

Faster DBSCAN and HDBSCAN in Low-Dimensional Euclidean Spaces*

Mark de Berg¹, Ade Gunawan¹, and Marcel Roeloffzen^{2,3}

¹Department of Computing Science, TU Eindhoven, P.O. Box 513, 5600 MB Eindhoven, the Netherlands, mdberg@win.tue.nl

²National Institute of Informatics, Tokyo, Japan, marcel@nii.ac.jp

³JST ERATO, Kawarabayashi Large Graph Project

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 ε (and assuming the second parameter MINPTS is set to a fixed constant, as is the case in practice). Experiments show that the new algorithm is not only fast in theory, but that a slightly simplified version is competitive in practice and much less sensitive to the choice of ε than the original DBSCAN algorithm. 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.

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 [20], 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.* [11] 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, ε 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 apposed to being noise); see Section 2 for a precise definition of the DBSCAN clustering. Typically MINPTS is a constant—in the original article [11] it is concluded that MINPTS = 4 works well—but finding the right value for ε 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 [15]) do so at the expense of computing a slightly different, and not clearly defined, clustering. Others (GridBSCAN [16]) 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 ε of p . Thus $\sum_{p \in D} |N_\varepsilon(p, D)|$ is a lower bound on the running time of the DBSCAN algorithm. In the worst case $\sum_{p \in D} |N_\varepsilon(p, D)| = \Theta(n^2)$,

*MdB is supported by the Netherlands Organization for Scientific Research under grant 024.002.003.

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 to time to report points.) In most practical instances the DBSCAN algorithm is much faster than quadratic. The reason is that ε 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 ε 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 ε will lead to large sets $N_\varepsilon(p, D)$ for points p near the center of the cluster.

Chen *et al.* [8] overcame the quadratic bottleneck of the standard approach, and designed an algorithm¹ with $O(n^{2-\frac{2}{d+2}} \text{polylog } n)$ worst-case running time. Note that for $d = 2$ the running time of the exact algorithm is $O(n^{1.5} \text{polylog } n)$. They also present an approximate algorithm that is more practical. 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 ε , 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 ε . Ankerst *et al.* [3] therefore introduced OPTICS, which can be seen as a hierarchical version of DBSCAN. Recently Campello *et al.* [7] 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 ε . (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 [19], a parallel algorithm that computes a similar (though 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?

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 ε used to define the DBSCAN clustering. (Here, and in our other results, we assume that the parameter MINPTS is a fixed constant. As mentioned this is the case in practice, where one typically uses 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 ε than the original DBSCAN algorithm. Some basic experimental results are provided in Section 6.

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.

Finally, we provide a slightly improved version of the approximate DBSCAN clustering algorithm by Chen *et al.* [8] and by Gan and Tao [12] (their results are discussed in more detail below). Specifically we improve the dependency on the approximation parameter δ . We then extend the concept of an approximate DBSCAN clustering as defined to the hierarchical version. We thus obtain δ -approximate HDBSCAN, an approximate version of the HDBSCAN hierarchy of Campello *et al.* [7], where the parameter δ specifies the accuracy of the approximation. (Intuitively, a δ -approximate HDBSCAN hierarchy has the same clusters as the standard HDBSCAN hierarchy at any level ε , except that

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

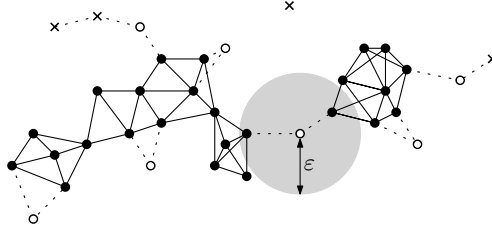


Figure 1: A neighborhood graph (with $\text{MINPTS} = 4$ and ε as indicated). Solid disks are core points, open circles are border points, and crosses are noise. Edges between core points are solid, other edges are dotted. The solid disks and edges form the core graph.

clusters at distance $(1 - \delta) \cdot \varepsilon$ from each other may be merged. See Section 5 for precise definition.) We show that a δ -approximate HDBSCAN hierarchy in \mathbb{R}^d can be computed in $O((n/\delta^{(d-1)/2}) \log n)$ time.

Further related work. This work should be viewed as the journal publication of the (so far unpublished) masters thesis of the second author [14], which contained the results on DBSCAN, extended with results on HDBSCAN. In the meantime, Gan and Tao [12] published a paper in which they extend the work from the masters thesis to \mathbb{R}^d , resulting in an algorithm for DBSCAN with a running time of $O(n^{2 - \frac{2}{\lceil d/2 \rceil + 1} + \gamma})$; 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 [10]. Finally, Gan and Tao show that a δ -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.* [8], 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 show in Section 5 our algorithm can also be used to obtain a deterministic algorithm with $O(n \log n + n/\delta^{d/3})$ running time.

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 (points in the “interior” of a cluster), border points (points on the boundary of a cluster), and noise (points not in any cluster). The distinction is based on two global parameters, ε and 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_{core} , and the set of border points by D_{border} . The remaining points are *noise*. In DBSCAN* [7] border points are not part of a cluster but are considered noise.

Ester *et al.* [11] define the DBSCAN clusters based on the concept of density-reachability (a detailed description is given below). 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 ε ; see Fig. 1. In other words,

$$E = \{ (p, q) \in D \times D : q \in N_\varepsilon(p, D) \setminus \{p\} \}.$$

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 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 ε 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 assigns a border point p to the first cluster that finds p (clusters are constructed one by one); we assign border points to the cluster of their nearest core point.

2.1 The original definition of the DBSCAN clustering.

For comparison purposes only, we restate the original definition of the DBSCAN clustering. Ester *et al.* [11] define when two points are in the same cluster based on the concept of density-reachability, as explained next. A point $q \in D$ is *directly density-reachable* from a point $p \in D$ if $q \in N_\varepsilon(p, D)$ and p is a core point. We denote this by $p \rightarrow q$. A core point is always directly density-reachable from itself, since a point $p \in D$ is always in its own neighborhood. A point q is *density-reachable* from a core point p , denoted by $p \rightsquigarrow q$, if there is a sequence $p = r_0, \dots, r_k = q$ (for some $k \geq 0$) such that $r_0 \rightarrow r_1 \rightarrow \dots \rightarrow r_k$. Observe that if two points p and q are both core points, then $p \rightsquigarrow q$ if and only if $q \rightsquigarrow p$. Two points p and q are *density-connected* if there is a point r such that $r \rightsquigarrow p$ and $r \rightsquigarrow q$.

A cluster is now defined as a subset $C \subseteq D$ such that (i) if a core point p is in C then all points q that are reachable from p are in C , and (ii) any two points in C are density-connected to each other. As observed by Ester *et al.*, a cluster must contain at least one core point, and a cluster is uniquely defined by any of its core points. More precisely, if $p \in C$ is a core point then $C = \{q \in D : p \rightsquigarrow q\}$. Each core point belongs to exactly one DBSCAN cluster. Under the above definition border points can belong to multiple clusters, however. This is typically undesirable, so the original DBSCAN algorithm assigns each border point q to only one cluster, namely the cluster from which q is discovered first by their algorithm. This implies that the computed clustering depends on the order in which their algorithm happens to handle the points.

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}_{box} . Its nodes are disjoint rectangular boxes with a diameter of at most ε that together contain all the points in D , and its edges connect pairs of boxes within distance ε ; 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

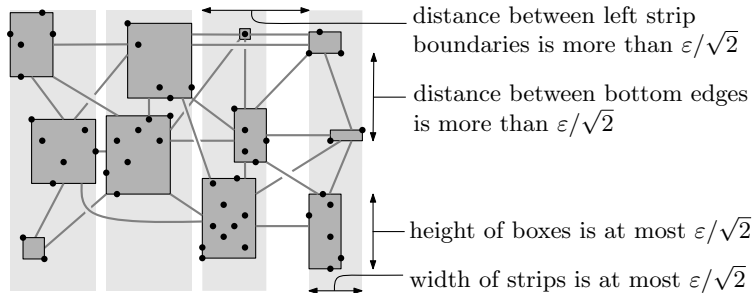


Figure 2: Example of a box graph.

core points when it contains at least 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}_{box} . To compute \mathcal{G}_{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 go through the remaining points from left to right, adding them 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 have 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}_{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 in \mathcal{G}_{box} . Let $N_\varepsilon(b, \mathcal{B})$ be the set of neighbors b .

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

Proof. The box graph obviously has at most n nodes. Next we give a precise analysis of the number of neighbors a box b in some strip \mathcal{B}_j can have. Consider the node corresponding to the box b . If there are two or more strips in between b and some other box b' then $\text{dist}(b, b') > 2(\varepsilon/\sqrt{2}) > \varepsilon$, so b can only have neighbors in \mathcal{B}_{j-2} , \mathcal{B}_{j-1} , \mathcal{B}_j , \mathcal{B}_{j+1} , or \mathcal{B}_{j+2} . We bound the number of neighbors in each of these strips separately.

- The number of neighbors in \mathcal{B}_j is at most four. Indeed, the boxes in \mathcal{B}_j can be ordered vertically, and if there are more than two boxes in between b and some other box $b' \in \mathcal{B}_j$, then the vertical distance between b and b' is more than ε .
- The number of neighbors in \mathcal{B}_{j-1} (and similarly in \mathcal{B}_{j+1}) is at most five. Suppose for a contradiction that b has six or more neighbors in \mathcal{B}_{j-1} . Let b' and b'' be the lowest and highest of these neighbors, respectively. Then the vertical distance between the top boundary of b' and the bottom boundary of b'' is more than $4\varepsilon/\sqrt{2}$. Since the height of b is at most $\varepsilon/\sqrt{2}$, this implies that the vertical distance between b and either b' or b'' is more than $(4\varepsilon/\sqrt{2} - \varepsilon/\sqrt{2})/2 > \varepsilon$, contradicting that both b' and b'' are neighbors of b .
- The number of neighbors in \mathcal{B}_{j-2} (and similarly in \mathcal{B}_{j+2}) is at most four. Suppose for a contradiction that b has five or more neighbors in \mathcal{B}_{j-2} . Let b' and b'' be the lowest and highest of these neighbors, respectively. Then the vertical distance between the top boundary of b' and the bottom boundary of b'' is more than $3\varepsilon/\sqrt{2}$. Since the height of b is at most $\varepsilon/\sqrt{2}$, this implies the vertical distance between b and either b' or b'' is more than $(3\varepsilon/\sqrt{2} - \varepsilon/\sqrt{2})/2 = \varepsilon/\sqrt{2}$. The horizontal distance is at least $\varepsilon/\sqrt{2}$ as well, because there is a strip in between. Hence, the total distance between b and either b' or b'' is more than ε , contradicting that both b' and b'' are neighbors of b .

Adding up the maximum number of neighbors in each of the strips gives us a maximum of 22 neighbors in total. \square

This also gives us an easy way to compute the edge set E_{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}_{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 ε . Hence, if a box contains more than MINPTS points, then all of them are core points. Thus the following simple algorithm suffices to determine 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 $O(\sum_{b' \in N_\varepsilon(b, \mathcal{B})} n_{b'})$ time for each point $p \in b$.

Lemma 3 *Given \mathcal{G}_{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

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

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

Step 3: Compute the cluster cores. In the previous step, we determined the core points. Next we wish to determine the clusters or, more precisely 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}_{box} . To this end, we first remove certain boxes and edges from \mathcal{G}_{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 [9]. After computing the connected components we 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 $|pp'| \leq \varepsilon$. If $n_b < \text{MINPTS}$ or $n_{b'} < \text{MINPTS}$ then we simply check all pairs of core points. We can argue as in the proof of Lemma 3 that this takes $O(\text{MINPTS} \cdot n) = O(n)$ time in total. If $n_b \geq \text{MINPTS}$ and $n_{b'} \geq \text{MINPTS}$ we have to be more careful, since checking all pairs of points can lead to a quadratic running time. For example, if both cells contain $n/2$ points, then checking all pairs would take $\Omega(n^2)$ time. When $n_b \geq \text{MINPTS}$ and $n_{b'} \geq \text{MINPTS}$ we therefore compute the Delaunay triangulation of $D(b) \cup D(b')$ in $O((n_b + n_{b'}) \log(n_b + n_{b'}))$ time [6], and check whether it has an edge (p, p') of length at most ε between points $p \in D(b)$, $p' \in D(b')$. This is sufficient because the pair of points defining the closest distance between $D(b)$ and $D(b')$ must define an edge in the Delaunay triangulation. 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 $|pp'| \leq \varepsilon$. As mentioned, the total time spent on pairs with $n_b < \text{MINPTS}$ or $n_{b'} < \text{MINPTS}$ is $O(n)$. Let \mathcal{B}^* be the set of boxes containing at least MINPTS points. Then the total time spent on the pairs of boxes from \mathcal{B}^* is

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

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

Remark. In practice we can also use a brute-force algorithm when $n_b \geq \text{MINPTS}$ and $n_{b'} \geq \text{MINPTS}$, because the number of points in boxes with more than MINPTS points is typically still not very large. Moreover, if both b and b' contain many points, then there are often many pairs of points within distance ε from each other, and we can stop when we find such a pair.

Step 4: Assigning border points to clusters. It remains to decide for non-core points p whether p is a border point or noise. (For DBSCAN* we can skip Step 4, since in DBSCAN* border points are considered 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 $b \in \mathcal{B}$ 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 1 *Let D be a set of n points in \mathbb{R}^2 , and ε and MINPTS be given constants. Then we can compute a DBSCAN clustering on D according to ε and 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 [12]. For completeness we describe the extension and

the resulting bound. One trivial modification is that in \mathbb{R}^d we need to use boxes in \mathcal{G}_{box} of width at most ε/\sqrt{d} along each axis to ensure their diameter is at most ε . The only other difference is in Step 3, where we have to decide for all pairs b, b' of neighboring boxes whether there are core points $p \in D(b)$ and $p' \in D(b')$ with $|pp'| \leq \varepsilon$. When $n_b \geq \text{MINPTS}$ and $n_{b'} \geq \text{MINPTS}$ we can no longer use the Delaunay triangulation for this, as it may have quadratic size. Instead we can use a known algorithm for bichromatic closest pair [2], which gives a running time of $O(n^{2 - \frac{2}{\lceil d/2 \rceil + 1} + \gamma})$, where $\gamma > 0$ is an arbitrarily small constant. (For $d = 3$ the n^γ factor can be replaced by a polylogarithmic factor.)

4 A fast algorithm for HDBSCAN in the plane

Campello *et al.* [7] 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* (for a fixed MINPTS) that arise as ε increases from $\varepsilon = 0$ to $\varepsilon = \infty$. Initially, when $\varepsilon = 0$, all points are noise. As ε 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 ε 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 ε such that $\varepsilon = |pq|$ for some pair of points $p, q \in D$. Also note that several events may happen simultaneously, not only when two pairwise distances are equal, but also when an event triggers other events. 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 ε at which the corresponding cluster was created, is the HDBSCAN hierarchy. Notice that we can easily extract the DBSCAN* clustering for any desired value of ε (with respect to the given MINPTS) from the dendrogram in linear time. 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}_{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}_{mr} . (Indeed, the cluster-growing process corresponds to the process of computing an MST on \mathcal{G}_{mr} using Kruskal's algorithm [9].) 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}_{mr} .
3. Convert \mathcal{T} into a dendrogram where each internal node stores the value of ε 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 ℓ -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 [21]. 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}_{mr} when computing \mathcal{T} . To overcome this problem we need the following generalization of Delaunay triangulations, introduced by Gudmundsson *et al.* [13]. 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 [13]. (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 6 *Let $\bar{\mathcal{G}}_{\text{mr}}$ be the subgraph of \mathcal{G}_{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}_{mr} .*

Proof. Imagine computing an MST \mathcal{T} on \mathcal{G}_{mr} using Kruskal's algorithm [9]. This algorithm treats the edges (p, q) of \mathcal{G}_{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 ε —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}_{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 p, q is a diametrical pair of C_{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 $|N_\varepsilon(r, D)| \geq |D(C_{pq})| + 2 = k + 3 = \text{MINPTS}$. 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, as shown in Figure 3. By induction, there is a path from p to r in $\bar{\mathcal{G}}_{\text{mr}}$ 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. \square

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 6 implies that after computing the core distances and the k -OD edges in $O(n \log n)$ time—recall that $k = \max(\text{MINPTS} - 3, 0) = O(1)$ —we can compute the MST for \mathcal{G}_{mr} by considering only $O(n)$ edges. Thus computing the MST can be done in $O(n \log n)$ time [9]. Since the rest of the algorithm takes linear time, we obtain the following theorem.

Theorem 2 *Let D be a set of n points in \mathbb{R}^2 and 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.*

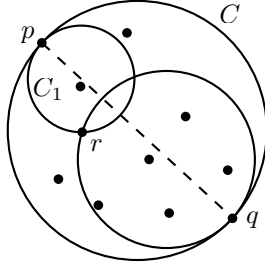


Figure 3: Illustration of the recursive argument in the proof of Lemma 6.

5 Approximate DBSCAN* and HDBSCAN

The approach from the previous section for computing the HDBSCAN hierarchy will not give a subquadratic bound in higher dimensions, since the number of Delaunay edges can be quadratic. The running time of the algorithm from Section 3 to compute a single DBSCAN clustering is subquadratic in any fixed dimension, but the running time quickly deteriorates as the dimension increases. In this section we introduce an approximate version of DBSCAN* and HDBSCAN, both of which can be computed in linear time. An approximation for DBSCAN* that runs in expected linear time was already provided by Chen *et al.* [8] Gan and Tao [12], however our approach runs has a slightly better dependency on the approximation factor δ . To the best of our knowledge this is the first linear time approximation algorithm for *hdbscan*.

5.1 Approximate DBSCAN*

First we define what exactly is an approximate DBSCAN* clustering. Our definition of approximate DBSCAN* is essentially the same as the definitions by Chen *et al.* [8] and Gan and Tao [12]. The main difference is that we base our definition on DBSCAN* instead of DBSCAN, which avoids some technical difficulties in the definition.

Let MINPTS be a fixed constant. Let $\mathcal{C}_\varepsilon(D)$ denote the set of clusters in the DBSCAN* clustering for a given value of ε . 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 ε increases, the DBSCAN* clusters merge or expand and new singleton clusters may appear, but clusters do not shrink or disappear. Hence, if $\varepsilon < \varepsilon'$ then² $\mathcal{C}_\varepsilon(D) \prec \mathcal{C}_{\varepsilon'}(D)$. An approximate DBSCAN* clustering is now defined as follows.

Definition 1 A δ -approximate DBSCAN* clustering of a data set D , for given parameters ε and MINPTS, and a given error $\delta > 0$, is now 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 δ sufficiently small, then a δ -approximate DBSCAN* clustering is very similar to the exact DBSCAN* clustering for the given parameter values.

The algorithm. As mentioned in the introduction, both Chen *et al.* [8] and Gan and Tao [12] already presented algorithms for this. We obtain a slightly better dependency on δ than Gan and Tao [12] by plugging in a better algorithm for approximate bichromatic closest pair.

Recall that the bottleneck in computing a DBSCAN* clustering lies in checking, for pairs b, b' of neighboring boxes, whether there is a pair $(p, p') \in D(b) \times D(b')$ with $|pp'| \leq \varepsilon$. We can perform this check approximately by computing an *approximate bichromatic closest pair* $(p, p') \in D(b) \times D(b')$ such that $|pp'| \leq (1 + \alpha) \cdot \text{dist}(D(b), D(b'))$ for $\alpha = \delta/(1 - \delta)$, where $\text{dist}(D(b), D(b'))$ denotes the

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

distance between the points of the true closest pair. This can be done in $O((1/\alpha)^{d/3}(n_b + n_{b'})) = O((1/\delta)^{d/3}(n_b + n_{b'}))$ time [4]. We add the edge (b, b') to the box graph \mathcal{G}_{box} when $|pp'| \leq \varepsilon$. This way we can obtain the following result.

Theorem 3 *Let D be a set of n points in \mathbb{R}^d , and let ε and MINPTS be given constants. Then, for any given $\delta > 0$, we can compute a δ -approximate DBSCAN* clustering on D with respect to ε and MINPTS for the Euclidean metric in $O(n \log n + (1/\delta)^{d/3}n)$ time.*

Proof. Let \mathcal{C} be the clustering computed using the approximate bichromatic closest pairs. When $\text{dist}(D(b), D(b')) \leq \varepsilon$, then the reported approximate bichromatic closest pair (p, p') has $|pp'| \leq \varepsilon$, so the set of edges added to the box graph is a subset of the edge set of the actual box graph at the given value of ε . Hence, $\mathcal{C} \prec \mathcal{C}_\varepsilon(D)$. On the other hand, when $\text{dist}(D(b), D(b')) \leq (1 - \delta)\varepsilon$ then the approximate closest pair (p, p') has $|pp'| \leq (1 + \alpha)(1 - \delta)\varepsilon = \varepsilon$. Hence, we are guaranteed to add all edges of the box graph for $(1 - \delta)\varepsilon$ and so $\mathcal{C}_{(1-\delta)\varepsilon}(D) \prec \mathcal{C}$. Thus \mathcal{C} is a δ -approximate DBSCAN* clustering.

Recall that the running time of our algorithm is $O(n \log n)$, plus the time needed to compute the edges of the box graph. In the approximate version the latter step takes time

$$\sum_{b \in \mathcal{B}^*} \sum_{b' \in N_\varepsilon(b, \mathcal{B}^*)} O((1/\delta)^{d/3} \cdot (n_b + n_{b'})) = O((1/\delta)^{d/3}n).$$

□

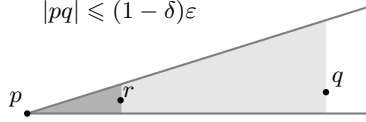
Approximate HDBSCAN. Our definition of an approximate HDBSCAN hierarchy is based on the definition of δ -approximate DBSCAN* clusterings: we say that a hierarchy is a δ -approximate HDBSCAN hierarchy if, for any value of ε , the clustering extracted from the hierarchy is a δ -approximate DBSCAN* clustering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

As in Section 4 we follow the algorithm by Campello *et al.* [7], and we speed up Step 2 of the algorithm by computing an MST on a subgraph of the mutual reachability graph \mathcal{G}_{mr} rather than on the whole graph. (Steps 1 and 3 can still be done in $O(n \log n)$ and $O(n)$ time, respectively.) 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 θ -graphs [18].

Let $p \in D$ be a point. We partition \mathbb{R}^d into simplicial cones with apex p and whose angular diameter is θ , where θ 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 Γ_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 ℓ_c with endpoint p that lies inside c . A θ -graph would now be obtained by projecting all points from $D(c)$ orthogonally onto ℓ_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 $O(k/\theta) = O(1/\theta)$ 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 θ -graph in \mathbb{R}^d [18, Chapter 5], leading to the following result.

Lemma 7 *$E(\theta)$ 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 θ 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 δ -approximate HDBSCAN clustering.



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 4: Illustration for the proof of Lemma 8.

Lemma 8 *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 δ -approximate DBSCAN* clustering for the given ε .*

Proof. For a weighted graph \mathcal{G} and threshold weight τ , let $\mathcal{G}[\tau]$ denote the subgraph obtained by removing all edges of weight greater than τ . 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}_{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 θ -graph has a small dilation.

Let (p, q) be an edge in $\mathcal{G}_{\text{mr}}[(1-\delta)\varepsilon]$. Consider the set Γ_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 ℓ_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) $|rq| < |pq|$.

Before we prove this claim, we first 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 $\bar{\mathcal{G}}_{\text{mr}}(\delta)[\varepsilon]$ by construction.

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

It remains to prove the claim. For this we use the following fact [18, Lemma 4.1.4], which is also used to prove that a θ -graph has small dilation.

Fact: Let s, t be any two points in a cone $c \in \Gamma_p$ such that, when projected onto the half-line ℓ_c , the distance from p to s is smaller than the distance from p to t . Then $|ps| \leq |pt|/\cos\theta$ and $|st| < |pt|$.

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 δ 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 ℓ_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. 4.

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 ℓ_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 θ 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 for part (i) of the claim and hence, of the lemma. \square

Combining the previous two lemmas we obtain the following theorem.

Theorem 4 *Let D be a set of n points in \mathbb{R}^d , and let ε and MINPTS be given constants. Then, for any given $\delta > 0$, we can compute a δ -approximate HDBSCAN clustering on D with respect to ε and MINPTS for the Euclidean metric in $O((n/\delta^{(d-1)/2}) \log^{d-1} n)$ time.*

6 Experimental evaluation

In this section we experimentally investigate the efficiency of our new algorithm and compare it to the original algorithm. The only goal of these experiments is to serve as a proof of concept to illustrate that indeed for very basic point distributions the original algorithm has a bad running time, whereas the new algorithm performs much better. We first describe some implementation details and we discuss our implementation of the original algorithm. We then describe the data sets and parameters for the tests. Finally we present the results and conclusions.

Implementation details. We implemented two versions of our new algorithm: the strip-based approach for Step 1 and the grid-based approach. In Step 3, where we have to decide for all pairs b, b' of neighboring boxes whether there are core points $p \in D(b)$ and $p' \in D(b')$ with $|pp'| \leq \varepsilon$, we use the following randomized brute-force approach. (This is instead of the theoretically better Delaunay triangulations or spherical emptiness queries.) Without loss of generality assume $n_b \leq n_{b'}$. For each core point $p \in D(b)$ we first test if $\text{dist}(p, b') \leq \varepsilon$. If not, no point of $D(b')$ can be within distance ε of p . If so, we test each point $p' \in D(b')$ to see if $|pp'| \leq \varepsilon$. If during this procedure we find two core points within distance ε from each other, we can stop. The randomization is obtained by considering the points in each box in random order; it ensures that if there are many pairs within distance ε , we expect to find such a pair early.

The original DBSCAN algorithm performs a spherical range query for each point $p \in D$ to find all $q \in D$ with $|pq| \leq \varepsilon$. To this end an indexing structure such as an R-tree is typically used. In our implementation we use the box-graph to answer the queries. Note that an R-tree also groups the points into boxes at the leaf level; the tree structure is then used find the leaf boxes intersecting the query range, after which the points inside these boxes are tested. The boxes in our box-graph can be seen as being optimized for the radius ε of the query range, and so the box-graph should be at least as efficient as a general-purpose R-tree.

Experimental set-up. We ran the algorithms on several synthetic data sets in 2D and 4D, each consisting of four clusters. (For other numbers of clusters the results are similar.) The clusters either have a uniform or Gaussian distribution, and their centers are placed roughly 700 units apart in a hypercube with edge length 1,000 units. Uniform clusters are generated within a ball of radius 300 around the cluster center, Gaussian clusters are generated with a standard deviation of 100. For several data sets we added 5% noise to the input, uniformly inside a slightly expanded bounding box around the clusters.

Parameters. We analyse the efficiency of the algorithms with respect to two parameters: the input size and the density within the clusters. As a measure of the density we use $n(r/\varepsilon)^d$, where r is the cluster radius; thus, for the uniform distribution, the density represents the expected number of points within distance ε from a core point, and also the expected number of points in the boxes of the box-graph within the clusters.

Measurements. We compare the algorithms in two different ways: we measure the actual execution times, and the number of pairs of points for which a distance computation is done. The latter is the main operation for finding the clusters in both the new and the original algorithm, so it provides a good implementation-independent measure. For the original algorithm we also count the sum of neighborhood sizes of all points. Following earlier work, we call this the number of *seeds*. This is a lower bound for the number of operations needed in the original algorithm, independent of the indexing structure used to find the neighborhoods.

Result for fixed input size. In these experiments we fix the input size and run the algorithms with different values of ε . In 2D we ran the algorithms on a data set of each type—uniform or Gaussian and with or without noise—in which the clusters contain 500,000 points each. In 4D we use the same types of data sets, but with 200,000 points per cluster. The results are shown in Fig. 5.

Result for fixed density. In our second set of experiments we use data sets of different sizes, where for each data set we pick ε such that the density (as defined above) remains constant. For each type of data we generated data sets ranging from 20,000 to 500,000 points per cluster in 2D, and from 20,000 to 300,000 points per cluster for 4D. We ran these experiments for two different values of ε : one that is roughly the smallest value needed to find the four clusters, and one that is roughly the highest value for which the four clusters do not start to merge. The results are given in Fig. 6.

Discussion. The running times in Fig. 5 (top row) show a lot of fluctuation, while the results on the number of distance computations are very stable. We suspect the fluctuation in running time is related to memory issues, or to other processes claiming resources. Hence, most of our conclusions will rely on the number of distance calculations and the sum of neighborhood-sizes.

The results from Fig. 5 clearly show the different dependency on density between the two algorithms. For the original algorithm the number of distance calculations (shown in the second row) grows linearly with the density. The sum of the neighborhood sizes (seeds) is slightly over 50% of the number of distance computations, so the linear dependency is inherent in the original algorithm, which reports the complete neighborhood of all points. The new algorithm, on the other hand, is insensitive to the density. Note that the bottom row in Fig. 5 shows that there is a fairly large range of values of ε at which the correct clusters are found. (The smallest range occurs for Gaussian clusters with noise, and even there the four clusters are still found for ε^2 roughly between 4 and 20.) The original algorithm can compete with the new algorithm only when ε is chosen very near the lower end of the correct range of ε ; for a “safe” value of ε in the middle of the range, our algorithm is significantly faster. This is important, as determining the right value of ε is hard and may require running the algorithm several times with different values of ε . Moreover, different clusters may have different densities. Choosing a value for ε that is near the low end of the range

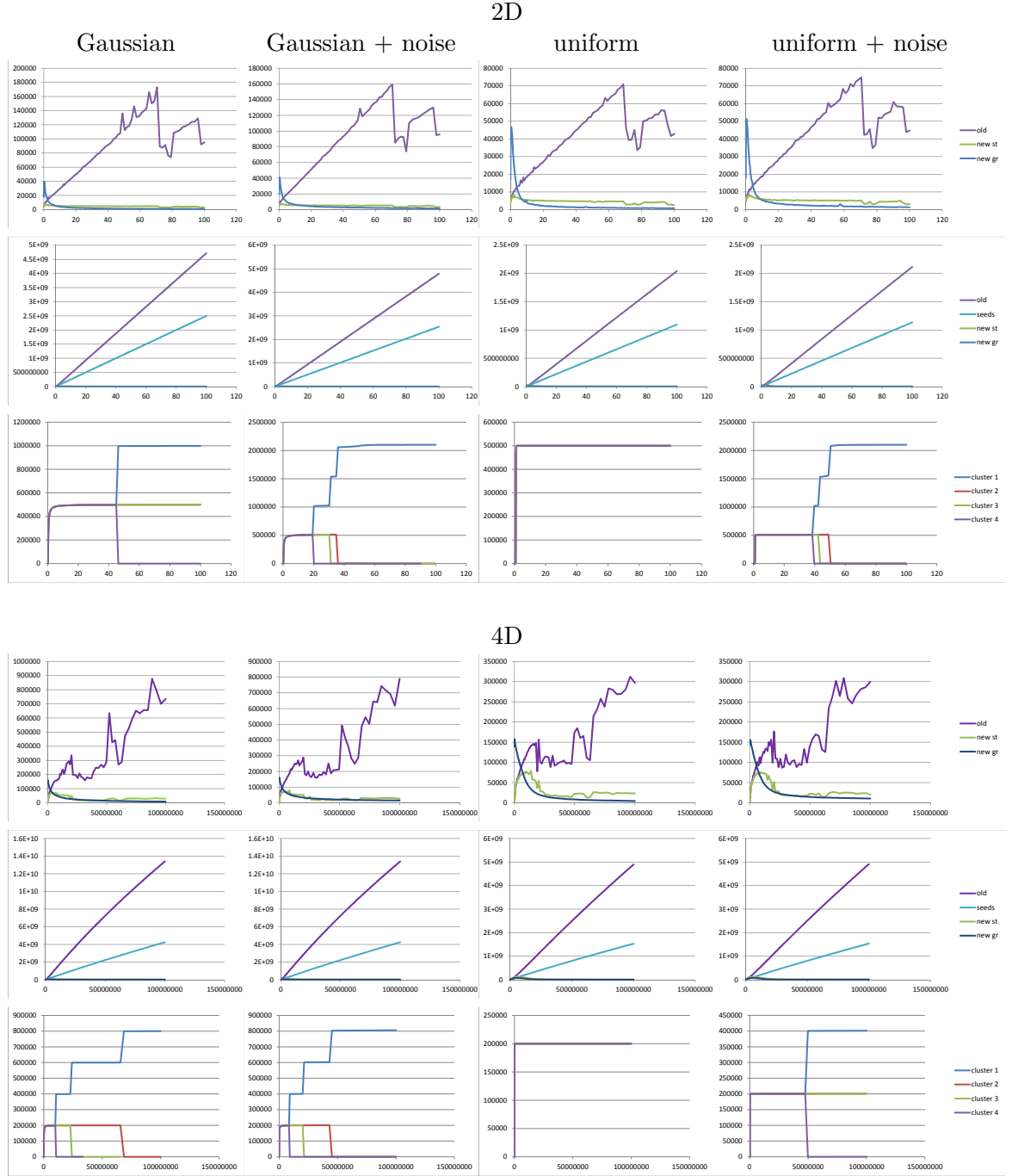


Figure 5: Results when varying ε . The x -axis denotes ε^2 and the y -axis denotes (from top to bottom row for both 2D and 4D) processing time in ms, number of distance calculations, and the number of points in the four largest clusters. The third row thus show for which values of ε the “correct” clustering (with four clusters of 500,000 points for 2D and 200,000 points in 4D) is found. The new method with a strip-based or grid-based construction is denoted by *new st* and *new gr* respectively. Note that in the second row these lines overlap.

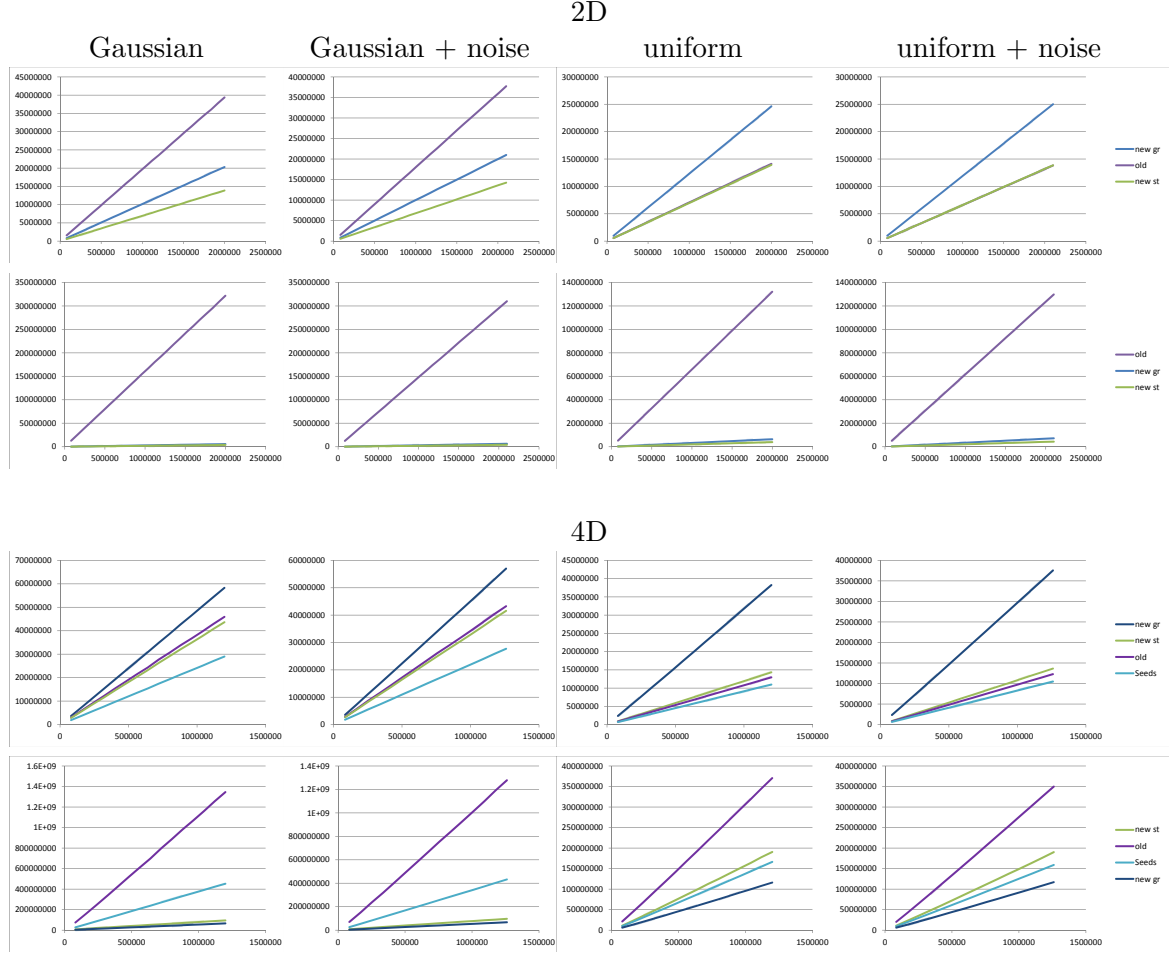


Figure 6: Results of running the algorithms on data-sets of varying sizes, but with a fixed density. The first and third rows are from experiments with a low density and the second and fourth rows with a high density. The x -axis denotes the total number of points in the point set and the y -axis the number of distance calculations. The new method with a strip-based or grid-based construction is denoted by *new st* and *new gr* respectively. Note that in the high-density experiments (second and fourth row), the lines for the two new methods almost overlap. However, in the low-density setting, the line for the old method and strip-construction (*new st*) overlap.

for some clusters may then either fail to find some other clusters or it may lead to an explosion in running time for the other clusters.

Fig. 6 shows that the strip-based approach tends to make fewer comparison than the grid-based approach, especially in the low-density setting shown in the upper row in the figure. (In Fig. 5 the low-density setting is in the far left in each graph, and therefore less visible.) This is caused by two things: the boxes are created in a data-driven way in the strip-based approach, and the boxes are smaller (being bounding boxes of the points in it). The former means that points within distance ε end up in the same box more often (in which case they are not compared), while the latter means that some pairs of bounding boxes can be at distance more than ε , while the “corresponding” grid cells have distance just below ε . Except for the low-density setting, the actual computation time (see the top row in Fig. 5) of the grid-based approach is better in 2D (even though the number of comparisons is not), due to the smaller constant factors in the approach.

7 Concluding remarks

We presented a new algorithm for DBSCAN in \mathbb{R}^2 , which runs in $O(n \log n)$ time in the worst case. We also presented an $O(n \log n)$ algorithm for HDBSCAN in \mathbb{R}^2 —this is the first subquadratic algorithm for this problem—and we presented near-linear algorithms for approximate HDBSCAN. It will be interesting to do a more thorough experimental evaluation of our new algorithms. Specifically to compare the approach to other implementations that are current available. It would also be interesting to implement the approximation algorithms and compare them to the exact ones. From the theoretical side, a main open problem is to see if we can compute the exact HDBSCAN hierarchy in subquadratic time in dimensions $d \geq 3$.

References

- [1] P. Afshani and T.M. Chan. Optimal halfspace range reporting in three dimensions. In *Proc. 20th ACM-SIAM Symposium on Discrete Algorithms*, 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 ε -dependencies for offline approximate nearest-neighbor search, Euclidean minimum spanning trees, and ε -kernels. In *Proc. 30th Symposium on Computational Geometry*, 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 Information Processing*, 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] 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.
- [8] 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)

- [9] T.H. Cormen, C.E. Leiserson, R.L. Rivest and C. Stein. *Introduction to Algorithms* (3rd edition), MIT Press, 2009.
- [10] J. Erickson. On the relative complexities of some geometric problems. In *Proc. 7th Canadian Conf. Comput. Geom. (CCCG)* pages 85–90, 1995.
- [11] 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.
- [12] 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.
- [13] J. Gudmundsson, M. Hammer and M. van Kreveld. Higher order Delaunay triangulations. *Computational Geometry: Theory and Applications* 23: 85–98 (2002).
- [14] A. Gunawan. A faster algorithm for DBSCAN. Master’s thesis, Technische University Eindhoven, March 2013.
- [15] B. Liu. A fast density-based clustering algorithm for large databases. In *Proc. Int. Conf. on Machine Learning and Cybernetics*, pages 996–1000, 2006.
- [16] 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.
- [17] J. Matoušek. Reporting points in halfspaces. *Computational Geometry: Theory and Applications* 2: 169–186 (1993).
- [18] G. Narasimhan and M. Smid. *Geometric Spanner Networks*. Cambridge University Press, 2007.
- [19] 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 Performance Computing, Networking, Storage and Analysis* pages 49:1–49:12, 2013.
- [20] P. Tan, M. Steinbach and V. Kumar. *Introduction to Data Mining*. Addison-Wesley (2006).
- [21] P.M. Vaidya. An $O(n \log n)$ algorithm for the all-nearest-neighbor problem. *Discr. Comput. Geom.* 4:101–115 (1989).