

Highlights

Clustering Based on Local Density Peaks and Graph Cut*

Zhiguo Long, Yang Gao, Hua Meng, Yuqin Yao, Tianrui Li

- A novel measure of similarity between local tree-structured clusters to capture complex structures of data is proposed.
- A new Ncut-based loss function is devised in consideration of tree structures.
- The effectiveness of the proposed algorithm is verified on both real-world and synthetic complex data sets.

Clustering Based on Local Density Peaks and Graph Cut

Zhiguo Long^a, Yang Gao^a, Hua Meng^{b,*}, Yuqin Yao^b, Tianrui Li^a

^a*School of Computing and Artificial Intelligences, Southwest Jiaotong University, Chengdu, 611756, China*

^b*School of Mathematics, Southwest Jiaotong University, Chengdu, 611756, China*

Abstract

Clustering by fast search and find of density peaks (DPC) is a widely used and studied clustering algorithm. In this article, we notice that DPC can achieve highly accurate clustering results when restricted to local neighborhoods. Therefore, by investigating density information in local neighborhoods, we propose to capture latent structures in data with *family trees*, which can reflect density dominations among nearest neighbors of data. A data set will then be partitioned into multiple family trees. In order to obtain the final clustering result, instead of exploiting the error-prone allocation strategy of DPC, we first elaborately design a novel similarity measure for family trees, characterizing not only the distance between data points, but also the structure of trees. Then, we adapt graph cut for the corresponding connection graph to also take global structural information into account. Extensive experiments on both real-world and synthetic data sets show that the proposed algorithm can outperform several prominent clustering algorithms for most of the cases, including the DPC and spectral clustering algorithms and some of their latest variants. We also analyze the robustness of the proposed algorithm w.r.t. hyper-parameters and its time complexity, as well as the necessity of its components through ablation study.

Keywords: Clustering, Density peaks, Spectral clustering, Local density, Similarity between trees

1. Introduction

Clustering [1] is an important unsupervised problem in the field of data mining and machine learning. It aims to divide a data set into multiple disjoint subsets, where the data in each subset are as similar as possible, and the data between subsets are as dissimilar as possible [2, 3]. Clustering can reveal inherent or latent knowledge and rules in data and has been widely used in scientific research and engineering applications [4, 5].

In the past decades, a large number of clustering algorithms have been proposed and applied in different scenarios. Traditional clustering algorithms can be roughly divided into

*Accepted version. Formal version available at: <https://doi.org/10.1016/j.ins.2022.03.091>

*Corresponding author.

Email address: menghua@swjtu.edu.cn (Hua Meng)

9 five categories: division based [6], hierarchical based [7], density based [8], model based [9],
10 and grid based [10]. In addition, there are also extensions like spectral clustering [11],
11 multi-view clustering [12], clustering ensemble [13] and deep clustering [14]. Each of these
12 algorithms has its own advantages and disadvantages for certain type of data or applications.

13 Clustering by fast search and find of density peaks (DPC; [15]) is an influential density
14 based clustering algorithm. Although shown to perform well on a variety of data, it has
15 several problems such as: it cannot properly process data sets with variant densities in
16 different part; its usage of Euclidean distance for calculating densities and searching of
17 density peaks is not appropriate for manifold structures; its allocation strategy of clusters can
18 suffer from the domino effect that the misallocation of one point will result in misallocation of
19 all subsequent points. A large number of improvements and extensions of the original DPC
20 algorithm were proposed to overcome those problems. For example, KNN-DPC [16] and
21 FDPC [17] try to improve DPC by optimizing the calculation of densities; FKNN-DPC [18]
22 and SNN-DPC [19] further improve allocation strategy; DLORE-DP [20] and DPC-GD [21]
23 introduce geodesic distance to adapt for manifold structures; FastDPeak [22] and DPCG [23]
24 optimize computational efficiency.

25 On the other hand, in this article, instead of improving density calculation, peak selec-
26 tion, or allocation strategy of DPC, we try to make use of certain merit in the idea of DPC
27 to devise a new clustering algorithm. In fact, although the allocation strategy of DPC is not
28 always reliable globally, it can achieve high *purity* in small local subsets of data (cf. Fig. 2),
29 i.e., locally, points that should be in the same cluster are also allocated to the same cluster.
30 By restricting DPC to local neighborhoods, we can obtain tree structures (cf. Fig. 1) that
31 reflect density domination relations in local neighborhood, and each of these trees has high
32 purity. The aforementioned DLORE-DP considers this kind of local structures as natu-
33 ral accumulations of data points, and makes use of DPC to cluster these accumulations in
34 the subsequent stage. FHC-LDP [24] uses these local trees as sub-clusters and combines
35 them with the help of hierarchical clustering. Besides, ADBC-KNN [25] and CDMC-IA [26]
36 exploit the idea of DPC to obtain sub-clusters, and propose strategies to achieve final clus-
37 tering, i.e., combining sub-clusters according to density reachability and quantity affinity
38 (in a hierarchical clustering way), respectively.

39 In this article, we further investigate the properties of such tree-based local structures,
40 and propose a novel clustering algorithm (LDP-SC), by designing a more sophisticated
41 similarity measure of these trees, and adopting an improved graph cut method. The use
42 of graph cut method can better determine whether to combine or to separate these trees,
43 because it takes into account more global structural information of data, than hierarchical
44 clustering, which mainly considers local information of sub-clusters. This graph cut method
45 also avoids the identification of cluster centers, which is an error-prone process, especially
46 for complex data.

47 This article has the following contributions:

- 48 1. We study the properties of local tree-structured clusters (called *family trees*), and pro-
49 pose a novel similarity measure for such local clusters, which helps to capture complex
50 structures of data.

- 51 2. We exploit the idea of spectral clustering on these local structures to obtain the final
52 clustering result, which avoids the error-prone steps of the identification of cluster centers
53 and the allocation of clusters.
- 54 3. We improve the loss function for graph cut to better suit local tree-structured clusters.
- 55 4. The proposed algorithm outperforms several prominent baseline algorithms, including
56 the DPC and spectral clustering algorithms and some of their latest variants.

57 The rest of this article is organized as follows. In Section 2, research works on improving
58 the DPC and spectral clustering algorithms in recent years are discussed, as well as work on
59 the combination of DPC and spectral clustering. In Section 3, preliminaries for DPC and
60 spectral clustering are introduced, as well as an analysis for the motivation of the improve-
61 ments proposed in this article. Section 4 gives details on the proposed LDP-SC algorithm.
62 Section 5 then presents a number of experiments for empirically evaluating the proposed
63 algorithm against some benchmark algorithms from various aspects and on different data
64 sets. Section 6 discusses the robustness of the LDP-SC algorithm, and performs ablation
65 study to confirm the significance of each component in the algorithm. Section 7 concludes
66 the article.

67 2. Related Work

68 In recent years, a large number of variants of DPC have been proposed to overcome its
69 defects. Generally, DPC consists of three essential parts, i.e., the calculation of densities,
70 the selection of cluster centers, and the cluster allocation strategy of non-center points.

71 KNN-DPC [16] uses K-nearest neighbors to replace the use of cutoff distance for the cal-
72 culation of densities. FDPC [17] and FKNN-DPC [18] exploit fuzzy theory to characterize
73 uncertainties introduced by the sampling of data points, and thus generalize the calculation
74 of densities. FKNN-DPC also modifies the allocation strategy by making use of neighboring
75 relations. Based on KNN-DPC and FKNN-DPC, DPC-DBFN [27] uses a fuzzy kernel to
76 extract a specific kind of structural information, called *density backbone*, to improve the
77 robustness of clustering. SNN-DPC [19] investigates local structures of data, and charac-
78 terizes similarities between data points in a more sophisticated way by considering shared
79 nearest neighbors. It further improves the calculation of densities by using the new similarity
80 measure, and the allocation strategy by distinguishing inevitable and possible subordinate
81 points. BPEC [28] replaces densities with beliefs and the allocation strategy with an op-
82 timization problem called *credal partition*, and can provide fuzzy or even rough partitions.
83 CFDPC [29] proposes an automatic method to filter out density peaks and thus several initial
84 clusters, conducts adaptive searching of core points, and finally fuses initial clusters based
85 on these core points in a hierarchical way. Considering manifold structures of some data,
86 DPC-GD [21] and DLORE-DP [20] replace Euclidean distance with geodesic distance to im-
87 prove the applicability to manifold structures. ADBC-KNN [25] and CDMC-IA [26] select
88 multiple centers (generally more than the number of target clusters), by setting a threshold
89 on the value of γ of DPC, and then combine the sub-clusters determined by these centers
90 to obtain the final clustering result, by density reachability and by quantity affinity in a

91 hierarchical way, respectively. Instead of selecting centers, LDP-MST [30] generates several
 92 minimum spanning trees and constructs connection relations between these trees. The final
 93 clustering result is obtained by removing some of the edges between trees. FHC-LDP [24]
 94 builds sub-clusters with local density structures, and then combines these sub-clusters by
 95 hierarchical clustering. It is worth noting that, identifying sub-clusters is a way to discover
 96 *local* structural information, and combining them with hierarchical clustering or density
 97 reachability emphasizes again on *local* information. Instead, here we propose to use graph
 98 cut method for the purpose of aggregating sub-clusters, which arguably better balances *local*
 99 and *global* information.

100 The spectral clustering algorithm [3] is a typical graph cut method, which is simple,
 101 effective and widely used. However, one of the most obvious disadvantages is its high time
 102 complexity $\mathcal{O}(m^3)$, where m is the number of data points, and in recent years, improvements
 103 of the spectral clustering algorithm have mainly focused on computational efficiency. To
 104 achieve the purpose of acceleration, Cheng et al. [31] randomly select p representative points
 105 to approximate the original data set and make use of approximate matrix decomposition. Cai
 106 et al. [32] propose the landmark-based spectral clustering algorithm (LSC), which obtains
 107 p representative points as landmarks by random selection or K-means pre-clustering, and
 108 then uses sparse linear combinations of these landmarks to approximate adjacency matrix.
 109 Huang et al. [33] use approximated K-nearest neighbors to construct a sparse adjacency
 110 submatrix and interpret connections between data points as a bipartite graph, which is then
 111 divided into clusters by the Ncut algorithm. This enables the spectral clustering algorithm
 112 to efficiently perform clustering on data sets with millions of data points. Our idea of
 113 performing graph cut on family trees obtained from local structural information improves
 114 both the computational efficiency and effectiveness of graph cut, because points in the same
 115 family tree will be considered as a whole and will never be separated.

116 There are also some researches on combining the ideas of the DPC algorithm with the
 117 spectral clustering algorithm. Li et al. [34] propose the DPSC algorithm, which replaces the
 118 K-means clustering with DPC in the second stage of spectral clustering, so that the decision
 119 graph of DPC could be used to determine the number of clustering and filter noises. Liu and
 120 Zhao [35] use changes of densities to improve the DPC, making it more effective for data
 121 with uneven density. They also apply the improved DPC algorithm to the second stage of
 122 spectral clustering.

123 3. Preliminaries

124 3.1. DPC

125 The idea of DPC to cluster data is first to identify several density peaks as cluster centers
 126 and then to allocate the other data points to the clusters of the corresponding centers. For
 127 cluster centers, DPC makes the following assumptions:

- 128 1. A cluster center has a higher density than the surrounding data points;
- 129 2. The distance between cluster centers is far.

Based on these assumptions, DPC calculates ρ_i and δ_i to determine which data points can be cluster centers, where ρ_i represents the density of a data point x_i and δ_i represents the distance between x_i and its nearest neighbor with higher density, defined as:

$$\rho_i = \sum_{i \neq j} \chi(d_{ij} - d_c), \text{ where } \chi(x) = \begin{cases} 1, & x < 0 \\ 0, & x \geq 0 \end{cases}, \text{ and} \quad (1)$$

$$\delta_i = \begin{cases} \min_{j: \rho_j > \rho_i} d_{ij} & \text{if } \exists j \text{ s.t. } \rho_j > \rho_i, \\ \max_j d_{ij} & \text{otherwise.} \end{cases} \quad (2)$$

130 Here, d_{ij} is the Euclidean distance between points x_i and x_j , and d_c refers to the cutoff
 131 distance, which is a given hyper-parameter. The original paper of DPC [15] suggests that
 132 d_c should be selected to cover 1-2% of the total number of points in data sets.

Alternatively, ρ_i can also be defined using a Gaussian kernel as follows:

$$\rho_i = \sum_{i \neq j} \exp \left[- \left(\frac{d_{ij}}{d_c} \right)^2 \right]. \quad (3)$$

133 After obtaining ρ_i and δ_i , there are two commonly used methods to identify C cluster
 134 centers, where C represents the number of target clusters. The first method is to draw a
 135 decision graph with δ_i and ρ_i as the x -axis and the y -axis, respectively, and then to manually
 136 select C cluster centers. The second method is to calculate $\gamma_i = \rho_i \times \delta_i$ and select the C
 137 data points with the first C largest values of γ_i as the cluster centers.

138 Once cluster centers are identified, DPC allocates each center to a cluster, and then each
 139 of the other data points is allocated to the cluster that its nearest neighbor with higher
 140 density belongs to.

141 The assumptions and strategies of DPC have been shown to be effective in many cases.
 142 However, DPC fails to find correct cluster centers sometimes. For example, when data have
 143 various densities [19], data points with highest γ_i could lie in the same ground-truth cluster
 144 but their distance could be large enough, which makes DPC choose cluster centers in the
 145 same ground-truth cluster and then fail with any allocation strategy. Also, since DPC is
 146 based on Euclidean distance, it performs poorly on data sets with complex structures [19],
 147 e.g., as illustrated in experiments of this article, it is not suitable for data sets with manifold
 148 structure. Finally, the choice of its parameter d_c , the cutoff distance, is also not an easy
 149 task.

150 3.2. Spectral Clustering

151 Spectral clustering is another widely used clustering algorithm. The idea of spectral
 152 clustering is: 1) to construct a connection graph with data points as nodes and their sim-
 153 ilarity measures as weights of edges; 2) to optimize the division of the connection graph
 154 into subgraphs s.t. the connection (similarity) between subgraphs is low and the connection
 155 of nodes within each subgraph is high; 3) to find an approximate solution to the former

156 optimization problem and obtain a new representation of data; 4) to perform clustering on
 157 the new representation.

Given a data set $X = \{x_1, x_2, \dots, x_m\}$, where m is the number of data points and the dimension of each data point is n , an adjacency matrix W is used to represent the connection graph, where W_{ij} represents the similarity value between x_i and x_j . A commonly used loss function in spectral clustering is the Ncut function [36]. Suppose C is the number of target clusters. Let $\text{cut}(A, \bar{A}) = \frac{1}{2} \sum_{x_i \in A, x_j \in \bar{A}} W_{ij}$ (where $\bar{A} = X \setminus A$) and $\text{vol}(A) = \sum_{x_i \in A, x_j \in X} W_{ij}$, for some $A \subseteq X$, then

$$\text{Ncut}(A_1, \dots, A_C) = \sum_{i=1}^C \frac{\text{cut}(A_i, \bar{A}_i)}{\text{vol}(A_i)}. \quad (4)$$

Finding optimal A_1, \dots, A_C for this objective function is an NP-hard problem and can be approximated by a relaxed problem, as follows. Let $h_j = (h_{1j}, h_{2j}, \dots, h_{mj})^T$, where

$$h_{ij} = \begin{cases} 0 & x_i \notin A_j \\ \frac{1}{\sqrt{\text{vol}(A_j)}} & x_i \in A_j \end{cases}. \quad (5)$$

Denote by D the diagonal matrix s.t. $D_{ii} = \sum_j W_{ij}$ and by L the Laplacian matrix $L = D - W$, then

$$\begin{aligned} h_i^T L h_i &= \frac{1}{2} \sum_{u,v} W_{ij} (h_{ui} - h_{vi})^2 \\ &= \frac{1}{2} \sum_{x_u \in A_i, x_v \in \bar{A}_i} W_{uv} \frac{1}{\text{vol}(A_i)} + \frac{1}{2} \sum_{x_u \in \bar{A}_i, x_v \in A_i} W_{uv} \frac{1}{\text{vol}(A_i)} \\ &= \frac{1}{2} \text{cut}(A_i, \bar{A}_i) \frac{1}{\text{vol}(A_i)} + \frac{1}{2} \text{cut}(\bar{A}_i, A_i) \frac{1}{\text{vol}(A_i)} \\ &= \frac{\text{cut}(A_i, \bar{A}_i)}{\text{vol}(A_i)}. \end{aligned} \quad (6)$$

In this way, the problem of minimizing the Ncut loss function is converted to:

$$\min_F \text{tr}(F^T D^{-1/2} L D^{-1/2} F) \quad \text{s.t.} \quad F^T F = I, \quad (7)$$

158 where $F = D^{1/2} H$ and the matrix $H = (h_1, \dots, h_C)$.

159 This optimization problem is a typical Rayleigh quotient problem [37] and the opti-
 160 mal F consists of the C eigenvectors corresponding to the first C smallest eigenvalues of
 161 $D^{-1/2} L D^{-1/2}$. The optimal solution H to the original problem can be obtained by cal-
 162 culating $H = D^{-1/2} F$. The final clustering result is obtained by K-means clustering on
 163 H , taking each row of H as a feature vector. Another commonly used method comes
 164 from Ng et al. [3], where instead of calculating H , it just standardizes F by row obtaining
 165 $Y_{ij} = F_{ij} / (\sum_{j'} F_{ij'}^2)^{1/2}$. Again, to obtain the final clustering result, K-means clustering is
 166 conducted on Y , taking each row of Y as a feature vector. This article will also adopt the
 167 latter strategy.

168 Spectral clustering does not make any assumptions about the structure of the data and,
 169 as a result, has good performance on data of various structures. However, spectral clustering
 170 has several significant problems: 1) it needs to calculate eigenvectors of the matrix and the
 171 time complexity is $\mathcal{O}(m^3)$ which is too high for large data sets; 2) a good similarity measure
 172 to form the adjacency matrix is crucial for its performance.

173 3.3. Analysis

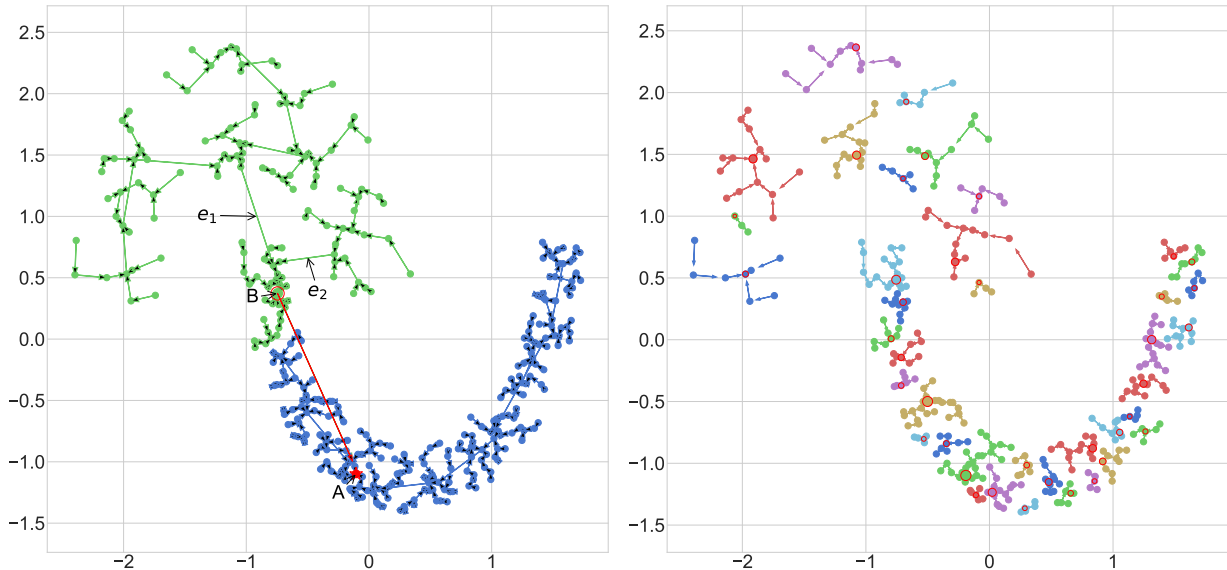


Figure 1: DPC clustering (left) and results with family trees (right).

174 DPC and many of its variants identify cluster centers based on the criteria that the
 175 centers should be points with globally higher density that are far apart from each other
 176 in terms of Euclidean distance. As densities could vary significantly across the whole data
 177 set and the data set might have a complex structure, this criteria is sometimes problematic.
 178 Consider Fig. 1 for example, where an arrow points from a data point to its nearest neighbor
 179 with higher density. On the left of the figure, point *A* is identified as the first cluster center
 180 by DPC; point *B* is a point with high density and is far enough from *A*, and thus DPC
 181 identifies it as the second cluster center, which is wrong. Moreover, as DPC uses Euclidean
 182 distance to measure the distance between a data point and its nearest neighbor with higher
 183 density, DPC allocates data points from the other cluster to the same cluster of *B*, e.g.,
 184 through arrows e_1 and e_2 . Both of the above make DPC have a poor clustering performance
 185 on this data set.

186 However, it is interesting to note that the strategy of DPC can have better clustering
 187 performance when restricted to local neighborhood. To illustrate this, we vary the number
 188 of target clusters from 2 to 29 for each of the four data sets, and measure the correctness
 189 of clusters by *purity*. The purity of a clustering result on a data set is defined as the
 190 $1 - \sum_i \frac{|c_i|}{m} \text{entropy}(c_i)$, where m is the number of data points, c_i stands for the i -th cluster in

191 the result, and $\text{entropy}(c_i)$ represents the information entropy of the real label distribution
 192 corresponding to i -th cluster. This measures the proportion of data points that are correctly
 193 clustered together. Fig. 2 shows the results. With the increase of the number of target
 194 clusters, or in other words, the decrease of the range of local neighborhood, the purity of
 195 clustering results increases constantly. This means that DPC is likely to correctly cluster
 data points when restricted to a local neighborhood.

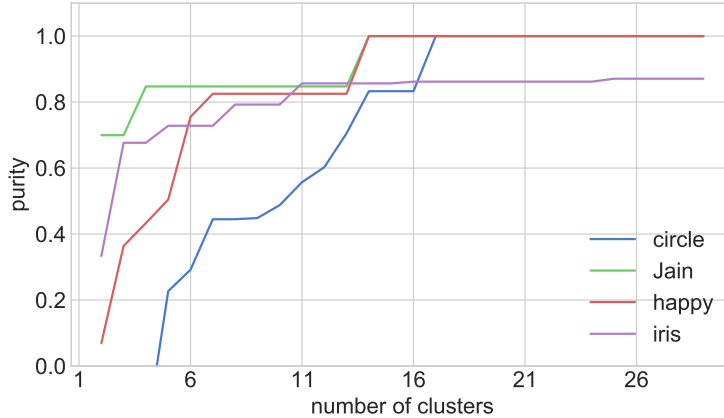


Figure 2: Purity change of DPC with respect to the number of target clusters.

196 The above observation inspires us to consider applying the idea of DPC in local proximity
 197 of data points with some necessary modifications, to first partition a data set into smaller
 198 local clusters. Here, we consider finding the nearest neighbor with higher density for a data
 199 point only from its k nearest neighbors, and if none is found then we treat this data point
 200 as a local density peak. The figure on the right of Fig. 1 shows the result of this idea. In
 201 the figure, the data set is partitioned into smaller local clusters (called *family trees*), where
 202 points with locally highest density are marked by circles with red boundary. Note that the
 203 incorrect inter-cluster connections disappear in the result.

204 The next important question is then how to combine these intermediate partitions into
 205 final target clusters. To this end, we exploit the idea of spectral clustering on family trees.
 206 In particular, we devise a novel similarity measure for these family trees to construct the
 207 adjacency matrix for spectral clustering. Unlike the original spectral clustering algorithm,
 208 this similarity measure also takes into account the contribution of the structure of family
 209 trees. In addition, we adapt the Ncut loss function for these family trees. As a result, we
 210 propose an algorithm combining the idea of local density peaks and spectral clustering (LDP-
 211 SC), and this algorithm retains several advantages of both DPC and spectral clustering.
 212

213 4. LDP-SC Algorithm

214 The LDP-SC algorithm consists of two stages: constructing family trees so that all the
 215 data points are *locally* clustered, and clustering of family trees so that all the data points
 216 are *globally* clustered.

217 *4.1. Construction of Family Trees*

218 As we have mentioned in the previous section, DPC performs well in identifying local
 219 clusters. The reason is that for local clusters, the sample space can be approximately
 220 regarded as an Euclidean space, which is exactly what DPC is good at. Therefore, instead
 221 of obtaining the final clustering result directly by finding cluster centers and allocating other
 222 points, we propose to first construct local clusters in the form of tree structures, known as
 223 *family trees*.

The idea to construct family trees for a data set is to find the parent for each data point in its local proximity in the sense of k nearest neighbors, and to identify the ones without parent as the root nodes of trees. To be more specific, let $\text{KNN}(x_i)$ denote the set of k nearest neighbors of a data point x_i (excluding x_i itself). The density of a data point x_i , denoted by $\rho(x_i)$, is given by

$$\rho(x_i) = \sum_{j=1}^k \exp(-\text{dist}_{ij}^2), \quad (8)$$

224 where dist is an $N \times k$ matrix and dist_{ij} is the Euclidean distance between x_i and its j -th
 225 nearest neighbor. In this article, dist is also normalized by $\text{dist} \leftarrow \text{dist} / \max(\text{dist})$ to reduce
 226 rounding errors. Compared with the definition of density in DPC (cf. Eq. 3), by exploiting
 227 k nearest neighbors, we avoid the cutoff distance d_c that is hard to determine and sensitive
 228 to specific distributions of data.

229 The *parent* $P(x_i)$ of a data point x_i can then be defined and the family trees of a data
 230 set will be automatically constructed with this information for all x_i .

Definition 1. The *parent* $P(x_i)$ is determined by

$$P(x_i) = \begin{cases} \arg \min_{x_j \in \text{higher}(x_i)} \text{dist}_{ij} & \text{if } \text{higher}(x_i) \neq \emptyset \\ \text{None} & \text{otherwise} \end{cases}, \quad (9)$$

231 where $\text{higher}(x_i) = \{x_j | x_j \in \text{KNN}(x_i), \rho(x_j) > \rho(x_i)\}$.

232 In other words, if there are points with higher density in the k nearest neighbors of x_i , then
 233 the nearest point x_j is selected as the *parent* of x_i (if there are multiple such x_j , we will
 234 select the first met one); otherwise, x_i has no parent ($P(x_i) = \text{None}$), and x_i is defined as
 235 a root and added to the set root . For any root $r_i \in \text{root}$, we will use $T(r_i)$ to refer to the
 236 family tree with r_i as its root.

237 **Definition 2 (Family tree).** A *family tree* is a tree with some data points x_i as its nodes,
 238 and has a directed edge (x_i, x_j) if $x_j = P(x_i)$.

239 Fig. 3 is an illustration of the family trees constructed for the data set Aggregation (see
 240 Section 5 for details) with the number of nearest neighbors $k = 20$. The data points with
 241 red boundary are root nodes. This data set has 788 data points while there are only 17
 242 family trees. From the figure, it is also worth noting that the local clusters given by these

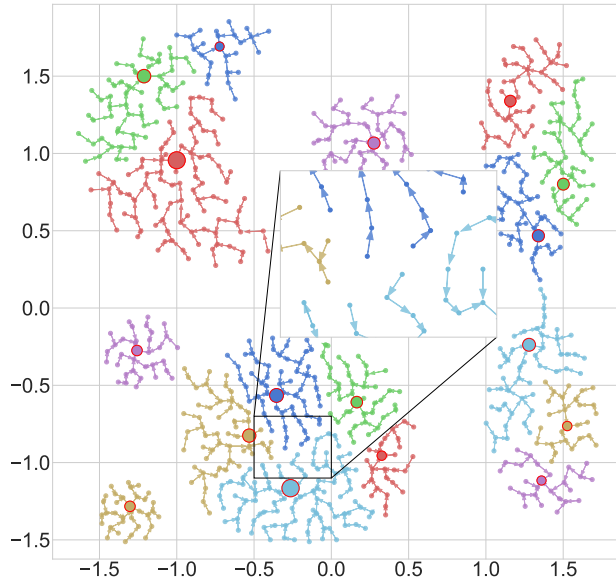


Figure 3: Illustration of family trees.

243 family trees have high purity, that is, the data points in the same tree are usually also in
 244 the same ground-truth cluster. A process to construct family trees for a data set is given in
 245 Algorithm 1.

246 Finally, the following theorem characterizes the correctness of Algorithm 1 and the rep-
 247 resentation power of family trees.

248 **Theorem 1.** *Given k the number of nearest neighbors and X a finite data set to be consid-
 249 ered, the returned P by Algorithm 1 forms a set of family trees as defined in Definition 2,
 250 and each of the data points in X is in exactly one of the family trees.*

251 *Proof.* To show that P returned by Algorithm 1 forms a set of family trees, it is enough to
 252 show that it forms a set of trees and each tree is a rooted tree. In fact, there is no cycle
 253 path in the graph induced by the parent-child relation P , so it forms a set of trees. Assume
 254 on the contrary that there is a cycle path $x_{i_1}, \dots, x_{i_s}, x_{i_1}$. Then by the construction of P ,
 255 $\rho(x_{i_j}) > \rho(x_{i_{j+1}})$ for $1 \leq j \leq s$, which implies $\rho(x_{i_1}) < \rho(x_{i_s})$, and $\rho(x_{i_s}) < \rho(x_{i_1})$, which is
 256 a contradiction. In addition, each data point x_i has either 1 or 0 out-edge, and each tree
 257 has exactly one x_i with 0 out-edge, because otherwise there will be a data point having at
 258 least 2 out-edges. This means each tree is a rooted tree.

259 To see that each data point x_i is in exactly one of the family trees, note that each family
 260 tree is a connected component if the edges are taken as undirected. If a data point is in two
 261 trees, then these two trees will be in a same connected component, which is not possible. \square

262 In other words, any finite data set can be represented as a unique set of family trees and
 263 such a set forms a partition of the data set. In practice, as we will see from the experimental
 264 results, this kind of partitions conform well with the ground-truth clusters of a data set, and
 265 the number of trees is much smaller than the number of data points. To obtain the final

Algorithm 1: FamilyTree

Input: A data set X ; the number of nearest neighbors k .

Output: Family trees represented by a parent-child relation P and a set of root nodes $root$.

```
1 foreach  $x_i \in X$  do
2   | Calculate  $KNN(x_i)$  and the distance matrix  $dist$ ;
3   | Calculate  $\rho(x_i)$ ;
4  $root \leftarrow \emptyset$ ;
5 foreach  $x_i \in X$  do
6   | Compute  $P(x_i)$  according to Eq. 9;
7   | if  $P(x_i) = \text{None}$  then
8     |    $root \leftarrow root \cup \{x_i\}$ ;
9 return  $P, root$ .
```

266 clustering result, we then need to combine the trees into larger clusters in a way it considers
267 both the root nodes and the structure of the trees, which will be discussed subsequently.

268 4.2. Family Tree Aggregation

269 In order to aggregate family trees into final clusters, it is necessary to consider not only
270 the roots of trees but also the characteristics of the tree structure. To this end, the strategy
271 of graph cut through spectral clustering is adopted, by considering each tree as a new data
272 point, and devising a novel similarity measure between trees and a more appropriate loss
273 function. The general steps of tree aggregation involve constructing an adjacency matrix,
274 conducting graph cut to extract abstract features for each tree, and finally performing a K-
275 means clustering on the new features to get the aggregation result of trees. In the following,
276 we will explain each of the steps in details.

277 4.2.1. Similarity Measure for Family Trees

278 An adjacency matrix for family trees represents the similarity between each pair of them.
279 A trivial solution is to calculate the distance between each pair of roots, which however,
280 cannot correctly reflects the similarity between trees, as two roots might be close to each
281 other but the trees as a whole are not. Therefore, we propose a more sophisticated similarity
282 measure for trees.

283 First, the *average distance* between trees is worth considering, which can reflect the
284 relative distance between trees to some extent.

Definition 3. The *average l -distance* between two trees $T(r_i)$ and $T(r_j)$ is

$$d_l(T(r_i), T(r_j)) = \frac{\sum_{t=1}^l \delta_t}{l}, \quad (10)$$

285 where $(\delta_1, \dots, \delta_s)$ is an ascendingly sorted list of all $d(a, b)$ for $a \in T(r_i)$ and $b \in T(r_j)$,
286 $d(a, b) = \|a - b\|_2$, l is a hyper-parameter, and if $s < l$ then l is directly set as s .

287 The structural information of a tree itself can also be interesting, and here we consider
 288 a measure that reflects the looseness of the distribution of data points in a tree. The reason
 289 is that when considering the distance between two trees, we need to take into account the
 290 internal distance for data points in each of the trees. For example, a distance δ between
 291 two trees in a case where the internal distance for them is much smaller than δ would be
 292 more significant than a case where that is larger than or similar to δ . Consider Fig. 4 as
 293 an example. Although the gap between the tree A and the tree B is larger than the gap
 between A and C , intuitively A is more similar to B than to C .

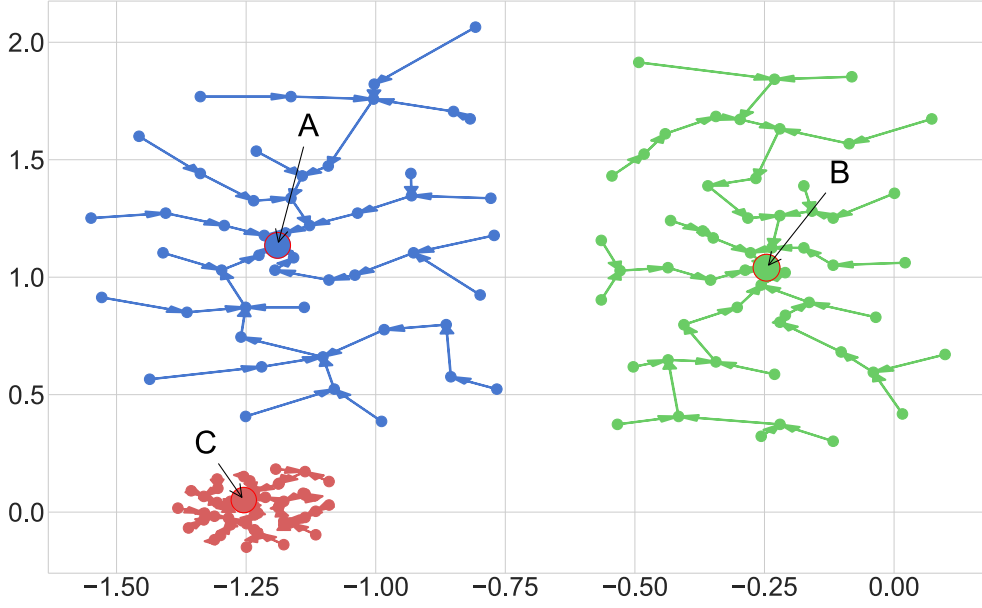


Figure 4: Illustration of the effect of looseness of trees on their similarities.

294

Definition 4. The *looseness* of a tree $T(r_i)$ is measured by

$$\iota(T(r_i)) = \sum_{x_j \in T(r_i)} g(x_j, r_i) \cdot \omega(x_j, r_i), \quad (11)$$

295 where $g(x_j, r) = 0$ if $T(r) = \{x_j\}$, and $g(x_j, r) = \min\{d(x_j, x_{j'}) | x_{j'} \neq x_j, x_{j'} \in T(r)\}$
 296 otherwise; $\omega(x_j, r) = \rho(x_j) / \sum_{x_{j'} \in T(r)} \rho(x_{j'})$.

297 This measure of looseness can be considered as the sum of the distances between x_j and
 298 its nearest neighbor in the tree for all x_j , weighted by the relative density of x_j in the tree.

299 Then, a measure of *separation* between two trees can be proposed as follows.

300 **Definition 5.** The *separation* between two trees $T(r_i)$ and $T(r_j)$ can be measured by
 301 $\sigma(T(r_i), T(r_j)) = d_l(T(r_i), T(r_j)) / \sqrt{\iota(T(r_i)) \cdot \iota(T(r_j))} + \epsilon$, where ϵ is a small number to
 302 avoid division by zero.

303 It is easy to verify in Fig. 4 that this separation measure conforms to intuition. In fact,
 304 when calculating the separation degrees for trees A , B and C , we get $d_l(A, B) > d_l(A, C)$
 305 and $\iota(A) \approx \iota(B) \gg \iota(C)$, which leads to $\sigma(A, B) < \sigma(A, C)$.

306 Finally, we would also consider the number of *shared nearest neighbors* of two trees as
 307 a contribution to the similarity measure of trees, which can be helpful for processing data
 308 with manifold structure.

309 **Definition 6.** The *shared nearest neighbors* of trees $T(r_i)$ and $T(r_j)$ $\text{SNN}(T(r_i), T(r_j)) =$
 310 $A(T(r_i)) \cap A(T(r_j))$, where $A(T(r_i))$ is the *affiliate set* of $T(r_i)$, which is $\bigcup_{x_j \in T(r_i)} \text{KNN}(x_j)$,
 311 i.e., the set of k nearest neighbors of data points in the tree.

312 Note that this definition of shared nearest neighbors is generalized from [19]. Particularly,
 313 the above definition is about the shared nearest neighbors of trees, instead of points as in
 314 [19]. Fig. 5 shows the shared nearest neighbor connections with the number of nearest
 315 neighbors being 5 on the Jain data set, where a red line connects two roots r_i and r_j if
 316 $|\text{SNN}(T(r_i), T(r_j))| > 0$. It can be seen from this example that this measure respects the
 manifold structure of the data.

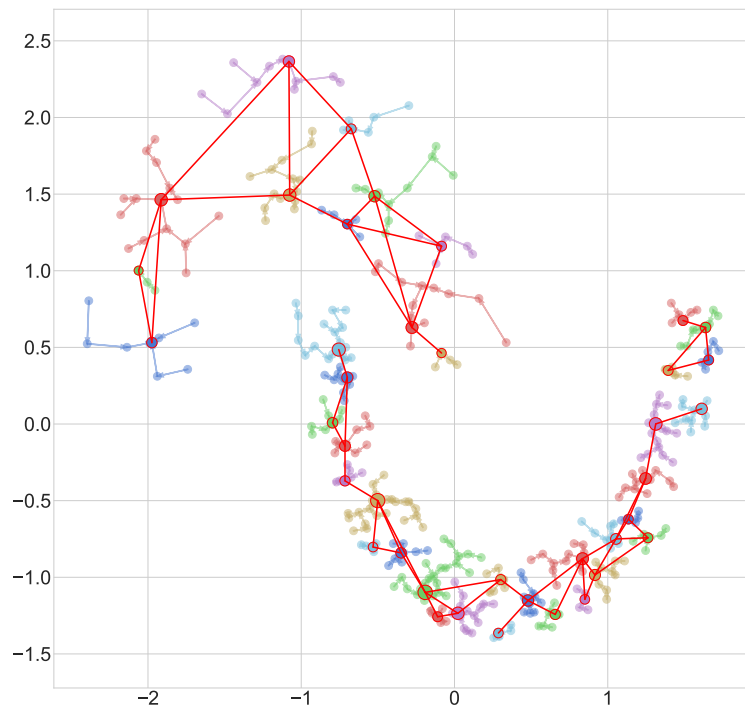


Figure 5: Illustration of shared nearest neighbor connections between trees.

317 As a result, the *similarity* between two trees can be obtained by synthesizing all of the
 318 above.
 319

Definition 7. The *similarity* between two trees $T(r_i)$ and $T(r_j)$ is given by

$$\text{sim}(T(r_i), T(r_j)) = \frac{|\text{SNN}(T(r_i), T(r_j))|}{(1 + \sigma(T(r_i), T(r_j))) \cdot (1 + d(r_i, r_j))}, \quad (12)$$

320 where $d(r_i, r_j) = \|r_i - r_j\|_2$.

321 The above similarity measure for two trees is determined together by the number of
 322 shared nearest neighbors, the separation measure between them, and the distance between
 323 roots. Although this measure makes use of shared nearest neighbors, it is significantly
 324 different from the similarity measure defined in [19]. Specifically, in addition to the difference
 325 in the use of shared nearest neighbors, we also take the distance between roots and the
 326 separation between trees into consideration. Moreover, the above similarity measure will
 327 be used for graph cut, whereas that in [19] is used for the calculation of densities, as a
 328 subroutine for a procedure similar to DPC.

329 For two trees $T(r_i)$ and $T(r_j)$, when their number of shared nearest neighbors is non-zero
 330 and fixed, if the densities of them have large difference (i.e., the separation degree is high),
 331 or if the distance between the roots is large, then the similarity will be low. Otherwise,
 332 the separation degree $\sigma(T(r_i), T(r_j))$ will be a balance against the distance between roots
 333 $d(r_i, r_j)$: even when the distance is relatively large, if the separation degree is relatively low,
 334 then the similarity measure can still be relatively large. If both the separation degree and the
 335 distance between roots are close to zero (i.e., two trees are very similar), then the similarity
 336 value will be close to $|\text{SNN}(T(r_i), T(r_j))|$. Finally, when two trees have similar densities and
 337 similar distances between roots, then the similarity value will be mainly determined by the
 338 number of shared nearest neighbors. It can be easily seen that the value of this similarity
 339 measure is between 0 and m (m is the number of data points in a data set).

340 The adjacency matrix W for the family trees of a data set is then calculated as $W_{ij} =$
 341 $\text{sim}(T(r_i), T(r_j))$ and normalized by $W_{ij} = W_{ij} / \max(W)$. Note that after normalization,
 342 $W_{ij} \leq 1$, and that the similarity between a tree and itself should be the highest. We thus
 343 set each W_{ii} to be 1. It is worth noting that the connection graph induced by the adjacency
 344 matrix (i.e., an edge exists if the corresponding entry in the matrix is nonzero) might not
 345 be connected (see Fig. 5 for example), and the number of connected components could be
 346 larger than the number of target clusters, which is not desirable [11]. Therefore, we construct
 347 a Gaussian kernel distance matrix W' between roots, which is $W'_{ij} = \exp(-\frac{d(r_i, r_j)^2}{2})$, and
 348 update W to be $W + \theta W'$. Because W' is added only to maintain the connectivity of the
 349 whole graph, usually θ is set to a small value, e.g., 0.001.

350 4.2.2. Graph Cut and Clustering of Trees

The idea to perform graph cut on the connection graph of trees is similar to that of
 spectral clustering. To simplify notations, in the connection graph, we use root r to represent
 the tree $T(r)$. Suppose \mathcal{T} is the set of family trees ($|\mathcal{T}| = p$) and the number of target clusters
 is C . We adapt the original Ncut loss function (Eq. 4) for trees as follows.

$$\text{Ncut}(A_1, \dots, A_C) = \sum_{i=1}^C \frac{\text{cut}(A_i, \bar{A}_i)}{\text{vol}(A_i) \cdot \sum_{r \in A_i} |T(r)|}, \quad (13)$$

351 where $A_i = \{r_{i_1}, \dots, r_{i_s}\}$ is a set of family trees ($\bigcup_i A_i = \mathcal{T}$), $\text{cut}(A_i, \bar{A}_i) = \frac{1}{2} \sum_{r_u \in A_i, r_v \in \bar{A}_i} W_{uv}$,
 352 and $\text{vol}(A_i) = \sum_{r_u \in A_i, r_j \in \mathcal{T}} W_{uj}$. Note that an important difference between the above equa-

tion and the original Ncut loss function is that we also take the size of trees into consideration, in order to balance the sizes of the final clusters of data points.

To find a solution that minimizes the above loss function, let $h_j = (h_{1j}, h_{2j}, \dots, h_{pj})^T$, where

$$h_{ij} = \begin{cases} 0 & v_i \notin A_j \\ \frac{1}{\sqrt{\text{vol}(A_j) \cdot \sum_{r \in A_j} |T(r)|}} & v_i \in A_j \end{cases} . \quad (14)$$

Write $U = DD^*$, where D^* and D are diagonal matrices, $D_{ii}^* = |T(r_i)|$, and $D_{ii} = \sum_j w_{ij}$. Similar to the derivation of Eq. 6, the final optimization problem becomes

$$\min_H \text{tr}(H^T L H) \quad \text{s.t.} \quad H^T U H = I, \quad (15)$$

where $H = (h_1, \dots, h_C)$ and $L = D - W$ is the Laplacian matrix. Again, we can convert the above problem to the following one, by setting $F = U^{1/2} H$.

$$\min_F \text{tr}(F^T U^{-1/2} L U^{-1/2} F) \quad \text{s.t.} \quad F^T F = I. \quad (16)$$

The above problem can be solved by computing the eigenvectors t_1, \dots, t_C corresponding to the first C smallest eigenvalues of $U^{-1/2} L U^{-1/2}$.

Following a similar fashion of spectral clustering, let $E = (t_1, \dots, t_C) \in \mathbb{R}^{p \times C}$ and $Y = (Y_{ij})_{p \times C}$ s.t. $Y_{ij} = E_{ij} / (\sum_{s=1}^C E_{is}^2)^{1/2}$. We can take each row of Y as a feature vector of each family tree, and the clustering result of the family trees can then be obtained by applying the K-means clustering algorithm on these feature vectors.

The effectiveness of the revised objective function is visually illustrated in Fig. 6. In these figures, the clustering results of several trees are shown, where the left figure is generated using the original Ncut loss function and the right is from the revised loss function for trees. In the results, each circle represents a tree and the size of a circle is proportional to the size of the tree. It can be seen that with the original Ncut function, there are three trees incorrectly clustered (tagged as ‘‘error’’), while with the revised function the three errors are corrected. By looking closer, the incorrectly clustered tree at the top has only a few nodes and has relatively low connection weights with other trees, so it is identified as an individual cluster; the other two incorrectly clustered trees at the bottom should be belong to the same cluster but are identified as parts of two separated clusters, because they have relatively high connection weights with other trees; the sizes of the trees are never considered as a decision factor in the original Ncut function.

4.3. The LDP-SC Algorithm and Its Time Complexity

The LDP-SC algorithm is shown in Algorithm 2. The main steps include: 1) constructing family trees by Algorithm 1; 2) calculating the adjacency matrix W for graph cut; 3) solving the graph cut problem to extract new features; 4) conducting K-means clustering on the new features to obtain the final clustering result.

Suppose the number of data points in a data set is m , the dimension of each data point is n , the number of target clusters is C , the number of nearest neighbors is k , and the number of family trees is p .

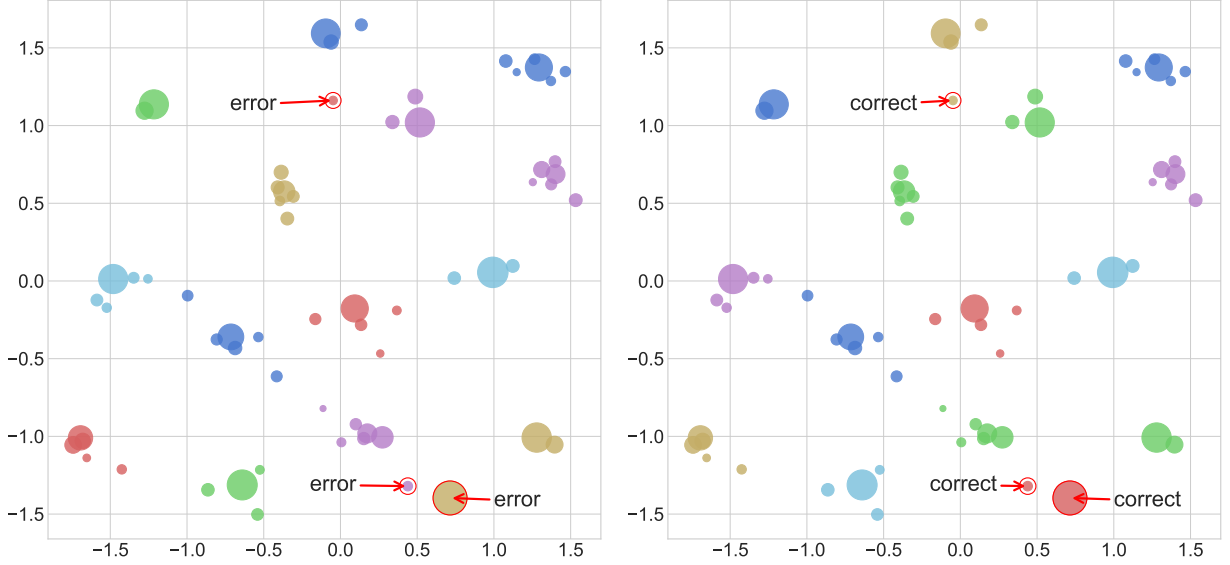


Figure 6: Illustration of the effectiveness of the revised objective function (right), compared with the original loss function (left).

381 In Algorithm 1, to get the distance matrix **dist** for the distances between points and their
382 k nearest neighbors, one can make use of a KD-tree and the time complexity is $\mathcal{O}((n +$
383 $k)m \log m)$. The computation of densities is only related to the k nearest neighbors of each
384 data point, so the time complexity is $\mathcal{O}(km)$. The family trees can then be constructed by
385 calculating $P(x_i)$ for all x_i in a total of time $\mathcal{O}(km)$, as to calculate $P(x_i)$ or to determine
386 root, one needs to check every k nearest neighbor of x_i . So Algorithm 1 requires $\mathcal{O}((n +$
387 $k)m \log m)$ time.

388 For the rest of the algorithm LDP-SC, in lines 2-3, when calculating the similarity mea-
389 sure, there are three parts in sequence.

390 For the part of $|\text{SNN}(T(r_i), T(r_j))|$ for all i, j , one needs first to get the k nearest neighbors
391 of each data point, which takes $\mathcal{O}(km)$ time in total; then the intersections of affiliate sets
392 of each pair of trees can be obtained in time $\mathcal{O}(\sum_{i=1}^p |T(r_i)|(p-1)) = \mathcal{O}(pm)$. Therefore,
393 for all $|\text{SNN}(T(r_i), T(r_j))|$, it costs $\mathcal{O}((k+p)m)$ time in total.

394 For the part of $\sigma(T(r_i), T(r_j))$ for all i, j , one needs $\mathcal{O}(n \cdot |T(r_i)| \cdot |T(r_j)|)$ time to get each
395 $d_l(T(r_i), T(r_j))$, by first calculating $d(a, b)$ for all $a \in T(r_i), b \in T(r_j)$ and then scanning
396 through the list of these values to get the smallest l values. So in total one needs $\mathcal{O}(\sum_{i,j} n \cdot$
397 $|T(r_i)| \cdot |T(r_j)|) = \mathcal{O}(n \cdot \sum_i |T(r_i)| \cdot m) = \mathcal{O}(nm^2)$ time to get all $d_l(T(r_i), T(r_j))$. To get
398 $\iota(T(r_i))$ for all i , it can be completed in $\mathcal{O}(km)$, as for each $\iota(T(r_i))$ one needs $\mathcal{O}(|T(r_i)|)$
399 time to get all $\omega(x_j, r_i)$ and $\mathcal{O}(m^2)$ time to get all $g(x_j, r_i)$. So, all $\sigma(T(r_i), T(r_j))$ can be
400 obtained in time $\mathcal{O}(nm^2)$.

401 The third part $d(r_i, r_j)$ for all i, j can be obtained in a total time of $\mathcal{O}(np^2)$.

402 Therefore, lines 2-3 in the algorithm LDP-SC take $\mathcal{O}((k+p)m + nm^2 + np^2)$ time.

403 Lines 4-6 take $\mathcal{O}(p^2)$ time.

404 Lines 7-14 of Algorithm 2 take the same time as a spectral clustering algorithm which is

405 $\mathcal{O}(p^3)$ [36].

406 Lines 15-17 can be completed in time $\mathcal{O}(m)$ as one only needs to go through all the data
407 points once.

Algorithm 2: LDP-SC

Input: A data set X ; the number of nearest neighbors k ; the number of clusters C .

Output: Clustering results label.

```
1  $P, \text{root} \leftarrow \text{FamilyTree}(X, k)$ ;  
2 foreach  $(r_i, r_j) \in \text{root} \times \text{root}$  and  $r_i \neq r_j$  do  
3   |  $W_{ij} \leftarrow \text{sim}(T(r_i), T(r_j))$  by Eq. 12;  
4  $W \leftarrow W / \max(W)$  and set each  $W_{ii}$  to 1;  
5 Compute the Gaussian kernel distance matrix  $W'$ ;  
6  $W \leftarrow W + \theta W'$ ;  
7 foreach  $r_i \in \text{root}$  do  $D_{ii} \leftarrow \sum_j W_{ij}$  ;  
8  $L \leftarrow D - W$ ;  
9 foreach  $r_i \in \text{root}$  do  $D_{ii}^* \leftarrow |T(r_i)|$  ;  
10  $U \leftarrow DD^*$ ;  
11  $E \leftarrow (t_1, \dots, t_C)$ , the  $C$  eigenvectors of  $U^{-1/2}LU^{-1/2}$  corresponding to its  $C$  smallest  
    eigenvalues;  
12 foreach  $(r_i, r_j) \in \text{root} \times \text{root}$  do  
13   |  $Y_{ij} \leftarrow E_{ij} / (\sum_{s=1}^C E_{is}^2)^{1/2}$ ;  
14 Apply K-means on  $Y$  to obtain cluster label  $c_i$  for each  $T(r_i)$ ;  
15 foreach  $r_i \in \text{root}$  do  
16   | foreach  $x_i \in T(r_i)$  do  
17     | |  $\text{label}(x_i) \leftarrow c_i$ ;  
18 return label.
```

408 In summary, the worst time complexity of the LDP-SC algorithm is $\mathcal{O}((n+k)m \log m +$
409 $(k+p)m + nm^2 + np^2 + p^3)$ and is bounded by $\mathcal{O}(km \log m + nm^2 + p^3)$. Usually, the
410 value of dimension n is considered to be bounded by a much smaller value than m and the
411 time complexity can be simplified as $\mathcal{O}(m^2 + p^3)$. A comparison of the time complexities of
412 different algorithms is shown in Table 1 (without considering dimension). In practice, the
413 time needed for LDP-SC is usually much lower than the worst case and it is actually quite
414 efficient. For example, LDP-SC takes about 3 seconds for data sets with about 9000 data
415 points (the Pendigits and USPS data sets), 160 seconds for data sets with about 70 000 data
416 points (MNIST and Fashion MNIST), and 3300 seconds for a data set with about 580 000
417 data points (coverttype). As an intuition on the efficiency of LDP-SC compared to others,
418 DPC needs about 140 seconds for Pendigits and USPS, and cannot finish in 5 hours for
419 MNIST, Fashion MNIST, and coverttype.

420 5. Empirical Evaluations

Table 1: Running time complexity comparison.

Algorithm	Date, Publication	Time Complexity
LDP-SC	Ours	$O(m^2 + p^3)$
SNN-DPC	2018, Information Sciences	$O(m^2)$
DPC	2014, Science	$O(m^2)$
LDP-MST	2021, IEEE Transactions on Knowledge and Data Engineering	$O(m \log m)$
FHC-LDP	2021, Neurocomputing	$O(m \log m)$
DPC-DBFN	2020, Pattern Recognition	$O(m^2)$
U-SPEC	2020, IEEE Transactions on Knowledge and Data Engineering	$O(mk(k + C) + p_1^3)^{**}$
LEC-K	2014, IEEE Transactions on Cybernetics	$O(p_1^3 + mp_1^2)$
SC	2000, IEEE Transactions on Pattern Analysis and Machine Intelligence	$O(m^3)$

* m is the number of data points for all the algorithms; p is the number of family trees ($p \ll m$).

** p_1 is the number of landmarks ($k \ll p_1 \ll m$).

421 5.1. Experiment Settings

422 To evaluate the performance of the LDC-SC algorithm, we will compare it with several
 423 baseline algorithms on both synthetic data sets and real-world data sets. Because the idea
 424 of the algorithm is based on local density peaks and graph cut through spectral clustering,
 425 the following prominent baseline algorithms are chosen:

- 426 1. DPC [15]: the original DPC algorithm.
- 427 2. SNN-DPC [19]: it improves density measurement, distance calculation, and allocation
 428 strategy of the DPC algorithm.
- 429 3. LDP-MST [30]: it improves the DPC algorithm by generating some local tree structures
 430 (minimum spanning trees) and constructing connection relations between these trees.
- 431 4. FHC-LDP [24]: it improves the DPC algorithm by building sub-clusters with local density
 432 structures, and then combining these sub-clusters by hierarchical clustering.
- 433 5. DPC-DBFN [27]: it improves the DPC algorithm by exploiting a fuzzy kernel to extract
 434 a specific kind of structural information to improve the robustness of clustering.
- 435 6. SC [36]: the original spectral clustering algorithm based on Ncut loss function.
- 436 7. U-SPEC [33]: an accelerated spectral clustering algorithm using approximated K-nearest
 437 neighbors and bipartite graphs.
- 438 8. LSC-K [32]: a variant of spectral clustering algorithms by using K-means to select land-
 439 marks.

440 The implementation of SC are from the scikit-learn library, the implementations of SNN-
 441 DPC, LDP-MST, FHC-LDP, and DPC-DBFN are from the source code provided by the
 442 authors, the implementations of DPC and LSC-K are based on the descriptions of the
 443 original paper, and the U-SPEC implementation is from [38].

444 Three commonly used clustering performance metrics, i.e., Adjusted Rand Index (ARI
 445 [39]), Normalized Mutual Information (NMI [40]) and Accuracy (ACC [41]) are selected for
 446 performance evaluation. All results are averaged over ten repeated runs of the algorithms.

447 For fair comparison, the hyper-parameters of each algorithm are tuned beforehand.
 448 Specifically, all algorithms use the number of ground-truth clusters as the number of target

449 clusters. The number k of nearest neighbors for LDP-SC, SNN-DPC, SC, U-SPEC, FHC-
450 LDP, DPC-DBFN, and LSC-K is searched from 2 to 50 with a step size of 1. DPC-DBFN
451 has two ways to calculate densities and we try both and select the best one. For DPC, we
452 search for an optimal value of the ratio of the number of neighbors against the number of all
453 data points, in a loop with a step size of 0.1 from 0.2 to 5. For identifying cluster centers of
454 SNN-DPC and DPC, this article selects the C data points with the first C largest values of
455 γ (cf. Section 3.1), which is more general and fairer compared with manual selection with a
456 decision graph. For LDP-MST, it recommends that data should have dimensionalities lower
457 than 10, so we search for an optimal dimensionality from 2 to $\min(50, n)$ by using PCA.
458 The adjacency matrix of SC is calculated based on k nearest neighbors. U-SPEC has a
459 parameter p and we set the value p as 1 000 if the number of data points $m > 10\,000$ and
460 as $\lfloor (m - 1)/10 \rfloor$ otherwise. For LSC-K, the parameter p is set to 500 if $m > 1000$ and to
461 $\lfloor m/2 \rfloor$ otherwise.

462 For LDP-SC, the parameter l for average l -distance is set as 4, the weight θ of Gaussian
463 kernel distance matrix is set as 0.001, and ϵ in Definition 5 is set as 10^{-6} .

464 For data sets, we utilize 8 widely used synthetic data sets, 8 real-world data sets, and 6
465 image data sets. The details of them are shown in Table 2, Table 3, and Table 4. Fig. 7 gives
466 an illustration of the image data sets. Data sets are all preprocessed in certain ways: for non-
467 image data sets, the SNN-DPC and FHC-LDP algorithms apply min-max standardization
468 to the data, and all the other algorithms applies z-score standardization; for image data
469 sets, if the pixel value is in range $[0, 255]$, then it is scaled to $[0, 1]$.

470 The experimental environment is: Python 3.9, Ubuntu 18.04, CPU i7-8700, and 64GB
471 RAM.

Table 2: Synthetic data sets.

Data set	#Instances	#Attributes	#Clusters
Aggregation	788	2	7
R15	600	2	15
S2	5000	2	15
Flame	240	2	2
Jain	373	2	2
Spiral	312	2	3
happy	266	2	3
circle	299	2	3

472 5.2. Results and Analysis on Synthetic Data Sets

473 Table 5 shows the results and the corresponding hyper-parameters, for LDP-SC and the
474 other eight comparison algorithms on the synthetic data sets. The best results are marked
475 in bold.

476 From Table 5, we can see that LDP-SC has the best or close to the best performance on
477 all of these data sets. On five out of the eight data sets, LDP-SC has 100 scores for all of
478 ARI, NMI, and ACC metrics. The cases where LDP-SC is slightly weaker than the optimal

Table 3: Real-world data sets.

Data set	#Instances	#Attributes	#Clusters
pengleukEW	72	7070	2
parkinsonsEW	195	22	2
seeds	210	7	3
spect	267	22	2
mfeat-fac	2000	216	10
landsatEW	6435	36	6
Pendigits	10992	16	10
convertype	581012	54	7

Table 4: Image data sets.

Data set	#Instances	#Attributes	#Clusters
MNIST	70000	784	10
Fashion MNIST	70000	784	10
COIL20	1440	1024	20
USPS	9298	256	10
Yale	2414	1024	38
ORL	400	4096	40

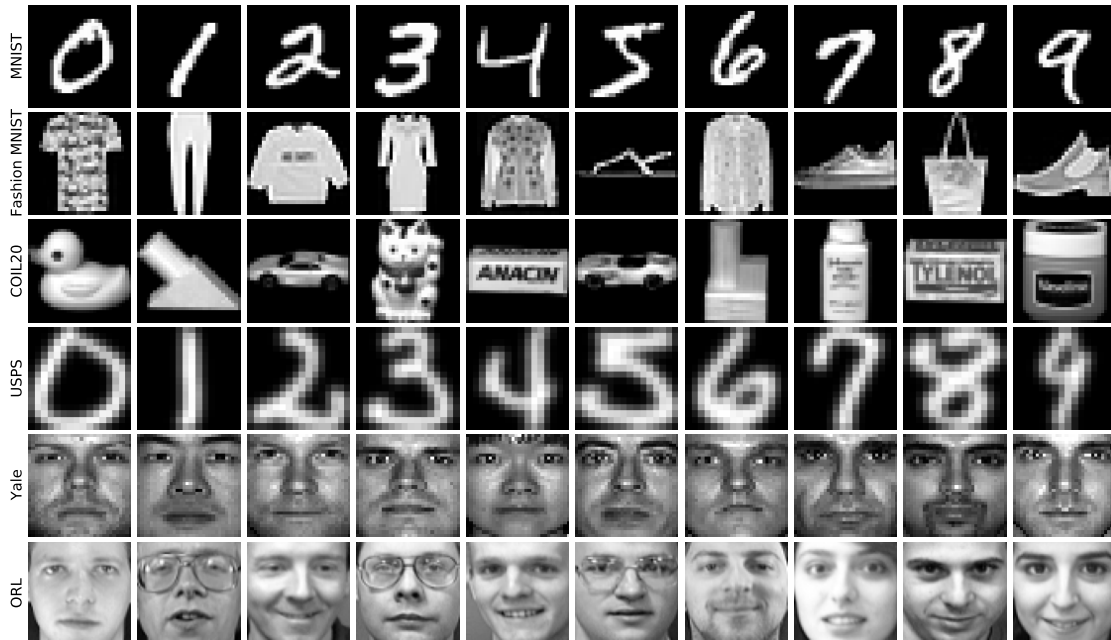


Figure 7: Examples of image data sets.

Table 5: Results on synthetic data set (%).

Data set		LDP-SC	SNN-DPC	DPC	LDP-MST	FHC-LDP	DPC-DBFN	U-SPEC	LEC-K	SC
Aggregation	ARI	99.56	95.94	99.56	99.56	100	99.27	90.57	99.20	98.98
	NMI	99.24	95.55	99.24	99.24	100	98.83	92.55	98.84	98.51
	ACC	99.75	97.84	99.75	99.75	100	99.62	85.41	99.62	99.49
	Par.	7	15	3.5	2	6	30	3	8	14
R15	ARI	99.28	99.28	99.28	98.91	99.28	99.64	98.91	99.28	99.28
	NMI	99.42	99.42	99.42	99.13	99.42	99.71	99.13	99.42	99.42
	ACC	99.67	99.67	99.67	99.50	99.67	99.83	99.50	99.67	99.67
	Par.	6	10	0.2	2	7	39	5	23	35
S2	ARI	92.94	91.91	92.49	91.88	92.86	93.61	93.57	93.52	93.41
	NMI	93.76	93.14	93.52	93.54	93.71	94.29	94.25	94.19	94.27
	ACC	96.58	96.08	96.36	96.04	96.54	96.92	96.90	96.88	96.82
	Par.	34	35	0.5	2	36	4	22	10	47
Flame	ARI	100	95.02	100	93.39	100	96.66	98.33	95.01	91.76
	NMI	100	89.94	100	87.52	100	92.69	96.35	89.91	85.19
	ACC	100	98.75	100	98.33	100	99.17	99.58	98.75	97.92
	Par.	23	5	2.7	2	7	6	4	5	5
Jain	ARI	100	59.35	71.46	100	100	80.08	100	100	100
	NMI	100	55.91	65.22	100	100	71.88	100	100	100
	ACC	100	88.74	92.49	100	100	95.17	100	100	100
	Par.	14	39	0.3	2	7	43	2	9	8
Spiral	ARI	100	100	100	100	100	14.87	2.64	100	100
	NMI	100	100	100	100	100	15.75	2.85	100	100
	ACC	100	100	100	100	100	55.45	44.87	100	100
	Par.	4	5	1.6	2	3	4	2	3	3
happy	ARI	100	100	71.79	100	100	52.21	100	100	100
	NMI	100	100	78.71	100	100	60.94	100	100	100
	ACC	100	100	89.85	100	100	81.58	100	100	100
	Par.	4	5	4.6	2	8	2	2	3	7
circle	ARI	100	100	21.41	100	100	46.06	100	100	100
	NMI	100	100	34.50	100	100	63.96	100	100	100
	ACC	100	100	56.86	100	100	73.91	100	100	100
	Par.	10	33	0.2	2	7	17	2	9	17

479 ones (FHC-LDP and DPC-DBFN) are the Aggregation, R15, and S2 data sets. It is worth
480 noting that on the data set Jain, the SNN-DPC algorithm performs much worse than in
481 the original paper, because in this article it identifies cluster centers automatically with γ
482 values, whereas in the original paper it was done manually.

483 We also visualize the results on synthetic data sets in Figs. 8 to 10. For the sake of space
484 limitation, we omit the visualization of the results of DPC, SC, and LSC-K. In all of these
485 figures, the cluster centers of SNN-DPC and DPC-DBFN are represented by pentagrams.

486 Figs. 8, 9, and 10 show the clustering results on three data sets with high density clusters.
487 SNN-DPC correctly finds the cluster centers and sometimes makes mistakes on borders due
488 to its allocation strategy. Note that U-SPEC sometimes puts two subsets of data points
489 that are far apart into one cluster. On the other hand, LDP-SC correctly separates the
490 boundaries and adapts well for clusters of complex shapes.

491 Figs. 11 to 15 show clustering results on five data sets with irregular shapes and manifold
492 structure. It can be seen that LDP-SC performs well in all cases. For Fig. 11, all the other
493 algorithms also perform well, except for some errors on the boundary. For Fig. 12, LDP-SC,
494 LDP-MST, FHC-LDP, and U-SPEC all correctly detect the manifold structure; SNN-DPC
495 and DPC-DBFN make mistakes on the boundary; the errors of SNN-DPC and DPC-DBFN

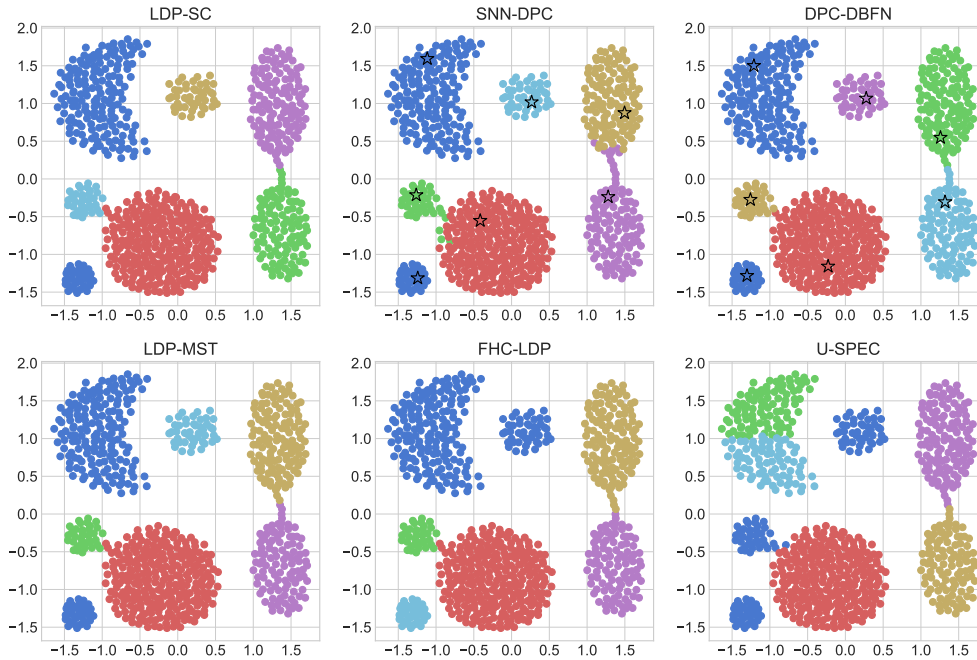


Figure 8: Aggregation

496 are due to large values for the number of nearest neighbors. For Fig. 13, LDP-SC, SNN-
 497 DPC, and LDP-MST perform well, while DPC-DBFN, FHC-LDP, and U-SPEC have poor
 498 performances on this data set. For Fig. 14, DPC-DBFN incorrectly includes data points
 499 from the curved strip for the clusters in the middle, and for Fig. 15, it identifies a wrong
 500 cluster center which makes it perform poorly.

501 5.3. Results and Analysis on Real-World Data Sets

502 Table 6 shows the results on eight commonly used real-world data sets, where the numbers
 503 of data points, clusters, and features vary in a large range, which allows us to evaluate our
 504 algorithm from different aspects. The proposed LDP-SC algorithm has the best performance
 505 in terms of all three metrics for five out of the eight data sets. Note that in four out of these
 506 five data sets, the performance of LDP-SC is significantly better than the second best. There
 507 are two data sets, i.e., spect and covertime, for which LDP-SC has the best scores for some
 508 of the three metrics. It is worth noting that for the covertime data set, the results of SNN-
 509 DPC, DPC, LDP-MST, FHC-LDP, DPC-DBFN, and SC cannot be obtained (shown as
 510 N/A in the table), because they exceed memory or time limit (5 hours). For the seeds data
 511 set, LDP-SC is slightly weaker than the best one. These results illustrate the superiority of
 512 LDP-SC on complex real-world data.

513 5.4. Results and Analysis on Image Data Sets

514 Table 7 shows the results on six commonly used image data sets. Note that for the
 515 MNIST and Fashion MNIST data sets, due to large number of data points (70,000), SNN-
 516 DPC, FHC-LDP, DPC-DBFN, and DPC exceed memory or time limit and have no results.

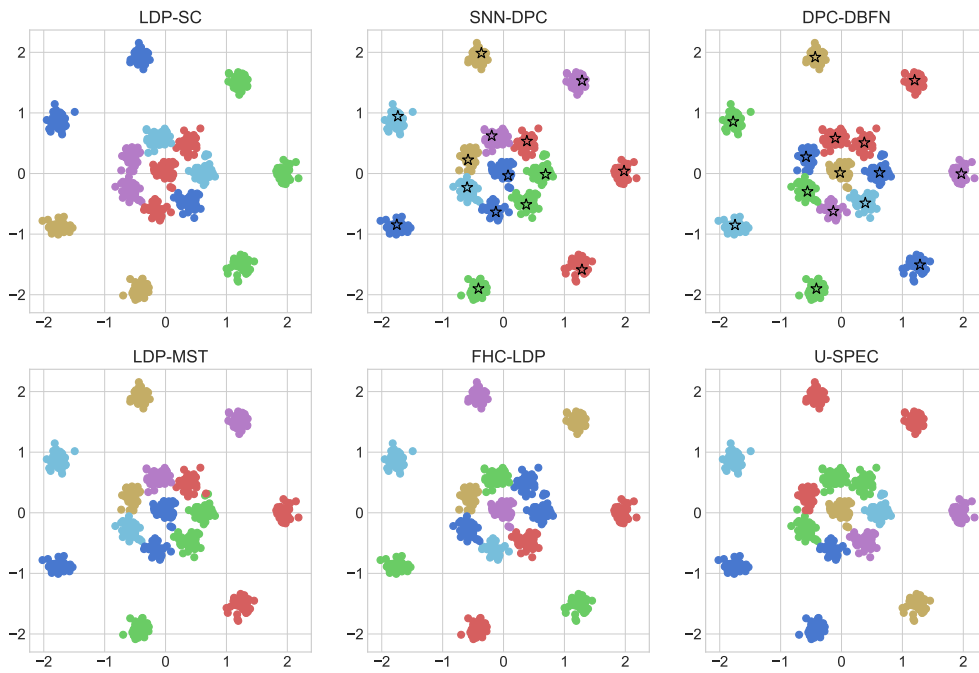


Figure 9: R15

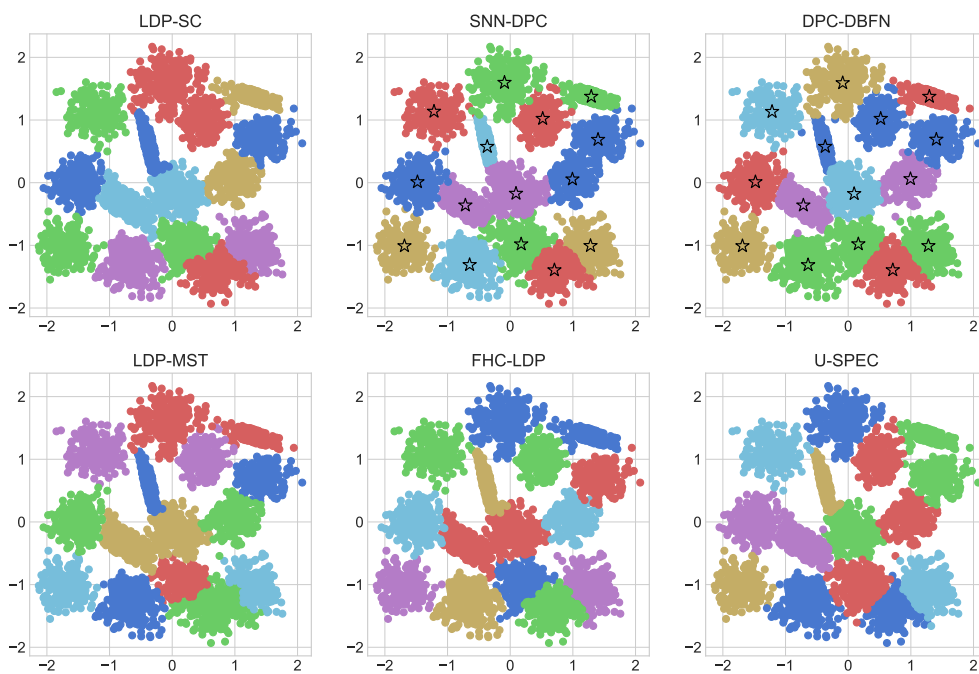


Figure 10: S2

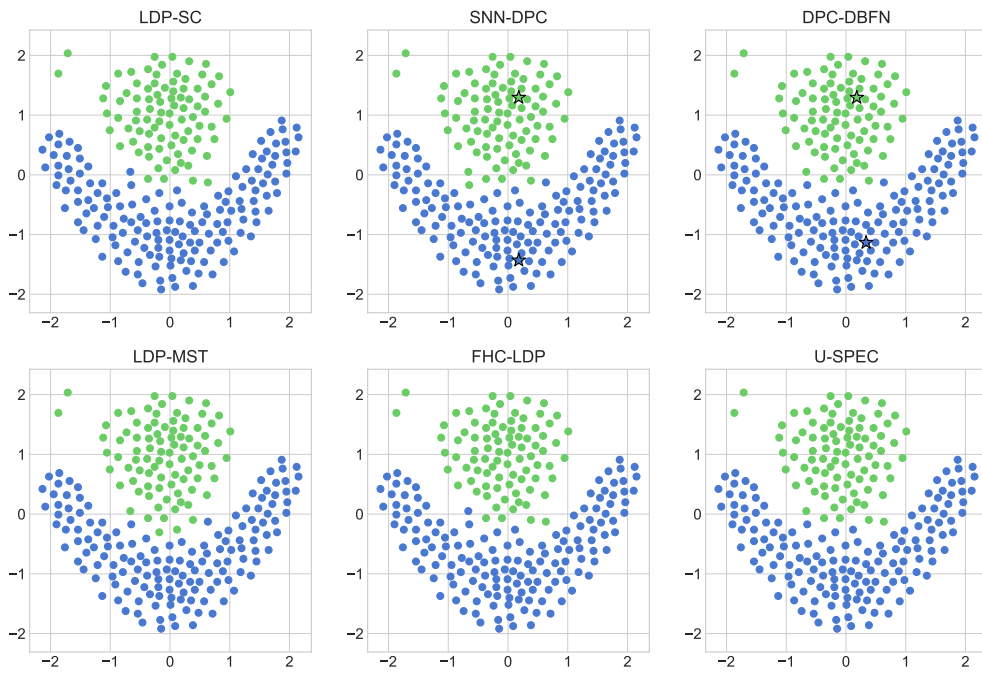


Figure 11: Flame

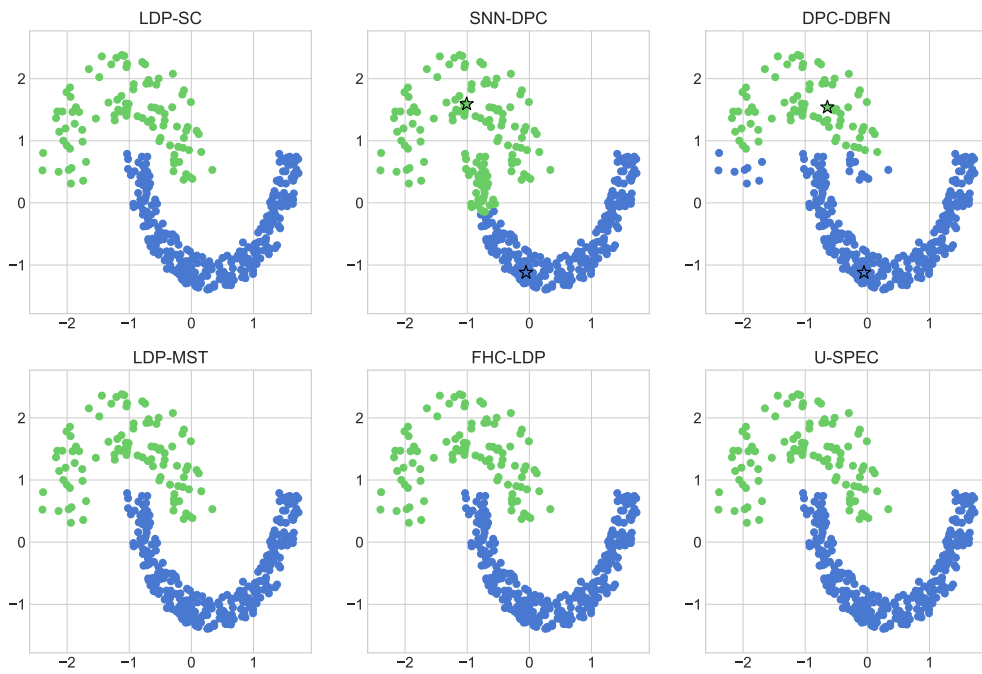


Figure 12: Jain

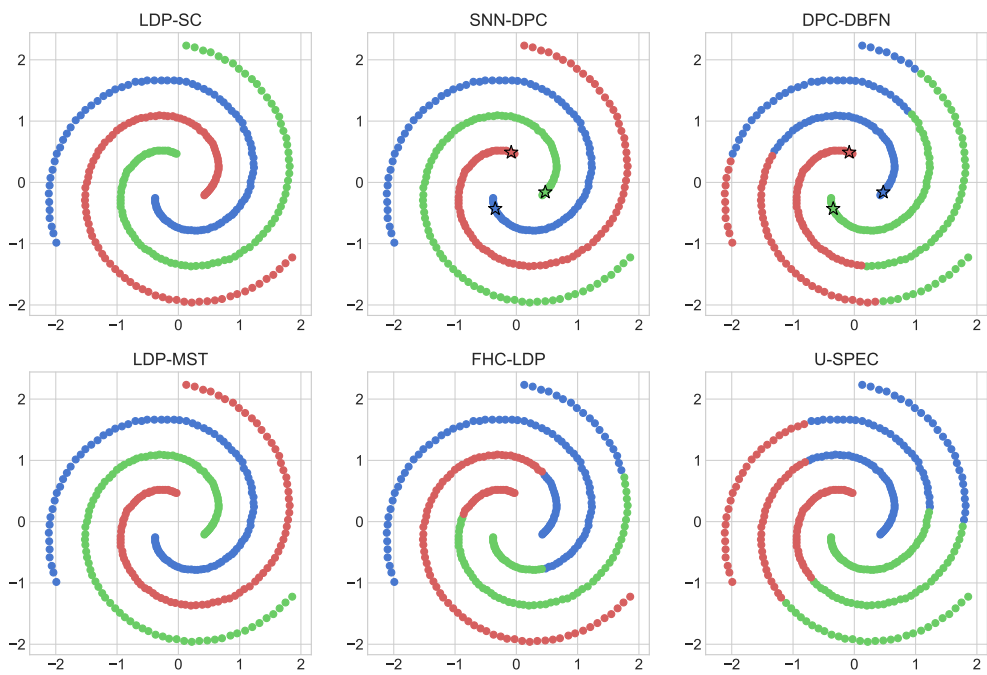


Figure 13: Spiral

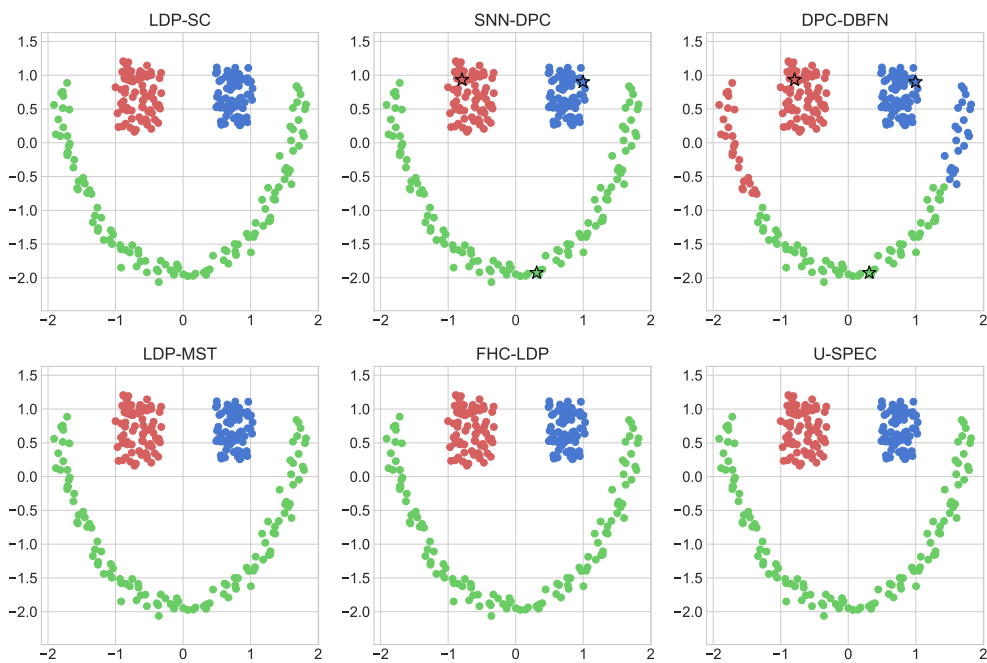


Figure 14: happy

Table 6: Results on real-world data sets (%).

Data set		LDP-SC	SNN-DPC	DPC	LDP-MST	FHC-LDP	DPC-DBFN	U-SPEC	LEC-K	SC
pengleukEW	ARI	50.41	36.21	13.48	29.22	38.17	18.30	13.48	29.61	46.46
	NMI	42.35	26.38	17.47	19.63	32.50	12.16	17.47	20.94	36.52
	ACC	86.11	80.56	72.22	77.78	81.94	72.22	72.22	77.78	84.72
	Par.	4	5	1	26	9	27	3	8	4
parkinsonsEW	ARI	46.23	45.21	40.28	41.71	39.10	9.77	-1.80	11.05	39.10
	NMI	35.30	31.87	30.38	33.60	35.14	5.15	15.65	11.70	35.14
	ACC	86.67	86.15	85.13	85.64	85.13	76.92	56.41	78.00	85.13
	Par.	7	25	0.4	13	5	8	4	2	5
seeds	ARI	78.48	77.76	79.97	57.02	72.89	76.64	82.31	82.25	82.18
	NMI	72.42	74.23	75.64	63.07	69.46	73.43	77.3	77.78	77.86
	ACC	92.38	91.90	92.86	81.90	90.00	91.43	93.81	93.81	93.81
	Par.	23	6	0.6	5	12	2	2	8	7
spect	ARI	30.62	35.98	-1.36	28.83	32.18	32.18	-1.99	-8.29	20.51
	NMI	20.59	19.33	0.06	16.36	19.93	19.93	13.03	5.03	12.72
	ACC	84.64	83.52	64.42	79.03	84.64	84.64	55.06	71.16	74.16
	Par.	5	15	0.20	3	23	24	4	6	4
mfeat-fac	ARI	88.38	71.86	64.15	55.97	81.69	48.85	83.48	85.01	85.90
	NMI	89.13	80.05	78.14	69.51	85.58	69.16	85.65	86.46	87.57
	ACC	94.60	78.40	69.95	63.65	90.95	59.25	92	92.88	93.25
	Par.	30	35	0.30	38	29	17	3	7	13
landsatEW	ARI	62.80	52.56	54.79	50.01	0.25	16.92	59.50	52.58	45.13
	NMI	65.02	57.80	58.24	59.83	1.64	30.40	63.61	62.83	59.04
	ACC	76.02	69.43	75.04	64.21	24.49	42.46	72.79	68.38	64.55
	Par.	6	8	0.20	6	3	2	6	3	7
Pendigits	ARI	80.04	64.17	65.10	70.11	68.36	54.50	71.93	64.31	77.67
	NMI	86.31	78.54	76.10	81.66	82.20	69.79	80.22	79.72	84.25
	ACC	89.30	74.15	76.72	78.02	79.08	69.92	84.95	74.65	88.05
	Par.	12	42	0.9	7	48	17	27	5	25
covertype	ARI	16.88	N/A	N/A	N/A	N/A	N/A	14.42	8.06	N/A
	NMI	18.24	N/A	N/A	N/A	N/A	N/A	12.21	12.12	N/A
	ACC	44.39	N/A	N/A	N/A	N/A	N/A	47.04	48.33	N/A
	Par.	21	-	-	-	-	-	3	6	-

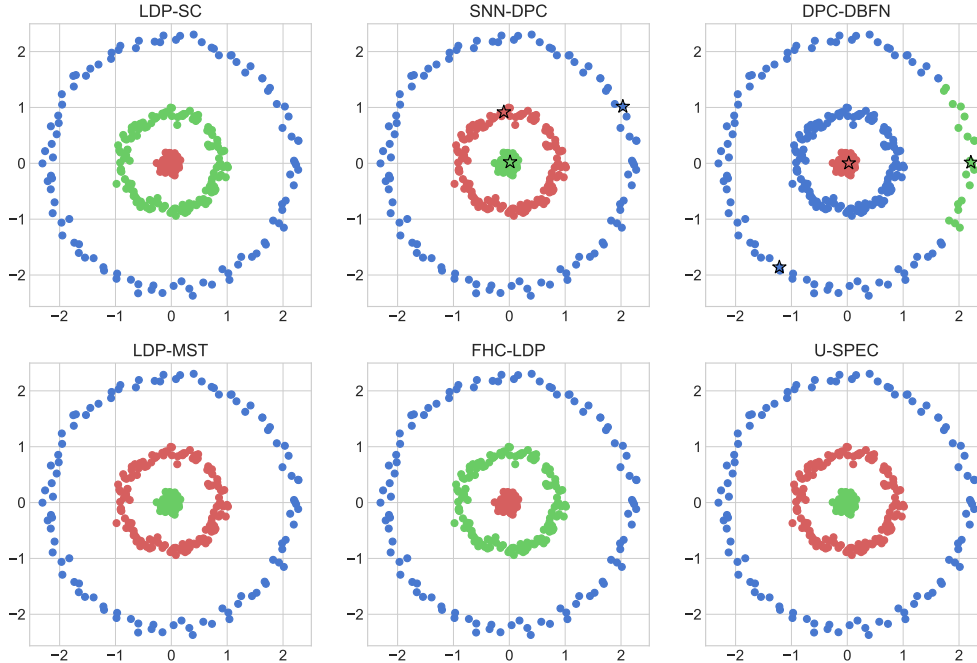


Figure 15: circle

517 For the Yale, COIL20, USPS, and MNIST data sets, LDP-SC has the best performance for
 518 all three metrics. For ORL, LDP-SC has the highest ARI and NMI scores, and has a only
 519 slightly lower ACC score than the best one. For Fashion MNIST data set, although LDP-SC
 520 is not the best, its performance is very close to the best. Again, for the data sets where
 521 LDP-SC performs best, the scores of it are significantly higher than the others, especially
 522 for the USPS data set. These results also show that LDP-SC is very promising in high
 523 dimensional and complex real-world data.

524 6. Discussion

525 6.1. Robustness

526 In the previous section, we discussed the optimal results of each algorithm with known
 527 data labels. However, unsupervised clustering can only be performed without labels in
 528 practice. Therefore, a good clustering algorithm should also be robust w.r.t. different values
 529 of parameters, so that its performance remains good enough without tuning parameters
 530 based on labels. Note that LDP-SC constructs a parent-child relationship by looking for
 531 the data point with higher density that is closest to a given data point. This relationship
 532 is expected to be insensitive to the only parameter k , i.e., the number of nearest neighbors:
 533 for a non-root point, when k changes within a reasonable range, its parent is likely to be
 534 the same one as before (because its parent is close to it and the relative ordering of their
 535 densities is likely to be unchanged); for a root point, k would need to become quite large to
 536 have a higher density point. To verify this empirically, Fig. 16 shows the fluctuation of ARI

Table 7: Results on image data sets (%).

Data sets		LDP-SC	SNN-DPC	DPC	LDP-MST	FHC-LDP	DPC-DBFN	U-SPEC	LEC-K	SC
ORL	ARI	52.76	31.36	38.59	28.69	37.33	28.87	31.66	45.96	49.66
	NMI	83.72	76.26	77.13	67.52	79.59	67.69	70.06	79.89	82.35
	ACC	65.00	52.50	50.50	49.2	58.75	47.25	45.50	63.05	66.50
	Par.	4	5	0.4	18	5	28	2	3	5
Yale	ARI	14.17	4.06	6.24	2.09	6.44	4.40	4.40	4.56	9.50
	NMI	47.04	36.00	37.83	12.46	37.65	31.50	20.81	26.34	41.74
	ACC	33.51	24.15	28.38	8.53	28.54	24.19	15.86	19.61	29.25
	Par.	3	6	3	42	9	2	2	2	3
COIL20	ARI	84.33	61.68	60.86	69.54	79.17	54.51	76.74	82.56	71.11
	NMI	96.30	82.07	84.19	89.42	92.76	75.59	89.69	95.38	86.24
	ACC	87.29	69.44	69.65	77.29	83.47	61.75	80.94	85.00	79.38
	Par.	3	41	2.6	13	9	26	2	3	13
USPS	ARI	91.91	61.82	35.65	80.29	68.91	51.16	80.72	75.20	68.02
	NMI	90.03	76.94	52.29	81.95	77.44	64.70	82.45	83.05	82.21
	ACC	95.81	66.41	46.86	86.47	73.13	63.49	84.76	78.92	68.36
	Par.	4	16	4.2	28	15	4	5	2	4
MNIST	ARI	75.75	N/A	N/A	62.66	N/A	N/A	64.95	63.16	62.95
	NMI	83.11	N/A	N/A	69.58	N/A	N/A	74.35	74.27	76.47
	ACC	79.75	N/A	N/A	68.57	N/A	N/A	74.57	73.72	68.31
	Par.	8	-	-	21	-	-	2	2	6
Fashion MNIST	ARI	46.92	N/A	N/A	40.36	N/A	N/A	47.01	40.16	46.8
	NMI	62.96	N/A	N/A	58.06	N/A	N/A	62.54	56.22	65.89
	ACC	58.69	N/A	N/A	47.08	N/A	N/A	58.85	58.13	58.33
	Par.	9	-	-	11	-	-	3	35	4

537 values of LDP-SC, LSC-K, SNN-DPC, U-SPEC, FHC-LDP and DPC-DBFN w.r.t. k on the
538 data sets spect, seeds, mfeat-fac and Pendigits. For DPC based algorithms, $d_c = k/10$. Note
539 that U-SPEC performs clustering on the specified core points, which might be smaller than
540 the given k and thus makes the algorithm have no result in that case (spect and seeds data
541 sets). From the results, it can be seen that LDP-SC algorithm is very stable w.r.t. changes
542 of k and can be considered as one of the most stable ones among the six algorithms.

543 6.2. Ablation Study

544 In order to verify the significance of each component in the proposed algorithm, an
545 ablation study was conducted here on seven real-world data sets. The results are optimal
546 by searching for the best parameters, and are shown in Table 8, where columns a, b, c, d
547 represents the following different settings, respectively:

- 548 a. Graph cut is changed to be on the roots only, using the original spectral clustering
549 algorithm. This is to verify the significance of the tree structures and the corresponding
550 similarity measure (Eq. 12) considered in this article.
- 551 b. The term of separation degree $(1 + \sigma(T(r_i), T(r_j)))$ in similarity measure (Eq. 12) is
552 removed. This is to verify that separation degree is useful for complex structures of
553 real-world data.
- 554 c. The term $|\text{SNN}(T(r_i), T(r_j))|$ in the similarity measure is replaced with the constant 1.
555 This is to verify that the neighboring information of trees is useful for clustering.
- 556 d. The total number of nodes in trees is removed from the loss function Ntcut (Eq. 13).
557 This is to verify the significance of tree sizes for clustering.

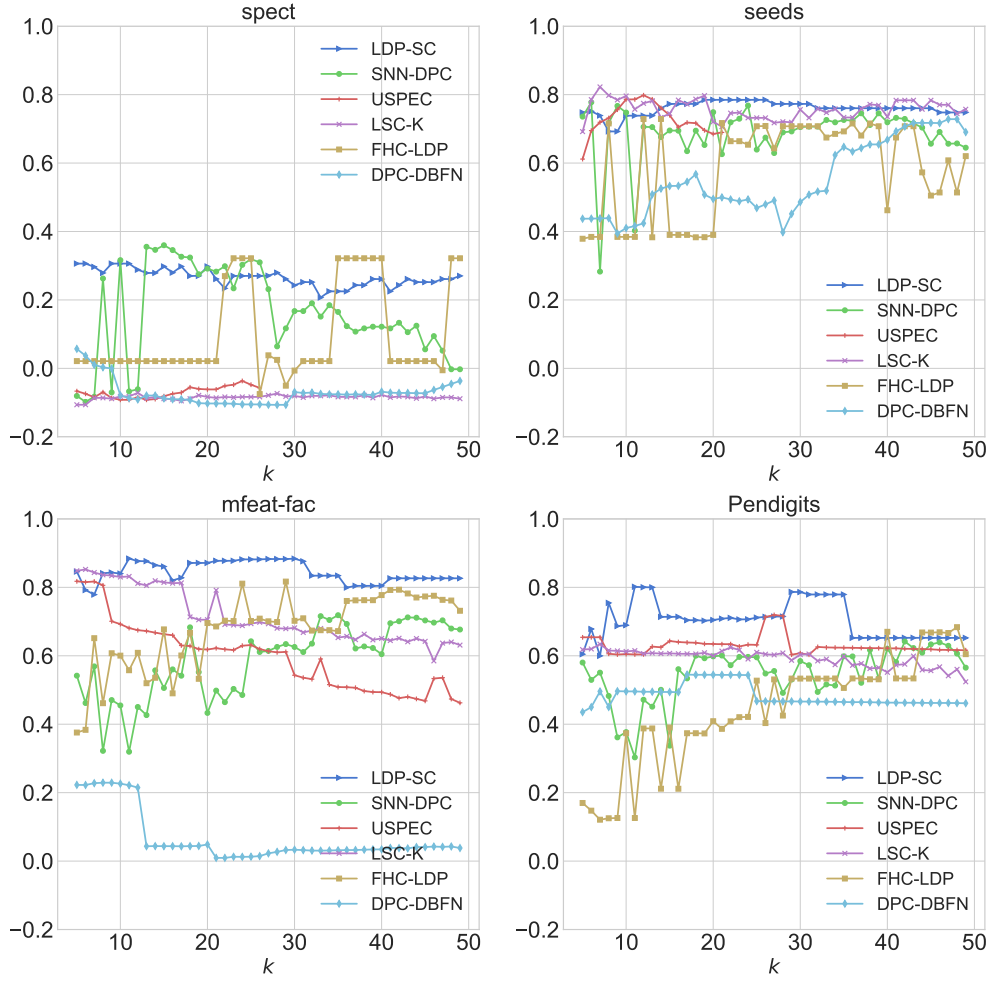


Figure 16: Performance changes with respect to the number of nearest neighbors k .

558 According to Table 8, as there will be a decrease in performance when a component is
 559 removed or changed in most data sets, it can be seen that each component of the proposed
 560 algorithm is useful.

561 7. Conclusion

562 This article proposes a new clustering algorithm based on the idea of local density peaks
 563 and graph cut. In particular, this algorithm takes the advantages of the DPC algorithm in
 564 local neighborhoods to capture latent structures of data to construct family trees, devises
 565 a novel similarity measure between family trees, and adapts the Ncut objective function to
 566 perform graph cut on the connection graph of these trees. In this way, a data set can be
 567 clustered by allocating data points in a family tree to the cluster of the tree. We performed
 568 a variety of empirical evaluations of the proposed algorithm on 22 challenging synthetic data
 569 sets and real-world data sets. The results demonstrated the superiority of our algorithm over

Table 8: Results of ablation study (in ARI, %).

Data set	LDP-SC	a	b	c	d
pengleukEW	50.41	0	50.41	13.76	50.41
spect	30.62	16.19	23.07	30.62	30.62
landsatEW	62.80	36.43	60.04	51.48	58.34
mfeat-fac	88.38	73.87	88.38	84.78	80.88
parkinsonsEW	46.23	39.10	39.10	39.10	46.23
Pendigits	80.04	56.23	79.16	72.68	79.50
seeds	78.48	78.48	78.48	78.48	78.48

570 several prominent clustering algorithms, in terms of three commonly used metrics (ARI,
571 NMI, and ACC). The analysis of its robustness w.r.t. hyper-parameters shows that the
572 proposed algorithm has stable performance with a large range of the hyper-parameter k ,
573 the number of nearest neighbors. This new algorithm is also very efficient compared to
574 several algorithms including SNN-DPC, LDP-MST, DPC-DBFN, FHC-LDP, DPC, and SC
575 on large data sets, and is competitive compared to accelerated spectral cluster algorithms.
576 Ablation study also proves that each of its components has a contribution to the overall
577 performance of the algorithm. Nevertheless, the current algorithm still has some defects,
578 e.g., the efficiency of the algorithm can be further improved, the number of target clusters
579 can only be determined manually, and the ability to detect and remove noise data is still
580 missing. All of these will be interesting to explore in the future.

581 Acknowledgment

582 This work was supported by the National Natural Science Foundation of China [grant
583 number 61806170]; the Humanities and Social Sciences Fund of Ministry of Education [grant
584 number 18XJC72040001]; and the National Key Research and Development Program of
585 China [grant number 2019YFB1706104].

References

- [1] A. K. Jain, Data clustering: 50 years beyond K-means, *Pattern Recognition Letters* 31 (2010) 651–666.
- [2] J. A. Hartigan, *Clustering algorithms*, John Wiley & Sons, Inc., 1975.
- [3] A. Y. Ng, M. I. Jordan, Y. Weiss, On spectral clustering: Analysis and an algorithm, in: *Advances in Neural Information Processing Systems*, 2002, pp. 849–856.
- [4] B. J. Frey, D. Dueck, Clustering by passing messages between data points, *Science* 315 (2007) 972–976.
- [5] U. von Luxburg, R. C. Williamson, I. Guyon, Clustering: Science or art?, in: *International Conference on Machine Learning: Workshop on Unsupervised and Transfer Learning*, 2012, pp. 65–79.
- [6] J. MacQueen, others, Some methods for classification and analysis of multivariate observations, in: *Berkeley Symposium on Mathematical Statistics and Probability*, volume 1, 1967, pp. 281–297.
- [7] T. Zhang, R. Ramakrishnan, M. Livny, BIRCH: An efficient data clustering method for very large databases, volume 25, 1996, pp. 103–114.
- [8] M. Ester, H.-P. Kriegel, J. Sander, X. Xu, A density-based algorithm for discovering clusters in large spatial databases with noise, in: *International Conference on Knowledge Discovery and Data Mining*, volume 96, 1996, pp. 226–231.

- [9] D. Alahakoon, S. K. Halgamuge, B. Srinivasan, Dynamic self-organizing maps with controlled growth for knowledge discovery, *IEEE Transactions on Neural Networks* 11 (2000) 601–614.
- [10] W. Wang, J. Yang, R. Muntz, STING: A statistical information grid approach to spatial data mining, in: *International Conference on Very Large Data Bases*, volume 97, 1997, pp. 186–195.
- [11] U. von Luxburg, A tutorial on spectral clustering, *Statistics and Computing* 17 (2007) 395–416.
- [12] C. Xu, D. Tao, C. Xu, A survey on multi-view learning, *arXiv preprint arXiv:1304.5634* (2013).
- [13] A. Strehl, J. Ghosh, Cluster ensembles—a knowledge reuse framework for combining multiple partitions, *Journal of Machine Learning Research* 3 (2002) 583–617.
- [14] J. R. Hershey, Z. Chen, J. Le Roux, S. Watanabe, Deep clustering: Discriminative embeddings for segmentation and separation, in: *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2016, pp. 31–35.
- [15] A. Rodriguez, A. Laio, Clustering by fast search and find of density peaks, *Science* 344 (2014) 1492–1496.
- [16] M. Du, S. Ding, H. Jia, Study on density peaks clustering based on k-nearest neighbors and principal component analysis, *Knowledge-Based Systems* 99 (2016) 135–145.
- [17] Z. Bian, F.-L. Chung, S. Wang, Fuzzy density peaks clustering, *IEEE Transactions on Fuzzy Systems* 29 (2021) 1725–1738.
- [18] J. Xie, H. Gao, W. Xie, X. Liu, P. W. Grant, Robust clustering by detecting density peaks and assigning points based on fuzzy weighted k-nearest neighbors, *Information Sciences* 354 (2016) 19–40.
- [19] R. Liu, H. Wang, X. Yu, Shared-nearest-neighbor-based clustering by fast search and find of density peaks, *Information Sciences* 450 (2018) 200–226.
- [20] D. Cheng, S. Zhang, J. Huang, Dense members of local cores-based density peaks clustering algorithm, *Knowledge-Based Systems* 193 (2020) 105454.
- [21] M. Du, S. Ding, X. Xu, Y. Xue, Density peaks clustering using geodesic distances, *International Journal of Machine Learning and Cybernetics* 9 (2018) 1335–1349.
- [22] Y. Chen, X. Hu, W. Fan, L. Shen, Z. Zhang, X. Liu, J. Du, H. Li, Y. Chen, H. Li, Fast density peak clustering for large scale data based on kNN, *Knowledge-Based Systems* 187 (2020).
- [23] X. Xu, S. Ding, M. Du, Y. Xue, DPCG: An efficient density peaks clustering algorithm based on grid, *International Journal of Machine Learning and Cybernetics* 9 (2018) 743–754.
- [24] J. Guan, S. Li, X. He, J. Zhu, J. Chen, Fast hierarchical clustering of local density peaks via an association degree transfer method, *Neurocomputing* 455 (2021) 401–418.
- [25] L. Yaohui, Z. Ma, Y. Fang, Adaptive density peak clustering based on k-nearest neighbors with aggregating strategy, *Knowledge-Based Systems* 133 (2017) 208–220.
- [26] G. Wang, Y. Wei, P. W. Tse, Clustering by defining and merging candidates of cluster centers via independence and affinity, *Neurocomputing* 315 (2018) 486–495.
- [27] A. Lotfi, P. Moradi, H. Beigy, Density peaks clustering based on density backbone and fuzzy neighborhood, *Pattern Recognition* 107 (2020) 107449.
- [28] Z. Su, T. Denoeux, BPEC: belief-peaks evidential clustering, *IEEE Transactions on Fuzzy Systems* 27 (2019) 111–123.
- [29] F. Fang, L. Qiu, S. Yuan, Adaptive core fusion-based density peak clustering for complex data with arbitrary shapes and densities, *Pattern Recognition* 107 (2020) 107452.
- [30] D. Cheng, Q. Zhu, J. Huang, Q. Wu, L. Yang, Clustering with local density peaks-based minimum spanning tree, *IEEE Transactions on Knowledge and Data Engineering* 33 (2021) 374–387.
- [31] W.-Y. Chen, Y. Song, H. Bai, C.-J. Lin, E. Y. Chang, Parallel spectral clustering in distributed systems, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 33 (2010) 568–586.
- [32] D. Cai, X. Chen, Large scale spectral clustering via landmark-based sparse representation, *IEEE Transactions on Cybernetics* 45 (2014) 1669–1680.
- [33] D. Huang, C.-D. Wang, J.-S. Wu, J.-H. Lai, C.-K. Kwok, Ultra-scalable spectral clustering and ensemble clustering, *IEEE Transactions on Knowledge and Data Engineering* 32 (2020) 1212–1226.
- [34] Y. Li, W. Liu, Y. Wang, D. Tao, Co-spectral clustering based density peak, in: *International Conference on Communication Technology*, 2015, pp. 925–929.

- [35] J. Liu, C. Zhao, Density gain-rate peaks for spectral clustering, *IEEE Access* 9 (2021) 46000–46010.
- [36] J. Shi, J. Malik, Normalized cuts and image segmentation, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 22 (2000) 888–905.
- [37] C. F. Van Loan, G. Golub, *Matrix computations*, Johns Hopkins Press, 1996.
- [38] A. Villar-Corrales, V. I. Morgenshtern, Scattering transform based image clustering using projection onto orthogonal complement, in: *Workshop on Intelligent Cross-Data Analysis and Retrieval*, 2021, pp. 24–32.
- [39] D. Steinley, Properties of the Hubert-Arable Adjusted Rand Index, *Psychological Methods* 9 (2004) 386–396.
- [40] W. Xu, X. Liu, Y. Gong, Document clustering based on non-negative matrix factorization, in: *Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2003, pp. 267–273.
- [41] Y. Yang, D. Xu, F. Nie, S. Yan, Y. Zhuang, Image clustering using local discriminant models and global integration, *IEEE Transactions on Image Processing* 19 (2010) 2761–2773.