

# RECOME: a New Density-Based Clustering Algorithm Using Relative KNN Kernel Density

Yangli-ao Geng<sup>a</sup>, Qingyong Li<sup>a</sup>, Rong Zheng<sup>b</sup>, Fuzhen Zhuang<sup>c,d</sup>, Ruisi He<sup>e</sup>, Naixue Xiong<sup>f</sup>

<sup>a</sup>*School of Computer and Information Technology, Beijing Jiaotong University, Beijing, China. {gengyla,liqy}@bjtu.edu.cn*

<sup>b</sup>*Department of Computing and Software, McMaster University, Hamilton, Canada. rzheng@mcmaster.ca*

<sup>c</sup>*Key Lab of Intelligent Information Processing of Chinese Academy of Sciences (CAS), ICT, CAS, Beijing, 100190, China.*

<sup>d</sup>*University of Chinese Academy of Sciences, Beijing 100049, China. zhuangfz@ics.ict.ac.cn*

<sup>e</sup>*State Key Laboratory of Rail Traffic Control and Safety, Beijing Jiaotong University, Beijing, China. ruisi.he@bjtu.edu.cn*

<sup>f</sup>*Department of Mathematics and Computer Science Northeastern State University, Tahlequah, USA. xiongnaixue@gmail.com*

---

## Abstract

Discovering clusters from a dataset with different shapes, densities, and scales is a known challenging problem in data clustering. In this paper, we propose the RElative COre MErge (RECOME) clustering algorithm. The core of RECOME is a novel density measure, i.e., Relative  $K$  nearest Neighbor Kernel Density (RNKD). RECOME identifies core objects with unit RNKD, and partitions non-core objects into atom clusters by successively following higher-density neighbor relations toward core objects. Core objects and their corresponding atom clusters are then merged through  $\alpha$ -reachable paths on a KNN graph. We discover that the number of clusters computed by RECOME is a step function of the  $\alpha$  parameter with jump discontinuity on a small collection of values. A fast jump discontinuity discovery (FJDD) method is proposed based on graph theory. RECOME is evaluated on both synthetic datasets and real datasets. Experimental results indicate that RECOME is able to discover clusters with different shapes, densities, and scales. It outperforms six baseline methods on both synthetic datasets and real datasets. Moreover, FJDD is shown to be effective to extract the jump discontinuity set of parameter  $\alpha$  for all tested datasets, which can ease the task of data exploration and parameter tuning.

*Keywords:* density-based clustering; density estimation;  $K$  nearest neighbors; graph theory.

---

## 1. Introduction

Clustering, also known as unsupervised learning, is a process of discovery and exploration for investigating inherent and hidden structures within

a large dataset [10]. It has been extensively applied to a variety of tasks [17, 32, 11, 45, 21, 18, 47, 30, 46, 41, 20]. Many clustering algorithms have been proposed in different scientific disciplines [13], and these methods often differ in the selection of objective functions, probabilistic models or heuristics adopted. Nonetheless, two difficulties, how to choose appropriate clustering number and how to discover clusters of an arbitrary shape, are faced by most methods. Density-based clustering approaches are characterized by aggregating mechanisms based on density [28]. They can handle data with irregular shapes and determine clustering number automatically. Ester *et al.* [9] and Sander *et al.* [35] pioneered two density-based methods, Density Based Spatial Clustering of Applications with Noise (DBSCAN) and Generalizing DBSCAN, to detect clusters in a spatial database according to density differences. Although both methods can detect clusters with different shapes, they face the challenge of choosing appropriate parameter values. Subsequently, many improved methods have been proposed [3, 24, 7, 26, 29]. Recently, a novel density based clustering method, named Fast search-and-find of Density Peaks (FDP) [33], was proposed. This algorithm assumes that cluster centers are surrounded by neighbors with lower local density and that they are at a relatively large distance from any point with higher density. FDP can recognize clusters regardless of their shape and of the dimensionality of the space in which they are embedded, but it lacks an efficient quantitative criterion for judging cluster centers. Accordingly, approaches such as 3DC [23] and STClu [42] have been proposed to improve FDP.

Density-based clustering methods have the advantages of discovering clusters with arbitrary shapes and dealing with noisy data, but they face two challenges. First, traditional density measures are not adaptive to clusters with different densities. Second, performances of traditional methods (e.g., DBSCAN and FDP) are sensitive to parameters, and it is non-trivial to set these parameters properly for different datasets.

Aiming to address these challenges, we propose the RElative COre MERge (RECOME) clustering algorithm, which is based on two density measures: the  $K$  nearest Neighbor Kernel Density (NKD) and Relative  $K$  nearest Neighbor Kernel Density (RNKD). RECOME firstly identifies *core objects* corresponding to objects with RNKD equal 1. A core object and its descendants, which are defined by a directed relation (i.e., *higher density nearest-neighbor*) based on the NKD, form an *atom cluster*. These atom clusters are then merged using a novel notion of  $\alpha$ -connectivity on a KNN graph. RECOME has been evaluated using both synthetic datasets and real world datasets. Experimental results demonstrate that RECOME outperforms six baseline methods. Furthermore, we find that the clustering results of RECOME can be characterized by a step function of its parameter  $\alpha$ , and therefore devise a fast jump discontinuity discovery (FJDD) algorithm to extract the small collection of jump discontinuity values. In summary, this work makes the following contributions.

1. We give a formal analysis showing that the density measure NKD enjoys some desirable properties. Furthermore, based on the NKD, we propose a

new density measure RNKD, which is instrumental in detecting clusters with different densities.

2. RECOME can avoid the “decision graph fraud” problem [23] of FDP and can handle clusters with different shapes, densities, and scales. Furthermore, RECOME has nearly linear computational complexity if the  $K$  nearest neighbors of each object are computed in advance.
3. FJDD can extract all jump discontinuity values of parameter  $\alpha$  for any dataset in  $\mathcal{O}(n \log n)$  time, where  $n$  is the number of objects. It will greatly benefit parameter selection in real-world applications.

This paper is organized as follows. Section 2 introduces the related work. Section 3 presents the new density measure RNKD and discusses the robustness of NKD and RNKD. Section 4 describes the proposed clustering method RECOME. Section 5 presents the auxiliary algorithm *FJDD*. Section 6 demonstrates experimental results. Finally, we conclude the paper in Section 7.

## 2. Related Work

Existing clustering methods can be categorized into partitional methods, hierarchical methods, grid-based methods, graph-based methods, density-based methods, etc [10]. Partitional methods such as K-means [27] and K-medoids [16], divide data to a number of partitions and a certain quantitative measure of the “goodness” of the resulting clusters is maximized iteratively. Hierarchical clustering methods can be agglomerative (bottom-up) or divisive (top-down). An agglomerative clustering (e.g., AGNES [14]) starts with one object for each cluster and recursively merges two or more of the most appropriate clusters. A divisive clustering (e.g., DIANA [15]) starts with the dataset as one cluster and recursively splits the most appropriate cluster. The process continues until a stopping criterion is reached. Grid-based methods such as STING [43] and CLIQUE [1], divide the original data space into grids, and then group the grids according to the statistical characters of objects in each grid. Graph-based methods, such as SCAN [44] and spectral clustering [37], first construct a similarity graph from a dataset, and then utilize the notion of structural-context similarity or the eigenvalues of Laplacian matrix to generate clusters. Density-based methods (e.g., DBSCAN [9] and DENCLUE [12]) first estimate the distribution density of objects in a feature space, and then recognize clusters as regions of high density separated by regions of lower density. In this paper, we focus on density-based methods because they are highly relevant to the proposed algorithm.

In [9], Ester et al. proposed the first density-based method DBSCAN. In DBSCAN, a *cut-off density* of an object  $o$  is defined as the number of objects falling inside a ball of radius  $\epsilon$  centered at  $o$ . If the cut-off density of  $o$  is higher than a threshold,  $MinPts$ ,  $o$  is regarded as a key object. When the distance between two key objects is less than  $\epsilon$ , they are called density-reachable. Density-reachable key objects form basic clusters. A non-key object is assigned to a basic cluster if it is within  $\epsilon$  distance to a key object in the respective

cluster; otherwise, the non-key object is treated as noise. DBSCAN is sensitive to the choice of parameters  $\epsilon$  and  $MinPts$ , and can hardly handle clusters with heterogeneous densities. To overcome these drawbacks, Ankerst *et al.* [3] proposed an enhanced density-connected algorithm OPTICS. OPTICS provides a visual tool to help users find the cluster structure and determine the parameters. Although OPTICS reduces the subjectivity in a parameter estimation, when dealing with a complex dataset, it is also difficult to determine how many  $\epsilon$ 's are needed to find potential clusters [7].

Kernel density [39] is a well-known alternative to cut-off density. It is continuous and less sensitive to parameter selection. DENCLUE [12] is a method based on kernel density, in which the local peaks (i.e., local density maxima) of the kernel density function are used to define clusters. Then, each object is assigned to a cluster by a hill-climbing procedure. However, traditional kernel density methods tend to give biased estimation when handling clusters with different scales. To overcome this difficulty, KNN kernel density [25] has been introduced in the KNN-kernel density-based Clustering (KNNC) [40]. In our work, the proposed RNKD estimation is inspired by KNN kernel density with further improvement allowing the inclusion of low-density clusters.

Detecting clusters with heterogeneous densities is another challenge for density-based approaches [19]. Some algorithms [8] [31] have been proposed to solve this problem. In [8], a Shared Nearest Neighbor (SNN) clustering algorithm was proposed. SNN first finds  $K$  nearest neighbors of each data object according to similarity, and then refines the similarity between pairs using the number of neighbors that the two objects share. Based on the new measure, a DBSCAN-like process is used to generate clusters. SNN has been shown the capacity in detecting clusters with complex distribution, whereas its excessive parameters and relatively high computational cost weaken its applicability in practice. DiscovErIng c-clusters of different dEnsities (DECODE) [31] assumes that the target dataset is generated by a series of point processes and tries to find clusters as connected regions of objects whose distances to their  $m$ -th nearest neighbor are similar. It has shown to be capable of determining the number of density types with little prior knowledge, but the high computational cost limits its application to large data. In [5], a density-based outlier detection approach was proposed. It relies on the local outlier factor (LOF) of each object, which is equal to the average of the ratios between the local density of an object and those of its  $K$  nearest neighbors. LOF can effectively distinguish outliers from normal clusters. However, it is not suitable for finding clusters with complex distribution. In this work, we introduce the novel density measure RNKD, which allows handling clusters with heterogeneous densities efficiently.

Rodriguez and Laio proposed a novel density-based clustering method by finding density peaks called FDP [33]. FDP discovers clusters by a two-phase process. First, local density is computed for each object according to the number of objects in its  $d_c$  neighborhood, and then a group decision method is applied to determine cluster centers, called density peaks. Second, remaining objects are assigned to the same cluster as its nearest neighbor with higher density. FDP is effective in finding clusters with different shapes. However, reasonable

cluster centers are hard to determine when several density peaks exist in a cluster. In this work, RECOME adopts an agglomerative procedure to merge atom clusters (analogous to those resulting from “density peaks”), which is feasible even encountering clusters with multi-peaks.

Key notations used in the paper are listed in Table 1.

Table 1: Main notations used throughout the paper.

Notation	Description
$ \cdot $	Absolute value of a scalar or cardinality of a set.
$V$	Dataset. Lower case symbols $u, v, w, u_i, v_i, w_i$ denote elements of $V$ .
$d(\cdot, \cdot)$	The distance function on $V$ . In the formal analysis, we assume it is a metric (i.e., satisfying non-negativity, symmetry and triangle inequality).
$N(u)$	The $K$ nearest neighbors set of $u$ in $V$ with respect to $d(\cdot, \cdot)$ .
$N_i(u)$	The $i$ -th nearest neighbor of $u$ in $V$ .
$d_i(u)$	The distance between $u$ and $N_i(u)$ , i.e., $d(u, N_i(u))$ .

### 3. Relative KNN Kernel density estimation

Though many density measures have been proposed, few considers the local relative density levels (will be shown in Section 3.1), which are crucial to detect clusters with various densities. In this section, we will first illustrate the disadvantage of classical density measures used in density clustering, and then introduce the proposed measure RNKD that homogenizes the density estimation across clusters with different densities.

#### 3.1. Density Estimation

The most commonly used density measure is cut-off density, which is defined as the number of objects in an  $\epsilon$ -ball centered at the respective object. However, it is highly sensitive to the parameter  $\epsilon$ . As shown in Figure 1a, a small variation in  $\epsilon$  can result in drastic differences in density estimation. Another classical measure is kernel density defined as

$$\bar{\rho}(u) = \sum_{v \in V} \text{ker} \left( \frac{d(u, v)}{h} \right),$$

where  $\text{ker}(\cdot)$  is usually a monotonically decreasing and continuous function, and  $h$  is a constant controlling the scale. Kernel density is continuous and less sensitive to parameter selection, but it tends to give biased estimation for objects in a small-size cluster because it considers contributions of all objects in the dataset (see Figure 1b).

Our proposed RNKD estimation is inspired by KNN kernel density [40], and only considers the objects in  $N(u)$  for the density estimation of object

$u$ . In addition, the density of an object should be positive and has a negative relation with the distances between itself and its neighbors. Thus, the K-nearest neighbor kernel density (NKD) of object  $u$  is defined as

$$\rho(u) = \theta \sum_{v \in N(u)} \exp\left(-\frac{d(u, v)}{\sigma}\right), \quad (1)$$

where  $\sigma = \frac{\sum_{u \in V} d_K(u)}{|V|}$  is the mean of the distance between  $v$  and its  $K$ -th nearest neighbors, and  $\theta$  is a normalizing factor.

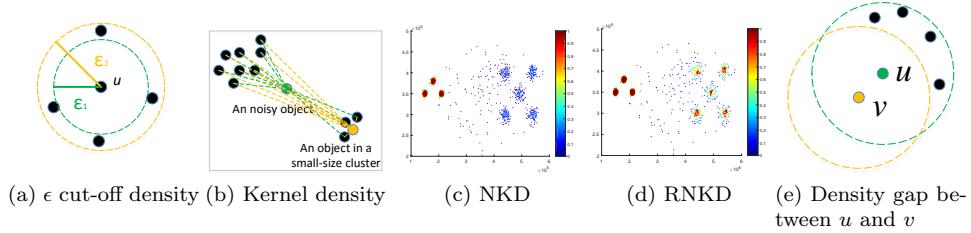


Figure 1: Comparison of different density measures. (a) Sensitivity to  $\epsilon$  in DBSCAN. (b) Non-parametric kernel density methods are sensitive to noise. (c) NKD estimates may mistake low-density clusters as noise. (d) RNKD allows discrimination of low-density clusters from noise. The temperature of data points in (c) and (d) indicates the NKD and RNKD, respectively.

NKD exhibits some good properties (as will be discussed in Section 3.2) and allows easy discrimination of outliers. However, it may mistake low-density clusters as outliers. One such example is shown in Figure 1c, where low-density clusters on the right may be mistaken as noise. The key observation is that though the NKD of a low-density cluster is low compared to high-density clusters, its NKD is still comparatively high in its surroundings. This observation motivates us to introduce a new density measure, called RNKD as follows.

**Definition 1.** Given  $u \in V$ , RNKD of  $u$ , denoted by  $\rho^*(u)$ , is defined as

$$\rho^*(u) = \frac{\rho(u)}{\max_{v \in N(u) \cup \{u\}} \{\rho(v)\}}. \quad (2)$$

From the definition, we can see RNKD is a normalized local relative density measure. Figure 1d shows that it can effectively accentuate both dense and sparse clusters, but diminish noises. Though the parameter  $K$  still needs to be tuned (the number of nearest neighbors), we will show that it is generally robust to  $K$  in Section 3.2.

### 3.2. Properties of NKD and RNKD

As discussed in Section 3.1, one shortcoming of cut-off density is its sensitivity to parameter  $\epsilon$ . Additionally, even for a fixed  $\epsilon$ , there can be significant

discontinuity in the cut-off densities of two adjacent objects (see Figure 1e), which will lead to undesirable results. Though non-parametric kernel density can avoid this gap, it may suffer from the problem illustrated in Figure 1b. Next, we show that NKD exhibits local continuity. In particular, the ratio of  $\rho(u)$  and  $\rho(v)$  between neighbor points  $u$  and  $v$  are bounded on both sides.

**Theorem 1.**  $\forall u, v \in V, u \neq v$ , we have  $\frac{|\rho(u) - \rho(v)|}{\rho(u) + \rho(v)} < 1 - \exp\left(-\frac{d(u,v)}{\sigma}\right)$ .

*Proof.* See appendix.  $\square$

Consequently, for any  $u, v \in V$ , the following results hold

$$\frac{\exp\left(-\frac{d(u,v)}{\sigma}\right)}{2 - \exp\left(-\frac{d(u,v)}{\sigma}\right)} \leq \frac{\rho(u)}{\rho(v)} \leq \frac{2 - \exp\left(-\frac{d(u,v)}{\sigma}\right)}{\exp\left(-\frac{d(u,v)}{\sigma}\right)}. \quad (3)$$

This implies that  $\rho(u)$  and  $\rho(v)$  are close to each other when  $\frac{d(u,v)}{\sigma}$  is small. In other words, NKD changes slowly in the interior of a cluster. Furthermore, (3) also shows that RNKD will be stable at a high level in the interior of a cluster.

Another important issue is how sensitive the density measures are to the choice of parameter  $K$ . Figure 2 shows the heat-map and the distribution of NKD and RNKD on the synthetic dataset in Figure 1 when  $K$  valued from  $0.4\sqrt{|V|}$  to  $\sqrt{|V|}$ . From this figure, we can observe that, when  $K \geq 0.5\sqrt{|V|}$ , both NKD and RNKD follow similar distributions despite varying of  $K$ . This indicates, for this representative dataset, NKD and RNKD are robust to the parameter  $K$  in the range  $[0.5\sqrt{|V|}, \sqrt{|V|}]$ . In Section 6.3, we will show that the proposed RECOME algorithm also achieves good performance when  $K$  is in this range.

#### 4. RECOME Clustering Algorithm

Now we are in the position to present RECOME based on NKD and RNKD. We first identify *core objects* corresponding to data points of peak relative density. These core objects serve as centers of sub-clusters, called *atom clusters*, which will be further merged through connected paths on a KNN graph. Thus, RECOME is, in essence, an agglomerative hierarchical clustering method.

##### 4.1. Finding core objects and atom clusters

From the definition of RNKD, we know that there exist objects with unit RNKD. Such objects are good candidates for cluster centers as they have local maxima values in density. Formally,

**Definition 2.** An object  $u$ ,  $u \in V$ , is called a **core object** if  $\rho^*(u) = 1$ . Denote the set of core objects by  $O = \{u | u \in V, \rho^*(u) = 1\}$ .

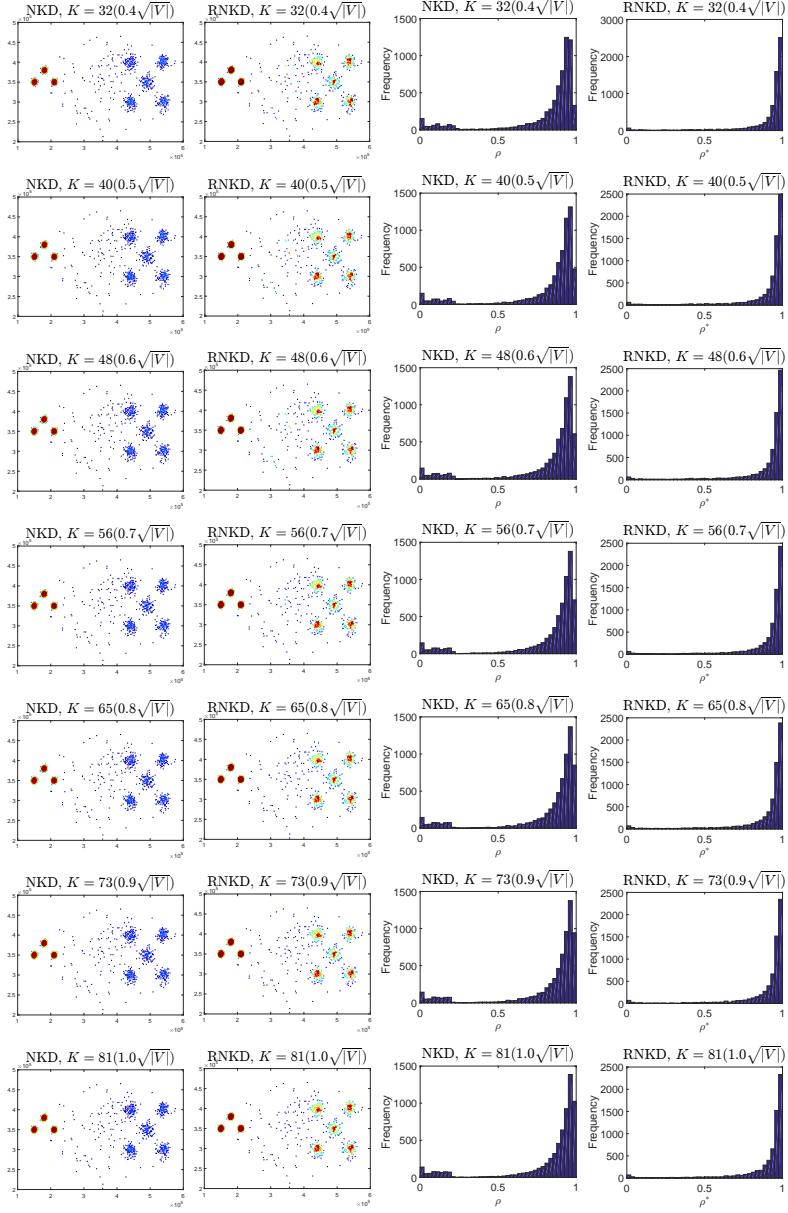


Figure 2: Columns from left to right are heat-map of NKD, heat-map of RNKD, distribution of NKD, and distribution of RNKD, respectively. Rows from top to bottom show the situations when  $K$  takes the value from small to large.

For a non-core object  $u$  (namely,  $\rho^*(u) < 1$ ), we define its *Higher Density*

*Nearest-neighbor* (HDN),  $\pi(u)$  as

$$\pi(u) = \arg \min_{\rho(v) > \rho(u), v \in V} \{d(u, v)\}. \quad (4)$$

By the definition of  $\rho^*(u)$ , we can see that  $\pi(u)$  exists in  $N(u)$ . In other words, the distance between an object and its HDN is small enough for most objects. So, HDN can be seen as a discrete approximation of gradients, which are hard to compute directly.

HDN allows us to construct a directed graph  $\mathcal{G} = (V, A)$ , where each vertex is an object and a direct edge exists from a non-core object to its higher density nearest neighbor, namely,  $A = \{\langle u, \pi(u) \rangle | u \in V \setminus O\}$ .

In  $\mathcal{G}$ , starting from any non-core object and following the directed edges, we will eventually reach a core object. In other words,  $\mathcal{G}$  can be partitioned into trees with disjoint vertices, where each tree is rooted at a core object. A core object and its descendants thus form a cluster, called an *atom cluster*. Due to the bijective relation between an atom cluster and its core object, we use the two terms interchangeably. Formally, for a core object  $o \in O$ , the atom cluster rooted at  $o$  is given by,

$$\{o\} \cup \{v \in V \setminus O | v \text{ is connected to } o \text{ in } \mathcal{G}\}. \quad (5)$$

Atom clusters form the basis of final clusters, however, they themselves tend to be too fine. A true cluster may consist of many atom clusters. This happens when many local maximals exist in one true cluster. Thus, a merging step is needed to selectively combine atom clusters into desirable clusters.

#### 4.2. Merging Atom Clusters

To merge the selected atom clusters to obtain a better clustering result, we treat each core object as the representative object of the atom cluster that it belongs to. The problem is thus transformed from merging atom clusters to merging core objects. To do so, we define the KNN graph as follows,

**Definition 3.** A **KNN graph**  $G_K = (V, E)$ , is an undirected graph, where  $V$  consists of all objects in the given dataset and  $E = \{\langle u, v \rangle | u \in N(v) \wedge v \in N(u)\}$ .

Therefore, two objects  $u$  and  $v$  are directly connected in graph  $G_K$  if they are the top  $K$  nearest neighbors of each other. Next, we define the notion of  $\alpha$ -reachability between vertices in  $G_K$ .

**Definition 4.** Given  $\alpha \in [0, 1]$  and  $u, v \in V$ , if there exists a path  $\langle u, w_1, \dots, w_s, v \rangle$  in  $G_K$  that satisfies  $\rho^*(w_i) > \alpha, i = 1, 2, \dots, s$ , then  $v$  is  **$\alpha$ -reachable from  $u$** , denoted by  $u \xrightarrow{\alpha} v$ .

Clearly,  $\alpha$ -reachability is reflexive, symmetric, and transitive for the core object set  $O$  in  $G_K$ . Thus, it can be used as an equivalence relation among the core objects to divide them into equivalence classes. All atom clusters associated with core objects in the same equivalent class are merged into a single cluster.

Alternatively, we can view  $\alpha$ -reachability as a way to prune edges in  $G_K$ . Core objects in the same equivalence class reside in the same connected component.

The choice of  $\alpha$  is expected to affect the partition of equivalence classes. When  $\alpha$  is small, say  $\alpha = 0$ , most core objects are merged together. On the other hand, when  $\alpha$  is large, in the extreme case when  $\alpha = 1$ , every core object forms an equivalence class and atom clusters are the final clusters. The property and the selection of  $\alpha$  will be analyzed in Section 6.3.

#### 4.3. Algorithm Description and Complexity Analysis

---

**Algorithm 1:** RECOME Clustering

---

```

Input:  $V, K, \alpha$ 
Output: Clusters
1 for  $u \in V$  do
2   | Compute the KNN set  $N(u)$  of  $u$ ;
3   |  $\rho(u) = \theta \sum_{v \in N(u)} \exp\left(-\frac{d(u,v)}{\sigma}\right)$ ;
4 end
5 for  $u \in V$  do
6   |  $\rho^*(u) = \frac{\rho(u)}{\max_{v \in N(u) \cup \{u\}} \{\rho(v)\}}$ ;
7 end
8 Let  $O = \{u | u \in V, \rho^*(u) = 1\}$ ;
9 for  $u \in V \setminus O$  do
10  |  $\pi(u) = \arg \min_{\rho(v) > \rho(u), v \in V} \{d(u, v)\}$ ;
11 end
12 Construct the directed graph  $\mathcal{G} = (V, A)$ , where  $A = \{\langle u, \pi(u) \rangle | u \in V \setminus O\}$ ;
13 for  $o \in O$  do
14  | Find the atom cluster  $C_o$  w.r.t.  $o$  by
    |  $C_o = \{o\} \cup \{u \in V \setminus O | u \text{ is connected to } o \text{ in } \mathcal{G}\}$ ;
15 end
16 Construct the undirected KNN graph  $G_K = (V, E)$ , where
    |  $E = \{\langle u, v \rangle | u \in N(v) \wedge v \in N(u)\}$ ;
17 Merge the core objects and their corresponding atom clusters by
    |  $\alpha$ -reachable relation in  $G_K$  (see Definition 4). Then get the final
        clustering result;

```

---

To this end, we have presented the four main steps of the proposed RECOME clustering algorithm, i.e., computing NKD, calculating RNKD, discovering atom clusters, and merging atom clusters. Algorithm 1 summarizes the details of RECOME.

Let the number of objects be  $n = |V|$ . Computing the KNN set is the most time-consuming step with computational complexity of  $\mathcal{O}(n^2 + Kn \log n)$ . Note that this step can be performed in parallel and accelerated using indexing

structures such as kd-tree [6] and R\*-tree [4]. Both  $\rho^*(u)$  and  $\pi(u)$  can be obtained from its  $K$  nearest neighbors in  $\mathcal{O}(K)$  time. After obtaining the KNN set  $N(u)$  and  $\rho^*(u)$  for each  $u$ , the KNN graph  $G_K$  can be built by a linear scan of  $N(u)$ . Merging core objects can be accomplished using depth-first-search on  $G_K$ . The computation complexity of Algorithm 1 is thus  $\mathcal{O}(n^2 + Kn \log n)$ . Specifically, when  $K$  is fixed, the KNN sets need to be computed only once and thus the time complexity will reduce to  $\mathcal{O}(Kn)$  for a different  $\alpha$ . The space complexity of RECOME is  $\mathcal{O}(Kn)$  because it needs not to store the distance matrix.

## 5. Fast Jump Discontinuity Discovery for Parameter $\alpha$

As evident from the description of RECOME, parameter  $\alpha$  is crucial to the clustering result for a fixed  $K$ . Take the example in Figure 3. Starting from the same set of core objects and atom clusters, different values of  $\alpha$  may result in different numbers of clusters. In particular, as  $\alpha$  increases, cluster granularity (i.e., the volume of clusters) decreases and cluster purity increases. Thus, a pertinent question is how to select a proper  $\alpha$ . For ease of presentation, we suppose parameter  $K$  is fixed in this section.

### 5.1. Problem Formalization

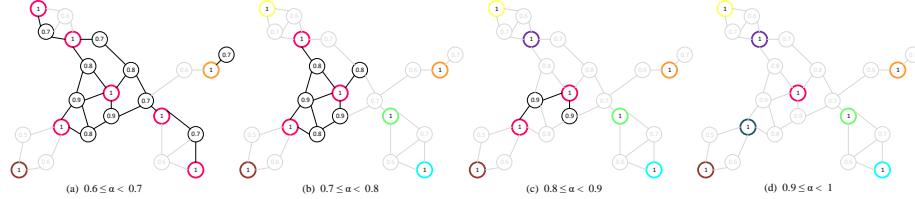


Figure 3: Example of merging core objects results for different values of  $\alpha$ . KNN graph is shown and each object (in circle) is marked with its relative density. Core objects of the same color belong to the same final cluster.

We observe that though  $\alpha$  varies continuously in the interval of  $[0, 1]$ , the number of clusters computed by RECOME is a step function of  $\alpha$ . For instance, in Figure 3, only five clustering outcomes are possible for the example dataset, corresponding to  $\alpha$  in the ranges of  $[0, 0.6)$ ,  $[0.6, 0.7)$ ,  $[0.7, 0.8)$ ,  $[0.8, 0.9)$ , and  $[0.9, 1)$ . The numbers of resulting clusters are 1, 3, 6, 7, and 8, respectively. It is desirable to have a small collection of  $\alpha$  values (or ranges) that affect the clustering result as the processes of parameter tuning by developers or parameter selection by domain experts can be simplified.

We formalize the above intuition by first introducing the notion of jump discontinuity set.

**Definition 5.** Given a data set  $V$  and an input parameter  $K$ , an ascending list  $L = \{\alpha_1, \alpha_2, \dots, \alpha_l\}$  is called a **jump discontinuity** (JD) set if the number of

resulting clusters from RECOME,  $\#(V, K, \alpha)$  is a step function of  $\alpha$  with jump discontinuity at  $\alpha_1, \alpha_2, \dots, \alpha_l \in [0, 1]$  from left to right.

Obviously,  $\#(V, K, \alpha)$  is a non-decreasing function of  $\alpha$ . By definition, each JD in  $L$  yields a unique clustering result. From all the JDs in  $L$ , we can produce all possible clusters using RECOME. Recall that  $O$  is the set of core objects and  $|O|$  is the maximum number of clusters attainable by RECOME. Trivially,  $|L| \leq |O|$ .

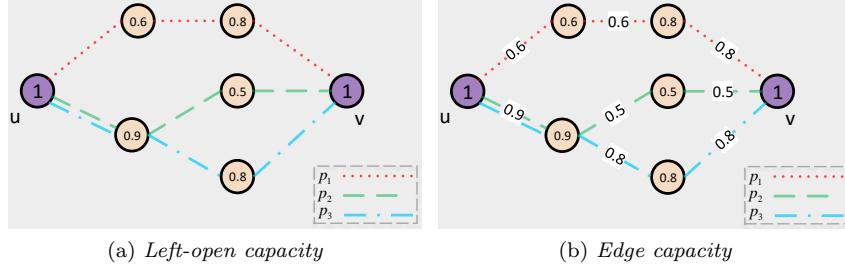


Figure 4: (a) *Left-open capacity*. Each object (in circle) is marked with its relative density. There are two core objects (i.e.,  $u$  and  $v$ ) and three paths (i.e.,  $p_1$ ,  $p_2$ , and  $p_3$ ) between them. It can be seen that  $c(p_1) = 0.6$ ,  $c(p_2) = 0.5$ , and  $c(p_3) = 0.8$ , so we have  $c(u, v) = \max\{0.6, 0.5, 0.8\} = 0.8$ . (b) *Edge capacity*. Each edge is marked with its weight. There are two core objects (i.e.,  $u$  and  $v$ ) and three paths (i.e.,  $p_1$ ,  $p_2$ , and  $p_3$ ) between them. It can be seen that  $c^*(p_1) = 0.6$ ,  $c^*(p_2) = 0.5$ , and  $c^*(p_3) = 0.8$ , so we have  $c^*(u, v) = \max\{0.6, 0.5, 0.8\} = 0.8$ .

As described in Section 4, the clustering number is only related to the merging of core objects, so we only need to discuss the effect of  $\alpha$  on the merging process. Recall that  $G_K$  is the KNN graph of dataset  $V$ . Without loss of generality, we suppose  $G_K$  is connected. Consider a path  $p = \langle u \equiv w_0, w_1, w_2, \dots, w_s \equiv v \rangle$ . Its *left-open path capacity* is defined as  $c(p) = \min_{i>0} \rho^*(w_i)$ . Suppose there are  $l$  paths from  $u$  to  $v$ , denoted by  $p_1, p_2, \dots, p_L$ . The *left-open capacity* between  $u$  and  $v$  is defined as  $c(u, v) = \max_{1 \leq l \leq L} c(p_l)$ . An example is shown in Figure 4a. In other words,  $v$  is  $\alpha$ -reachable from  $u$  iff  $\alpha < c(u, v)$ . Furthermore, we have the following proposition.

**Proposition 1.** *The JD set  $L$  equals  $\{0\} \cup \{c(u, v) | u, v \in O, u \neq v\}$ .*

*Proof.* See appendix. □

Considering that 0 belongs to  $L$  trivially, in the following section, we will focus on finding out  $L^* \triangleq \{c(u, v) | u, v \in O, u \neq v\}$ .

### 5.2. Fast Jump Discontinuity Discovery

Finding  $L^*$  can be reduced to the problem *All-Pairs Bottleneck Paths* in vertex weighted graphs [36]. Unfortunately, the state-of-the-art method for it has a time complexity of  $\mathcal{O}(n^{2.575})$ . To reduce the time complexity, we need to exploit the characteristics of our problem.

To transform our problem, we define the edge weight function  $w(\langle u, v \rangle) = \min\{\rho^*(u), \rho^*(v)\}$  for  $\langle u, v \rangle \in E(G_K)$ . For a path  $p = \langle u \equiv w_0, w_1, w_2, \dots, w_s \equiv v \rangle$ , its *path edge-capacity* is defined as  $c^*(p) = \min_{i>0} w(\langle w_{i-1}, w_i \rangle)$ . Suppose that there are  $l$  paths from  $u$  to  $v$ , denoted by  $p_1, p_2, \dots, p_L$ . The *edge capacity* between  $u$  and  $v$  is defined as  $c^*(u, v) = \max_{1 \leq l \leq L} c^*(p_l)$ . An example is shown in Figure 4b.

**Proposition 2.**  $\forall u, v \in O, u \neq v$ , the left-open capacity between  $u, v$  equals the edge capacity  $c^*(u, v)$  in the transformed graph, i.e.,  $c(u, v) = c^*(u, v)$ .

Proposition 2 implies that  $L^* \triangleq \{c(u, v) | u, v \in O, u \neq v\} = \{c^*(u, v) | u, v \in O, u \neq v\}$ . Therefore, it suffices to determine  $\{c^*(u, v) | u, v \in O, u \neq v\}$ . This problem can be reduced to the problem of *All-Pairs Bottleneck Paths* in edge weighted graphs (edge-APBP) [36]. For edge-APBP, the following property holds:

**Theorem 2.** Suppose  $G$  is a connected edge weighted graph and  $T$  is the **maximum spanning tree** of  $G$ . Then,  $\forall u, v \in V(G) = V(T), u \neq v$ ,  $c_G^*(u, v) = c_T^*(u, v)$ .

*Proof.* See [36]. □

$c_G^*(u, v)$  and  $c_T^*(u, v)$  denote the *edge capacity* between  $u$  and  $v$  in  $G$  and  $T$ , respectively. Suppose the maximum spanning tree of  $G_K$  is  $T_K$ , and then according to Theorem 2, we can extract  $L^*$  from  $T_K$  using Algorithm 2.

---

**Algorithm 2:** Extracting-JD

---

```

Input:  $T_K$ 
Output:  $L^*$ 
1 Sort  $E(T_K)$  by edge weight decreasingly and suppose the result is
    $\{e_1, e_2, \dots, e_{n-1}\}$ ;
2 Set  $L^* \leftarrow \emptyset$ ,  $T_K \leftarrow T_K - E(T_K)$  (remove all edges from  $T_K$ );
3 for  $i$  from 1 to  $n - 1$  do
4   | Suppose  $u$  and  $v$  are the two ends of  $e_i$ ;
5   | Suppose  $C_u$  and  $C_v$  are the two components in  $T_K$  that contain  $u$  and
       $v$ , respectively;
6   | if both  $C_u$  and  $C_v$  contain at least one core object then
7   |   |  $L^* \leftarrow L^* \cup \{w(e_i)\}$ ;
8   | end
9   |  $T_K \leftarrow T_K + e_i$ ;
10 end

```

---

In the Algorithm 2, the sorting step is most time consuming with computational complexity  $\mathcal{O}(n \log n)$ . The Prim's maximum spanning tree algorithm has a time complexity of  $\mathcal{O}(|E| + n \log n)$ , where  $|E| = \mathcal{O}(Kn)$ . Therefore, the computation complexity of Algorithm 3 is  $\mathcal{O}(Kn + n \log n)$ . Finally, the Fast Jump Discontinuity Discovery algorithm is summarized in Algorithm 3.

---

**Algorithm 3:** Fast Jump Discontinuity Discovery

---

**Input:** KNN graph  $G_K$  and core object set  $O$

**Output:** Sorted JD list  $L$

- 1 Define the edge weight function  $w(\langle u, v \rangle) = \min\{\rho^*(u), \rho^*(v)\}$  for  $\langle u, v \rangle \in E(G_K)$ ;
  - 2 Compute the maximum spanning tree of  $G_K$  using Prim's algorithm and denote it as  $T_K$ ;
  - 3 Set  $L \leftarrow \text{Extracting-JD}(T_K) \cup \{0\}$ ;
- 

## 6. Experiments

In this section, we evaluate RECOME over synthetic and real-world datasets, and compare it with six other representative algorithms. We introduce the experiment setup in Section 6.1, and then present the experimental results and analysis in Section 6.2. In section 6.3, we conduct the parameter analysis of RECOME. Our implementation uses Microsoft Visual C++ 2015 14.0.24720.00, and all experiments are conducted on a workstation (Windows 64 bit, 4 Intel 3.2 GHz processors, 4 GB of RAM). Our code has been released in: <https://github.com/gyla1993/RECOME-A-new-density-based-clustering-algorithm-using-relative-KNN-kernel-density>.

### 6.1. Experiment Setup

#### 6.1.1. Datasets

*Two-dimensional synthetic datasets:* Four representative datasets S1, S2, S3, and S4 are used. S1 comprises 5,000 objects and 15 Gaussian clusters. S2 is an unbalanced dataset, which contains 8 classes of different density and size, 6500 objects with 117 noisy objects injected. S3 is comprised of 6 classes of non-convex shape and 8000 objects; S4 is the mixture of S2 and S3, which contains 14 classes of different shapes, densities, and scales. See Figure 5 for details.

*MNIST training set* [22]: It is a real-world dataset containing 60,000 examples of handwritten digits from 0 to 9, which has been widely used in data mining and machine learning. We select the subsets of different digits for clustering tasks, including M367 (all examples of digits “3, 6, 7”), M3467, M0:8 (all examples of digits from 0 to 8), M0:9, sM3467 (small edition of M3467 for the visualization, generated by randomly selecting 900 examples of each digit), and sM0:8.

*Olivetti face database* [34]: It is a real-world dataset containing 40 subjects of human face and each contains 10 examples. This dataset poses a challenge to the density estimation since the real number of clusters is comparable with the number of objects in each cluster (10 different pictures for each people). Two subsets, face10 (the examples of first 10 subjects) and face40 (the examples of all 40 subjects), have been selected for the evaluation.

#### 6.1.2. Baseline Methods and Settings

We select 6 representative density-based clustering algorithms as baseline methods, i.e., DBSCAN [9], SNN [8], KNNC [40], FDP [33], 3DC [23], STClu

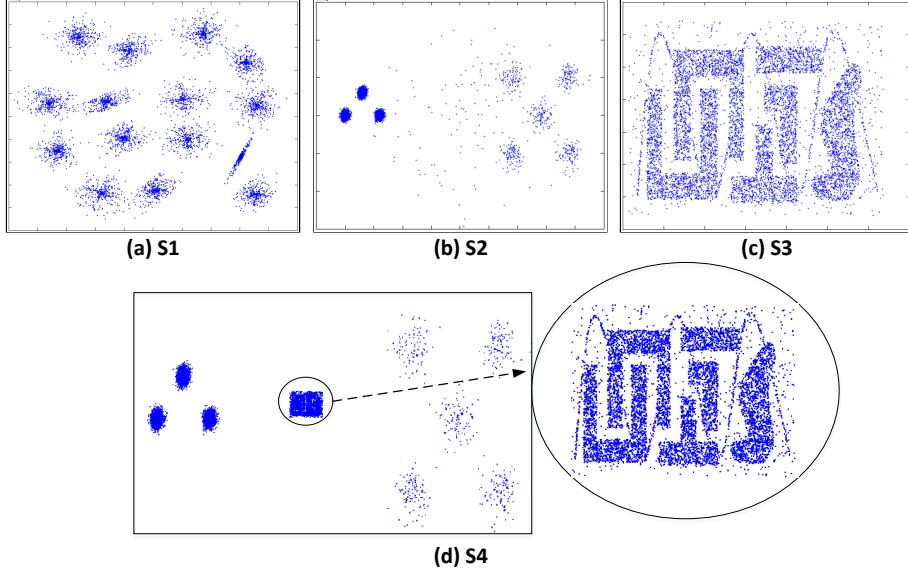


Figure 5: Two-dimensional synthetic datasets.

[42]. Since many methods share a common density measure, in the implementation, we consider the two density measures.

(i)  $\epsilon$  cut-off density. The parameter  $\epsilon$  is set to the value of the top  $\beta$  percent distance among all object pairs. We vary  $\beta$  value in  $\{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}$  according to the recommendation in [33].

(ii) Density based on  $K$  nearest neighbors. The parameter  $K$  is searched from  $(0, \sqrt{|V|})$  with a step size of 10% of the range, where  $V$  is the number of objects. This range covers the recommended parameters of different clustering methods and is shown to provide a good trade-off between evaluating clustering performance and computation time.

The detailed settings of all methods are listed in TABLE 2.

The best performance is reported for all methods by searching through the afore-mentioned parameter space.

#### 6.1.3. Performance Metrics

For the two-dimensional synthetic datasets, due to the lack of ground truth labels, we compare clustering results visually. For real-world datasets, two metrics are calculated based on the ground truth of the datasets. One is the normalized mutual information (NMI) [38], which is one of the most widely used measures of clustering quality. Given a dataset  $V$  of size  $n$ , suppose there are  $C$  clusters and  $C_t$  actual classes. Let  $n_i$ ,  $n^{(j)}$  and  $n_i^{(j)}$  denote the number of objects in cluster  $i$ , actual class  $j$ , and both cluster  $i$  and actual class  $j$ , respectively.

Table 2: The detailed settings of all methods.

Methods	Density measure	Other settings
DBSCAN [9]	(i)	The parameter $MinPts$ is determined by $\lambda \cdot \bar{\rho}$ , where $\lambda$ is searched in $\{0.05, 0.1, 0.2, 0.4, 0.8, 1.6\}$ and $\bar{\rho}$ is the mean density. Each outlier is assigned to its nearest cluster.
SNN [8]	(ii)	For remained parameters $\epsilon$ and $MinPts$ , both of them are set to $\lambda K$ , where $\lambda$ is chosen from $\{0.2, 0.4, 0.5, 0.6, 0.8\}$ , according to the recommendation in [8].
KNNC [40]	(ii)	No other parameters.
FDP [33]	(i)	The objects with top $C_t$ gamma values are chosen as cluster centers, where $C_t$ is the real class number.
3DC [23]	(i)	No other parameters.
STClu [42]	(ii)	No other parameters.
RECOME	(ii)	For parameter $\alpha$ , enumerate all JD values extracted by $FJDD$ algorithm. Euclidean distance is specified as the distance function.

Then, NMI can be computed by

$$NMI = \frac{\sum_{i=1}^C \sum_{j=1}^{C_t} \frac{n_i^{(j)}}{n} \log \frac{n n_i^{(j)}}{n_i n^{(j)}}}{\sum_{i=1}^C \frac{n_i}{n} \log \frac{n_i}{n} \sum_{j=1}^{C_t} \frac{n^{(j)}}{n} \log \frac{n^{(j)}}{n}}.$$

The other metric is the F value. Given a dataset  $V = \{v_1, v_2, \dots, v_n\}$  with cluster labels  $\{c_1, c_2, \dots, c_n\}$  and actual class labels  $\{l_1, l_2, \dots, l_n\}$ , define

$$\text{Correctness}(v_i, v_j) = \begin{cases} 1 & \text{if } l_i = l_j \Leftrightarrow c_i = c_j \\ 0 & \text{otherwise} \end{cases}.$$

Then, define  $PB$  and  $RB$  as:

$$PB = \frac{1}{n} \sum_{i=1}^n \frac{\sum_{i \neq j, c_i = c_j} \text{Correctness}(v_i, v_j)}{|\{v_j | i \neq j, c_i = c_j\}|},$$

$$RB = \frac{1}{n} \sum_{i=1}^n \frac{\sum_{i \neq j, l_i = l_j} \text{Correctness}(v_i, v_j)}{|\{v_j | i \neq j, l_i = l_j\}|}.$$

$PB$  and  $RB$  refer to precision b-cubed and recall b-cubed [10], respectively. Then F value is computed by  $F = \frac{2 \times PB \times RB}{PB + RB}$ . We use this measure because  $PB$  and  $RB$  have been shown superior than other indices [2].

Both NMI and F fall in  $[0, 1]$ , and a higher value denotes better clustering performance.

## 6.2. Experimental Results

### 6.2.1. Results on Synthetic Dataset

Figure 6 shows the clustering results on the two-dimensional datasets. From the last row of this figure, we can observe that, RECOME achieves desirable results for all the four datasets.

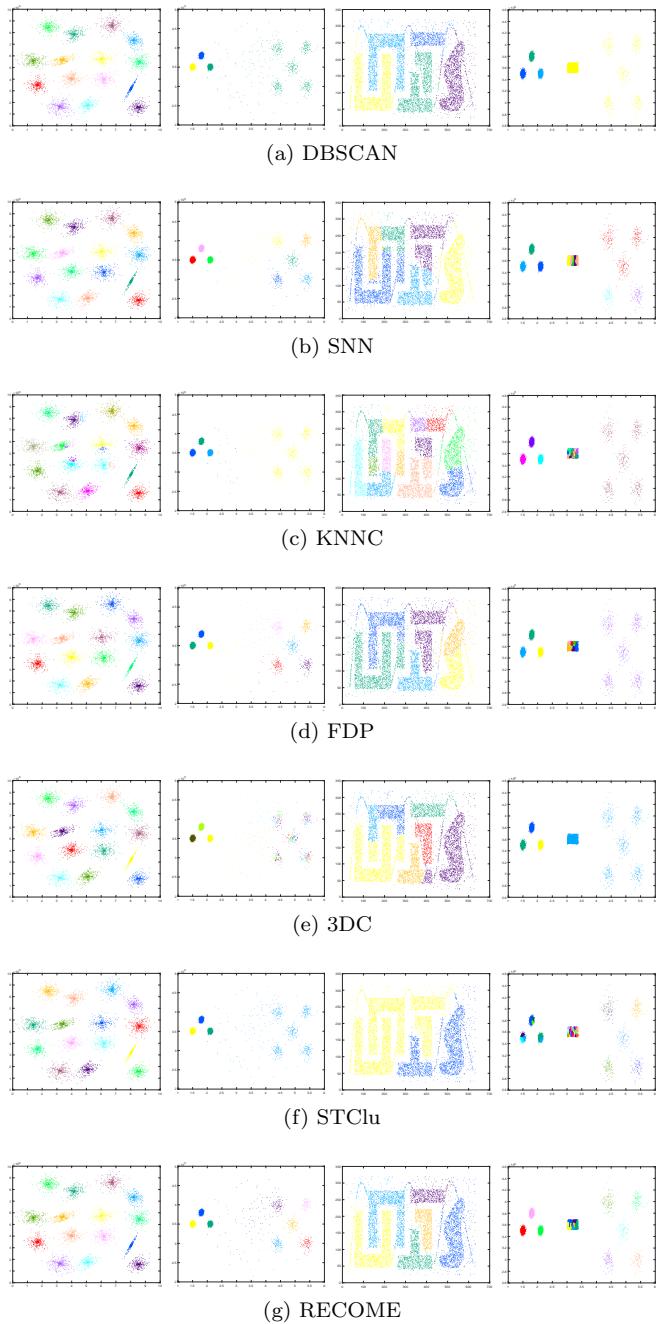


Figure 6: Clustering results on two-dimensional synthetic datasets.  
Different colors are used to represent different output clusters.

For dataset S1 (the first column) which contains only trivial Gaussian clusters, desirable results are achieved by all methods except for KNNC. For dataset S2 (the second column), which comprises unbalanced convex clusters and sparse noises, SNN, FDP and RECOME detect correct clusters. DBSCAN, KNNC and STClu overlook the clusters with little size. For dataset S3 (the third column), which is comprised of 6 clusters of nonconvex shape, a desirable result is output by RECOME. DBSCAN achieve a comparable and slightly poor performance. However, other methods rarely produce satisfactory results, mis-merge objects from different true clusters or subdivide a true cluster into different parts. For dataset S4 (the last column), which contains more complex clusters and proposes a big challenge to all baseline methods, only RECOME correctly detects the 14 clusters in it. These results indicate that RECOME is robust to the structure of cluster and own the potential to simultaneously discover clusters with different shapes, density, and scales.

#### 6.2.2. Results on MNIST Datasets

The clustering results on the MNIST datasets are presented in TABLE 3 and the visualization for sM3467 and sM0:8 are shown in Figure 7.

From TABLE 3, we can see that RECOME has the best performance for all datasets except M367, on which a comparable result is produced by STClu. SNN achieves slightly lower scores for almost all datasets than RECOME at much longer running time. Among all the algorithms, KNNC has the shortest running time for all datasets, but its performance is not steady—it achieves passable performance on M0:8 and M0:9 but poor performance on the remaining datasets. For FDP, we error on the conservative side by feeding it the true cluster number. Even so, RECOME achieves higher NMI and F values than FDP and the gap becomes obvious as the increasing of class number. Similar to FDP, the results of 3DC are not desirable for the datasets with many classes. The performance of STClu is unstable—it achieves comparable results on M367 to M3467 and moderate results on sM0:8 and M0:8, but fails to work well on M0:9. This may be due to it based on statistical testing, which is largely depended on data quality. Despite our best efforts in parameter tuning, DBSCAN fail to perform well over all datasets but M367.

As shown in Figure 7, for dataset sM3467, SNN and RECOME output almost correct results. FDP and STClu detect four main classes, but cluster excessive objects wrong classes. DBSCAN and 3DC fail to find the true cluster number for this dataset. For sM0:8, we can see that SNN and RECOME achieve the best performance but they both threat the two clusters of the digits “3” and “8” (the first column of the second row and the last column of the third row) as one. Besides, for the cluster of the digit “5”, SNN divides it into several sub-clusters. In contrast, RECOME detects this cluster but merge it into other clusters. Unfortunately, all other methods fail to output satisfactory results for this dataset, detecting few meaningful clusters or attributing many objects into wrong clusters.

Table 3: Performance comparison of the seven methods on the MNIST datasets.  $C$  is the cluster number identified by the algorithms (or utilized as prior knowledge for FDP) and  $C_t$  gives the true cluster number.

		DBSCAN	SNN	KNNC	FDP	3DC	STClu	RECOME
M367 ( $C_t = 3$ )	$C$	3	3	6	(3)	3	3	3
	NMI	.89	.70	.77	.93	.93	.96	<b>.98</b>
	F	.95	.79	.78	.97	.97	<b>.98</b>	.96
sM3467 ( $C_t = 4$ )	Time(s)	545	587	26	96	124	57	64
	$C$	1	5	13	(4)	3	5	4
	NMI	.00	.91	.72	.87	.76	.85	<b>.94</b>
M3467 ( $C_t = 4$ )	F	.40	.95	.67	.92	.81	.88	<b>.97</b>
	Time(s)	21	18	1	3	4	2	2
	$C$	1	4	10	(4)	3	5	4
sM0:8 ( $C_t = 9$ )	NMI	.00	.90	.76	.91	.80	.90	<b>.94</b>
	F	.40	.95	.74	.95	.83	.92	<b>.97</b>
	Time(s)	956	1006	46	170	235	99	116
M0:8 ( $C_t = 9$ )	$C$	2	26	24	(9)	7	17	9
	NMI	.20	.76	.68	.53	.35	.68	<b>.80</b>
	F	.33	.77	.54	.54	.42	.62	<b>.80</b>
M0:9 ( $C_t = 10$ )	Time(s)	106	95	5	16	25	10	14
	$C$	2	29	15	(9)	7	19	14
	NMI	.14	.78	.75	.47	.04	.73	<b>.82</b>
Time(s)	F	.30	.80	.69	.47	.22	.68	<b>.84</b>
	$C$	2	40	20	(10)	8	2	12
	NMI	.09	.72	.74	.43	.29	.26	<b>.80</b>
Time(s)	F	.24	.71	.70	.41	.37	.19	<b>.80</b>
	$C$	2	40	20	(10)	8	2	12
	NMI	.09	.72	.74	.43	.29	.26	<b>.80</b>
Time(s)	F	.24	.71	.70	.41	.37	.19	<b>.80</b>
	$C$	2	40	20	(10)	8	2	12
	NMI	.09	.72	.74	.43	.29	.26	<b>.80</b>
Time(s)	F	.24	.71	.70	.41	.37	.19	<b>.80</b>

### 6.2.3. Results on Olivetti Face Database

The clustering results on face10 and face40 are presented in TABLE4. The visualized results for face10 are shown in the rightmost column of Figure 7 (The visualization for face40 can be found in the supplementary material).

For face10, we can see that RECOME gets the highest score on both NMI and F metrics from TABLE 4. The other methods with the exception of STClu achieve similar scores, which are slightly lower than RECOME. Though the correct cluster number is found by RECOME, as shown in Figure 7, it mistakenly merges two clusters (the second row) and subdivides a true cluster (the left subject of the third row) into two classes.

For dataset face40, as shown in TABLE 4, KNNC achieves the best result in NMI and F measures, but this comes at the cost of an abnormally large cluster number. RECOME and SNN get comparable scores lower than KNNC by one percent and two percent in the NMI value, respectively. Besides, DBSCAN and FDP achieve reasonable performance and have a gap of 0.06 in the F metric compared with RECOME. Despite our best efforts in parameter tuning, 3DC and STClu fail to perform well for this dataset. This may be explained by



(a) DBSCAN



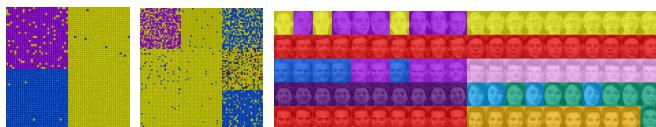
(b) SNN



(c) KNNC



(d) FDP



(e) 3DC



(f) STClu



(g) RECOME

Figure 7: Visualized clustering results on real datasets. Columns from left to right are visualizations for sm3467, sm0:8, and face10, respectively. Different clusters are marked by different colors.

Table 4: Performance comparison of the seven methods on the Olivetti face database, where  $C$  means the cluster number identified by algorithm (or utilized as prior knowledge, for FDP) and  $C_t$  means the true cluster number.

		DBSCAN	SNN	KNNC	FDP	3DC	STClu	RECOME
(C <sub>t</sub> = 10)	$C$	11	10	19	(10)	9	2	10
	NMI	.87	.88	.87	.89	.83	.42	<b>.94</b>
	F	.81	.81	.73	.82	.73	.32	<b>.90</b>
(C <sub>t</sub> = 40)	$C$	42	50	66	(40)	12	2	37
	NMI	.80	.87	<b>.89</b>	.83	.54	.27	.88
	F	.60	.64	<b>.67</b>	.60	.29	.09	.66

statistical errors since there are only 10 pictures for each people as discussed in [23, 42].

### 6.3. Parameter Analysis

As discussed previously, parameters  $K$  and  $\alpha$  play fundamental roles in RECOME. Specifically,  $K$  determines the density estimation and the structure of the KNN graph. The number of the final clusters are largely determined by  $\alpha$ . In this section, we present quantitative results on the impact of parameters  $K$  and  $\alpha$  on the performance of RECOME.

For a given dataset with  $|V|$  objects, the possible value of the number of  $K$  can vary from 1 to  $|V| - 1$ . However, we find that when  $K$  is large enough (up to  $\sqrt{|V|}$ ), increasing  $K$  has little or negative impacts on the performance of RECOME. Thus, in this set of experiments,  $K$  is taken values from  $(0, \sqrt{|V|}]$  with a step size of  $\sqrt{|V|}/40$ .

For parameter  $\alpha$ , as discussed in Section 5, only a few values will lead to different results. To validate this claim, we search  $\alpha$  in the range  $[0, 1]$  with a step size of 0.01 to make a better showcase and a evident staircase can be observed from the presentation. To evaluate the clustering results, both NMI and F metrics are used. Next, we present impacts of parameters  $K$  and  $\alpha$  on NMI and F values for all datasets.

For synthetic datasets, due to the lack of ground truth for S2, S3, and S4, we label them according to the results shown in Figure 6g as they agree with our intuitive understanding. As shown in Figure 8a, for all the four datasets, both NMI and F approach one when  $K$  is above  $\sqrt{|V|}/4$  (in the figure,  $k = 10$  for  $K = \sqrt{|V|}/4$ ) for appropriate  $\alpha$ 's, and then change very slowly despite the increase of  $K$ . Meanwhile, when  $K$  falls in the range  $[\sqrt{|V|}/4, \sqrt{|V|}]$ ,  $\alpha$  valued greater than 0.8 will lead to desirable results. In addition, it is noted that, compared with S3 and S4, the maximum NMI and F values can be attained easier for S1 and S2. This is due to the fact that clusters in both S1 and S2 are convex and well separated, but for the other two datasets, irregular shapes and scales make it difficult to detect the true clusters. Regardless of the complex shapes and scales, RECOME can find the correct cluster numbers and output the desirable results in all cases with properly selected parameters.

Figure 8e to Figure 8j show the influences of parameters  $K$  and  $\alpha$  on performance for MNIST training set. We can see that, similar to the synthetic

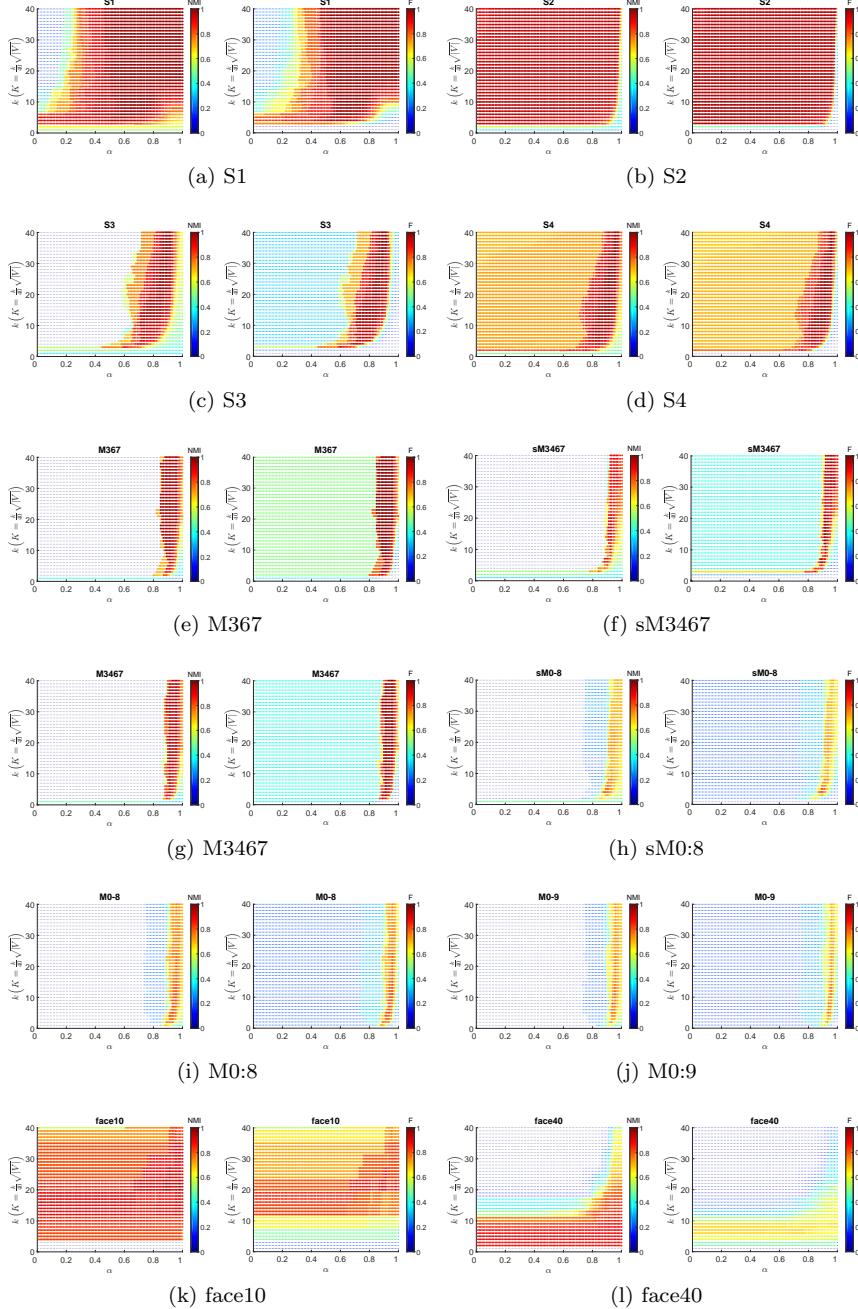


Figure 8: (a)-(d) The clustering performance on 2-D synthetic datasets as a function of  $K$  and  $\alpha$ . (e)-(j) The clustering performance on MNIST training set as a function of  $K$  and  $\alpha$ . (k)-(l) The clustering performance on Olivetti face database as a function of  $K$  and  $\alpha$ .

datasets, both NMI and F values plateau out as  $K$  take around  $\sqrt{|V|}/4$  for all 6 datasets. At the same time, when  $K$  is in the interval  $[\sqrt{|V|}/2, \sqrt{|V|}]$ , the  $\alpha$  value in the range  $[0.85, 0.95]$  gives the best performance. On the other hand, as the number of clusters increases, the clustering performance degrades. For example, NMI over 0.9 can be easily attained for M3467 but not for M0:9. This is because more clusters tend to have more complex sample distribution and more overlapping among different classes. Furthermore, compared with two-dimensional datasets, we can observe that the region in the  $\alpha$ - $K$  plot that reaches high NMI and F values shrinks. This is due to the fact that the MNIST training set with high dimension is far harder to be clustered well. However, RECOME still achieves better performance across almost all six datasets compared with baseline methods.

For the Olivetti Face Database, the trends NMI and F are different from that in the other datasets. Performance on face10 follows a step-wise pattern with regard to  $K$  due to the small sample size (i.e.,  $\sqrt{|V|} = 10$ ). For face40, both the NMI and F values remain stable when  $K$  varies from  $\sqrt{|V|}/40$  to  $\sqrt{|V|}/4$ . Then, with the increasing of  $K$ , the two indices drop rapidly and become quite small. This is because the ground truth partition constitutes 40 clusters and only 10 pictures in each cluster. The statistical error of the estimated density on such a small set of pictures is large [33]. Therefore, for datasets consisting of clusters with few objects,  $K$  should be carefully tuned.

In summary, the performance of RECOME is dependent on both parameters of  $K$  and  $\alpha$ . For most datasets,  $K$  in the range of  $[\sqrt{|V|}/2, \sqrt{|V|}]$  leads to a stable condition. Under this condition, the  $\alpha$  value that produces a good performance falls in the range of  $[0.8, 1]$ . This implies that RECOME is not very sensitive to small changes in parameters  $K$  and  $\alpha$ . Furthermore, as discussed in Section 5, the algorithm FJDD can quickly determine the JD set with tiny size. Thus, those who fall in  $[0.8, 1]$  can then be used as the candidate set for  $\alpha$ .

#### 6.4. FJDD validation

As discussed in Section 5, the number of clusters produced by RECOME is a step function of  $\alpha$ . In this experiment, we fix  $K$  at  $0.5\sqrt{|V|}$  and search  $\alpha$  in  $[0, 1]$  with a step size 0.01. The resulting cluster number as a function of  $\alpha$  is plotted in Figure 9. It can be observed that all curves are step-wise.

We further verify the correctness of Algorithm 3 on all datasets used in this paper with the jump points shown in the magenta dashed lines. From Figure 9, it can be seen that FJDD extracts all jump discontinuity correctly for all datasets. Interestingly, the number of jump discontinuity is less than 20 for all datasets, which means that to obtain all possible clustering results, only at most 20% of run times is needed compared with searching for  $\alpha$  in  $[0, 1]$  with the step size 0.01.

## 7. Conclusion

In this paper, we presented a new density-based clustering method RECOME. RECOME exploits a novel density measure RNKD to detect core ob-

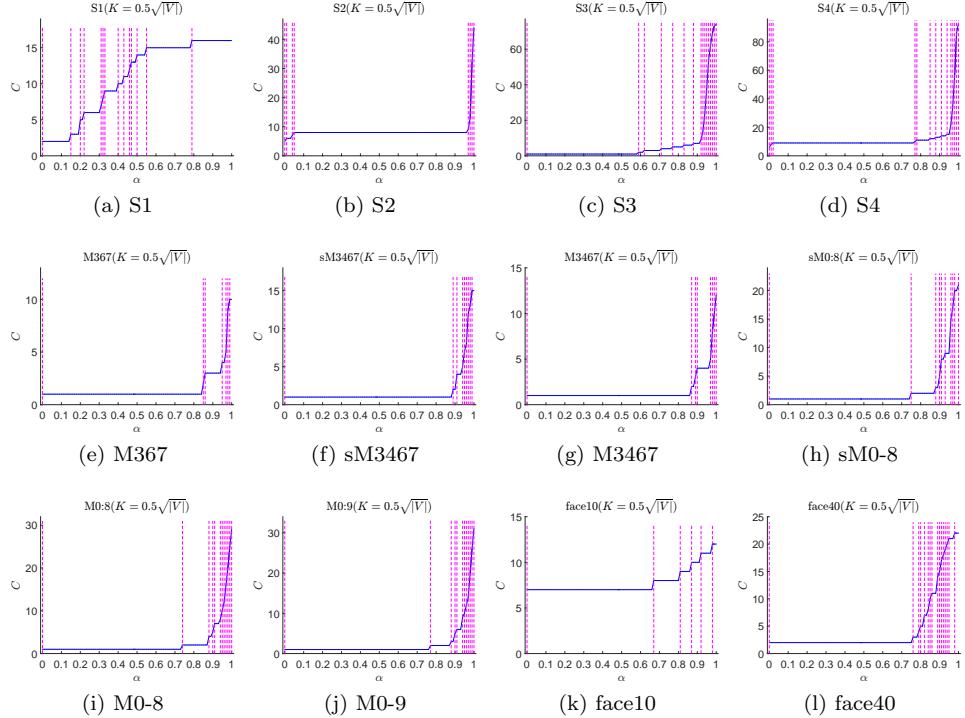


Figure 9: The cluster number  $C$  as a step function of  $\alpha$  when  $K = 0.5\sqrt{|V|}$ . The elements of JD set  $L$  produced by Algorithm 3 are marked by magenta dashed line.

jects with various densities and generate atom clusters. A merging strategy based on KNN graph has been introduced to refine the atom clusters. In addition, we discovered that the number of clusters obtained by RECOME is a monotonic step function of  $\alpha$ , and therefore proposed an auxiliary algorithm FJDD to help end users to select parameter  $\alpha$ . Experiment evaluations showed that RECOME can discover clusters of different shapes, densities, and scales with parameters chosen in stable ranges. FJDD can significantly reduce the number of effective choices for the parameter  $\alpha$ .

In future research, we will extend the current work in three aspects. Firstly, we will design a distributed RECOME in cluster-computing frameworks (e.g. Apache Spark) to handle large volume of data. Secondly, when facing massive data, the number of jump points can grow significantly. We will investigate how to quickly identify a relevant range of  $\alpha$ . Lastly, it is attractive to design an automatic strategy for parameter selection to make RECOME parameter-free.

## Acknowledgments

This work is partly supported by the Fundamental Research Funds for the Central Universities (2016JBZ006), Beijing Natural Science Foundation (No. J160004), and National Science and Engineering Council, Canada.

## Appendix A. Proof of Theorem 1

**Lemma 1.**  $\forall u, v \in V$ , the difference between  $d_i(u)$  and  $d_i(v)$  is not greater than  $d(u, v)$ , i.e.,  $|d_i(u) - d_i(v)| \leq d(u, v)$  ( $i = 1, 2, \dots, K$ ).

*Proof.* It is equivalent to show that the following holds,

$$d_i(v) - d(u, v) \leq d_i(u) \leq d_i(v) + d(u, v). \quad (\text{A.1})$$

If  $u = v$ , the inequality holds trivially. If  $u \neq v$ , we first show

$$d_i(u) \leq d_i(v) + d(u, v). \quad (\text{A.2})$$

Assume the conclusion is false, namely,  $d_i(u) > d_i(v) + d(u, v)$ . Then, we can see  $d_i(u) > d(u, v)$  according to the non-negativity of distance. On the other hand, we have  $d_i(v) = \max_{j \leq i} \{d(v, N_j(v))\}$ . From triangle inequality, we have

$$\begin{aligned} d_i(u) &> d_i(v) + d(u, v) \\ &= \max_{j \leq i} \{d(v, N_j(v))\} + d(u, v) \\ &\geq \max_{j \leq i} \{d(u, N_j(v))\}. \end{aligned}$$

Denote  $A = \{v, N_1(v), N_2(v), \dots, N_i(v)\}$ . Therefore,  $\forall w \in A$   $d_i(u) > d(u, w)$ . But  $A - \{u\}$  contains at least  $i$  elements, a contradiction with the definition of  $d_i(u)$ . Thus inequality (A.2) holds.

Similar argument holds for the left-side of inequality (A.1).  $\square$

Base on Lemma 1, we can prove Theorem 1.

*Proof.*

$$\begin{aligned} |\rho(u) - \rho(v)| &= \theta \left| \sum_{k=1}^K \exp \left( -\frac{d_k(u)}{\sigma} \right) - \sum_{k=1}^K \exp \left( -\frac{d_k(v)}{\sigma} \right) \right| \\ &\leq \theta \sum_{k=1}^K \left| \exp \left( -\frac{d_k(u)}{\sigma} \right) - \exp \left( -\frac{d_k(v)}{\sigma} \right) \right|. \end{aligned}$$

Let  $h(k) = \arg \min_{w \in \{u, v\}} \{d_k(w)\}$ . According to Lemma 1 and the monotonicity of the exp function, we have

$$\begin{aligned}
& \theta \sum_{k=1}^K \left| \exp\left(-\frac{d_k(u)}{\sigma}\right) - \exp\left(-\frac{d_k(v)}{\sigma}\right) \right| \\
& \leq \theta \sum_{k=1}^K \left( \exp\left(-\frac{d_k(h(k))}{\sigma}\right) - \exp\left(-\frac{d_k(h(k)) + d(u, v)}{\sigma}\right) \right) \\
& = \left(1 - \exp\left(-\frac{d(u, v)}{\sigma}\right)\right) \theta \sum_{k=1}^K \exp\left(-\frac{d_k(h(k))}{\sigma}\right) \\
& < \left(1 - \exp\left(-\frac{d(u, v)}{\sigma}\right)\right) \theta \sum_{k=1}^K \left( \exp\left(-\frac{d_k(u)}{\sigma}\right) + \exp\left(-\frac{d_k(v)}{\sigma}\right) \right) \\
& = \left(1 - \exp\left(-\frac{d(u, v)}{\sigma}\right)\right) (\rho(u) + \rho(v)).
\end{aligned}$$

Since  $\rho(u) + \rho(v) > 0$ , the result follows.  $\square$

## Appendix B. Proof of Proposition 1

*Proof.* Denote  $A = \{0\} \cup \{c(u, v) | u, v \in O, u \neq v\} = \{\alpha_1, \alpha_2, \dots, \alpha_t\}$ , where  $0 \equiv \alpha_1 < \alpha_2 < \dots < \alpha_t$ . For any  $\alpha$  with  $0 \leq \alpha \leq 1$ , let  $G_K^{>\alpha}$  denote the remaining graph after removing all non-core nodes with weights not larger than  $\alpha$  from  $G_K$ . In addition, let  $\#(G_K^{>\alpha})$  denote the number of components containing at least one core object in  $G_K^{>\alpha}$ . The key of the proof is the following equation.

$$\#(V, K, \alpha_i) = \#(G_K^{>\alpha_i}). \quad (\text{B.1})$$

To prove the conclusion, only two things need to be shown. First,  $\forall i \neq j$ ,  $\#(V, K, \alpha_i) \neq \#(V, K, \alpha_j)$ . Second,  $\forall 0 \leq \alpha \leq 1$ ,  $\exists \alpha_i \in A$  such that  $\#(V, K, \alpha) = \#(V, K, \alpha_i)$ .

Now we show the former. Obviously,  $\#(V, K, \alpha)$  is a non-decreasing function with respect to  $\alpha$ , which implies that showing “ $\forall 1 \leq i < t, \#(V, K, \alpha_i) < \#(V, K, \alpha_{i+1})$ ” is enough. Besides, for any  $u, v \in O$  with  $u \neq v$ , “ $u$  and  $v$  are located in different components in  $G_K^{>\alpha_i}$ ” implies “ $u$  and  $v$  are located in different components in  $G_K^{>\alpha_{i+1}}$ ”. On the other hand, from the definition of  $A$ , there exists a pair of core objects  $u$  and  $v$  such that  $u$  and  $v$  are connected in  $G_K^{>\alpha_i}$  but disconnected in  $G_K^{>\alpha_{i+1}}$ . This implies “ $u$  and  $v$  are located in the same component in  $G_K^{>\alpha_i}$ ” and “ $u$  and  $v$  are located in different components in  $G_K^{>\alpha_{i+1}}$ ”. Therefore, we have  $\#(G_K^{>\alpha_i}) < \#(G_K^{>\alpha_{i+1}})$ . According to (B.1), the conclusion follows.

Now we show the latter. For any  $\alpha$  with  $0 \leq \alpha \leq 1$ , suppose  $\alpha_i$  is the maximal element of  $A$  that is not greater than  $\alpha$ . According to the definition of  $A$ , we have  $\#(G_K^{>\alpha}) = \#(G_K^{>\alpha_i})$ . Thus, according to (B.1), the desirable result follows.  $\square$

## References

### References

- [1] R. Agrawal, J. Gehrke, D. Gunopulos, P. Raghavan, Automatic subspace clustering of high dimensional data for data mining applications, volume 27, ACM, 1998.
- [2] E. Amig, J. Gonzalo, J. Artiles, F. Verdejo, A comparison of extrinsic clustering evaluation metrics based on formal constraints, *Information Retrieval* 12 (2009) 461–486.
- [3] M. Ankerst, M. M. Breunig, H.-P. Kriegel, J. Sander, Optics: ordering points to identify the clustering structure, in: *ACM SIGMOD*, volume 28, 1999, pp. 49–60.
- [4] N. Beckmann, H.-P. Kriegel, R. Schneider, B. Seeger, The r\*-tree: an efficient and robust access method for points and rectangles, in: *ACM Sigmod Record*, volume 19, Acm, 1990, pp. 322–331.
- [5] M. M. Breunig, H. P. Kriegel, R. T. Ng, J. Sander, Lof: identifying density-based local outliers, *Acm Sigmod Record* 29 (2000) 93–104.
- [6] R. A. Brown, Building a balanced k-d tree in  $o(kn \log n)$  time, *CoRR* abs/1410.5420 (2014).
- [7] C. Cassisi, A. Ferro, R. Giugno, G. Pigola, A. Pulvirenti, Enhancing density-based clustering: Parameter reduction and outlier detection, *Information Systems* 38 (2013) 317–330.
- [8] L. Ertöz, M. Steinbach, V. Kumar, Finding clusters of different sizes, shapes, and densities in noisy, high dimensional data, in: *Siam International Conference on Data Mining*, San Francisco, Ca, Usa, May, 2003.
- [9] M. Ester, H.-P. Kriegel, J. Sander, X. Xu, A density-based algorithm for discovering clusters in large spatial databases with noise., in: *ACM SIGKDD*, volume 96, 1996, pp. 226–231.
- [10] J. Han, M. Kamber, J. Pei, *Data mining: concepts and techniques*, Elsevier, 2011.
- [11] R. He, Q. Li, B. Ai, Y. L. A. Geng, A. F. Molisch, V. Kristem, Z. Zhong, J. Yu, A kernel-power-density-based algorithm for channel multipath components clustering, *IEEE Transactions on Wireless Communications* 16 (2017) 7138–7151.
- [12] A. Hinneburg, D. A. Keim, An efficient approach to clustering in large multimedia databases with noise, *KDD* 98, (1999).
- [13] A. K. Jain, Data clustering: 50 years beyond k-means, *Pattern recognition letters* 31 (2010) 651–666.

- [14] L. Kaufman, P. J. Rousseeuw, Agglomerative nesting (program agnes), *Finding Groups in Data: An Introduction to Cluster Analysis* (2008) 199–252.
- [15] L. Kaufman, P. J. Rousseeuw, Divisive analysis (program diana), *Finding Groups in Data: An Introduction to Cluster Analysis* (2008) 253–279.
- [16] L. Kaufmann, P. J. Rousseeuw, Clustering by means of medoids, in: *Statistical Data Analysis Based on the L1-norm & Related Methods*, 1987, pp. 405–416.
- [17] M. Kim, J. Han, A particle-and-density based evolutionary clustering method for dynamic networks, *PVLDB* 2 (2009) 622–633.
- [18] N. R. Kisore, C. B. Koteswaraiah, Improving atm coverage area using density based clustering algorithm and voronoi diagrams, *Information Sciences* 376 (2016) 1–20.
- [19] H.-P. Kriegel, P. Kr?ger, J. Sander, A. Zimek, Density-based clustering, *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* 1 (2011) 231–240.
- [20] J. Kubalk, P. Tich?, R. ?indel?, R. J. Staron, Clustering methods for agent distribution optimization, *IEEE Transactions on Systems Man and Cybernetics Part C* 40 (2010) 78–86.
- [21] S. Laohakiat, S. Phimoltares, C. Lursinsap, A clustering algorithm for stream data with lda-based unsupervised localized dimension reduction, *Information Sciences* 381 (2016) 104C123.
- [22] Y. Lecun, L. Bottou, Y. Bengio, P. Haffner, Gradient-based learning applied to document recognition, *Proceedings of the IEEE* 86 (1998) 2278–2324.
- [23] Z. Liang, P. Chen, Delta-density based clustering with a divide-and-conquer strategy: 3dc clustering, *Pattern Recognition Letters* 73 (2016) 52–59.
- [24] P. Liu, D. Zhou, N. Wu, Vdbscan: varied density based spatial clustering of applications with noise, in: *2007 International conference on service systems and service management*, IEEE, 2007, pp. 1–4.
- [25] D. O. Loftsgaarden, C. P. Quesenberry, A nonparametric estimate of a multivariate density function, *Annals of Mathematical Statistics* 36 (1965) 1049–1051.
- [26] W. Loh, H. Yu, Fast density-based clustering through dataset partition using graphics processing units, *Information Sciences* 308 (2015) 94–112.

- [27] J. MacQueen, et al., Some methods for classification and analysis of multivariate observations, in: Proceedings of the fifth Berkeley symposium on mathematical statistics and probability, volume 1, Oakland, CA, USA., 1967, pp. 281–297.
- [28] H. Miller, J. Han, Spatial clustering methods in data mining: a survey, *Geographic data mining and knowledge discovery*, Taylor and Francis (2001).
- [29] S. J. Nanda, G. Panda, Design of computationally efficient density-based clustering algorithms, *Data and Knowledge Engineering* 95 (2015) 23–38.
- [30] M. E. Newman, Modularity and community structure in networks, *Proceedings of the national academy of sciences* 103 (2006) 8577–8582.
- [31] T. Pei, A. Jasra, D. J. Hand, A.-X. Zhu, C. Zhou, Decode: a new method for discovering clusters of different densities in spatial data, *Data Mining and Knowledge Discovery* 18 (2009) 337–369.
- [32] Z. Qiu, H. Shen, User clustering in a dynamic social network topic model for short text streams, *Information Sciences* 414 (2017) 102 – 116.
- [33] A. Rodriguez, A. Laio, Clustering by fast search and find of density peaks, *Science* 344 (2014) 1492–1496.
- [34] F. S. Samaria, A. C. Harter, Parameterisation of a stochastic model for human face identification, in: *Proceedings of 1994 IEEE Workshop on Applications of Computer Vision*, 1994, pp. 138–142. doi:10.1109/ACV.1994.341300.
- [35] J. Sander, M. Ester, H. P. Kriegel, X. Xu, Density-based clustering in spatial databases: The algorithm gdbcscan and its applications, *Data Mining and Knowledge Discovery* 2 (1998) 169–194.
- [36] A. Shapira, R. Yuster, U. Zwick, All-pairs bottleneck paths in vertex weighted graphs, *Algorithmica* 59 (2011) 621–633.
- [37] J. Shi, J. Malik, Normalized cuts and image segmentation, *IEEE Trans.pattern Anal.mach.intell* 22 (2000) 888–905.
- [38] A. Strehl, J. Ghosh, Cluster ensembles — a knowledge reuse framework for combining multiple partitions, *Journal of Machine Learning Research* 3 (2002) 583–617.
- [39] G. R. Terrell, D. W. Scott, Variable kernel density estimation, *The Annals of Statistics* (1992) 1236–1265.
- [40] T. N. Tran, R. Wehrens, L. M. Buydens, Knn-kernel density-based clustering for high-dimensional multivariate data, *Computational Statistics and Data Analysis* 51 (2006) 513–525.

- [41] C. D. Wang, J. H. Lai, J. Y. Zhu, Graph-based multiprototype competitive learning and its applications, *IEEE Transactions on Systems Man and Cybernetics Part C* 42 (2012) 934–946.
- [42] G. Wang, Q. Song, Automatic clustering via outward statistical testing on density metrics, *IEEE Transactions on Knowledge & Data Engineering* 28 (2016) 1–1.
- [43] W. Wang, J. Yang, R. R. Muntz, Sting: A statistical information grid approach to spatial data mining, in: VLDB, volume 97, 1997, pp. 186–195.
- [44] X. Xu, N. Yuruk, Z. Feng, T. A. J. Schweiger, Scan: a structural clustering algorithm for networks, in: ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2007, pp. 824–833.
- [45] Q. Ye, W. Gao, W. Zeng, Color image segmentation using density-based clustering, in: Multimedia and Expo, 2003. ICME '03. Proceedings. 2003 International Conference on, volume 2, 2003, pp. II–401–4 vol.2. doi:10.1109/ICME.2003.1221638.
- [46] J. Yin, J. Wang, A dirichlet multinomial mixture model-based approach for short text clustering, in: ACM SIGKDD, 2014, pp. 233–242.
- [47] H. Zhang, H. Zhai, L. Zhang, P. Li, Spectral–spatial sparse subspace clustering for hyperspectral remote sensing images, *IEEE Transactions on Geoscience and Remote Sensing* 54 (2016) 3672–3684.