepsilon$ . Thus, including the point  $p$ , we know that  $r$  has at least  $\text{MINPTS} - 1$  points within distance  $(1 - \delta)\varepsilon$ .

In the latter case we will argue that  $|r_i r| \leq (1 - \delta)\varepsilon$  for all  $\text{MINPTS} \leq i \leq 2 \cdot \text{MINPTS} - 3$ . Since by part (iii) of the claim we have  $|rq| \leq (1 - \delta)\varepsilon$ , we conclude that also in the latter case  $r$  has at least  $\text{MINPTS} - 1$  points within distance  $(1 - \delta)\varepsilon$ . To argue that  $|r_i r| \leq (1 - \delta)\varepsilon$  we first note that for any point  $s \in c$  we have  $|ss^*| \leq \sin\theta \cdot |ps|$ , where  $s^*$  denotes the orthogonal projection of  $s$  onto  $\ell_c$ . Thus

$$\begin{aligned}
|rr_i| &\leq |rr^*| + |r^*r_i^*| + |r_i r_i^*| \\
&\leq \sin\theta \cdot |pr| + |r^*q^*| + \sin\theta \cdot |pr_i| \\
&\leq 2\sin\theta \cdot |pq|/\cos\theta + |r^*q^*| \\
&= 2\sin\theta \cdot |pq|/\cos\theta + |pq^*| - |pr^*| \\
&\leq 2\sin\theta \cdot |pq|/\cos\theta + |pq| - |pr|\cos\theta \\
&\leq \left(2\frac{\sin\theta}{\cos\theta} + 1 - \cos\theta\right) \cdot (1 - \delta)\varepsilon
\end{aligned}$$

where the last inequality uses  $|pq| \leq (1 - \delta)\varepsilon$  and that we are now considering the case  $|pr| > (1 - \delta)\varepsilon$ . Since we can assume that  $\theta$  is small enough to ensure  $2\sin\theta < \cos^2\theta$ , we conclude that, indeed,  $|rr_i| \leq (1 - \delta)\varepsilon$ . This finishes the proof part (i) of the claim and hence, of the lemma.  $\blacktriangleleft$

Combining the previous two lemmas we obtain the following theorem.

► **Theorem 12.** *Let  $D$  be a set of  $n$  points in  $\mathbb{R}^d$ , and let  $\varepsilon$  and  $\text{MINPTS}$  be given constants. Then, for any given  $\delta > 0$ , we can compute a  $\delta$ -approximate HDBSCAN clustering on  $D$  with respect to  $\varepsilon$  and  $\text{MINPTS}$  for the Euclidean metric in  $O((n/\delta^{(d-1)/2}) \log^{d-1} n)$  time.*

---

## References

- 1 P. Afshani and T.M. Chan. Optimal halfspace range reporting in three dimensions. In *Proc. 20th ACM-SIAM Symp. on Discr. Alg.*, pages 180–186, 2009.
- 2 P.K. Agarwal, H. Edelsbrunner, and O. Schwarzkopf. Euclidean minimum spanning trees and bichromatic closest pairs. *Discr. Comput. Geom.* 6:407–422 (1991).
- 3 M. Ankerst, M.M. Breunig, H.-P. Kriegel, and J. Sander. OPTICS: ordering points to identify the clustering structure. *SIGMOD Rec.* 28:49–60 (1999).
- 4 S. Arya and T.M. Chan. Better  $\varepsilon$ -dependencies for offline approximate nearest-neighbor search, Euclidean minimum spanning trees, and  $\varepsilon$ -kernels. In *Proc. 30th Symp. on Comput. Geom.*, pages 416–425, 2014.
- 5 B. Borah and D. Bhattacharyya. An improved sampling-based DBSCAN for large spatial databases. In *Proc. Int. Conf. on Intelligent Sensing and Inf. Proc.*, pages 92–96, 2004.
- 6 M. de Berg, O. Cheong, M. van Kreveld and M. Overmars. *Computational Geometry: Algorithms and Applications (3rd edition)*. Springer-Verlag, 2008.
- 7 M. de Berg, A. Gunawan, M. Roeloffzen. Faster DB-scan and HDB-scan in Low-Dimensional Euclidean Spaces *CoRR:abs/1702.08607*, 2017.
- 8 R.J.G.B. Campello, D. Malouvi and J. Sander. Density-based clustering based on hierarchical density estimates. In *Proc. 17th Pacific-Asia Conference on Knowledge Discovery and Data Mining*, LNCS 7819, pages 160–172, 2013.
- 9 D.Z. Chen, M.H. Smid and B. Xu. Geometric Algorithms for Density-based Data Clustering. *Int. J. Comput. Geometry Appl.* 15:239–260 (2005)
- 10 T.H. Cormen, C.E. Leiserson, R.L. Rivest and C. Stein. *Introduction to Algorithms* (3rd edition), MIT Press, 2009.
- 11 J. Erickson. On the relative complexities of some geometric problems. In *Proc. 7th Canadian Conf. Comput. Geom. (CCCG)* pages 85–90, 1995.
- 12 M. Ester, H.-P. Kriegel, J. Sander, and X. Xu. A density-based algorithm for discovering clusters in large spatial databases with noise. In *Proc. 2nd Int. Conference on Knowledge Discovery and Data Mining (KDD)*, pages 226–231, 1996.
- 13 J. Gan and Y. Tao. DBSCAN revisited: Mis-claim, un-fixability, and approximation. In *Proc. 2015 ACM SIGMOD Int. Conf. on Management of Data*, pages 519–530.
- 14 J. Gudmundsson, M. Hammer and M. van Kreveld. Higher order Delaunay triangulations. *Computational Geometry: Theory and Applications* 23: 85–98 (2002).
- 15 A. Gunawan. A faster algorithm for DBSCAN. Master’s thesis, TU Eindhoven, March 2013.
- 16 B. Liu. A fast density-based clustering algorithm for large databases. In *Proc. Int. Conf. on Machine Learning and Cybernetics*, pages 996–1000, 2006.
- 17 S. Mahran and K. Mahar. Using grid for accelerating density-based clustering. In *8th Int. Conf. on Computer and Information Technology*, pages 35–40, 2008.
- 18 J. Matoušek. Reporting points in halfspaces. *Computational Geometry: Theory and Applications* 2: 169–186 (1993).

- 19 G. Narasimhan and M. Smid. *Geometric Spanner Networks*. Cambridge Univ. Press, 2007.
- 20 M.M.A. Patwary, D. Palsetia, A. Agrawal, W.-K. Liao, F. Manne, and A. Choudhary. Scalable parallel OPTICS data clustering using graph algorithmic techniques. In *Proc. Int. Conf. on High Perf. Computing, Networking, Storage and Analysis* pages 49:1–49:12, 2013.
- 21 P. Tan, M. Steinbach and V. Kumar. *Introduction to Data Mining*. Addison-Wesley (2006).
- 22 P.M. Vaidya. An  $O(n \log n)$  algorithm for the all-nearest-neighbor problem. *Discr. Comput. Geom.* 4:101–115 (1989).