

# Depth based local center clustering

SIYI WANG

2022/10/21

# 1 Background

## 1.1 Data depth

Ordering is an important concept in analyzing data. In one dimensional case, the arrangement of real numbers provides a natural ordering. However, in multivariate case, there is no natural ordering of vectors. Data depth provides a non-parametric approach to give a central-outward ordering of data.

Assume a data set  $\mathbf{X}$  follows a distribution  $F$  in a  $d$  dimensional real space,  $\mathbb{R}^d$ , a depth function for  $\mathbf{X}$  is defined as

$$D(\bullet, F) : \mathbb{R}^d \rightarrow [0, 1] : z \mapsto D(z \mid F),$$

where  $z \in \mathbb{R}^d$ . The following restrictions should, in theory, be satisfied by data depth [1]:

1. Affine invariance,
2. Null at infinity,
3. Monotone on rays,
4. Upper semi-continuous.

In real applications, the inherent distribution  $F$  is unknown, and hence depth is calculated by the empirical distribution of  $\mathbf{X}$ . For a more intuitive notation, depth function is written as  $D(z \mid \mathbf{X})$  instead of  $D(z \mid F)$  in this report.

## 1.2 Literature Review

Numerous applications of data depth have been developed since Tukey's 1975 introduction of the first depth function, such as multivariate density estimation [2], tests for multivariate scale difference [3], robust and affine equivariant estimations of location [4], classification [5, 6, 7] and so on.

Much research has been done using data depth in classification, while less attention is given to the problem of clustering. Jornsten [5] proposed an algorithm called DDclust, improves the

output of other clustering algorithms. It takes as input the result from a clustering algorithm like partition around medoids (PAM) clustering and calculates the depths of observations regarding each cluster based on the initial partition to see if relocating some observations can achieve a better outcome. In 2016, Jeong et al. [8] proposed a density-based clustering algorithm called data depth based clustering analysis (DBCA), which can be regarded as an alternative version of the well-known DBSCAN algorithm [9]. DBCA modifies Mahalanobis depth to a local version as follows:

$$D_M(\mathbf{z} \mid \mathbf{x}_i) = \left[ 1 + (\mathbf{z} - \mathbf{x}_i)' \hat{\Sigma}^{-1} (\mathbf{z} - \mathbf{x}_i) \right]^{-1}, \quad (1)$$

where  $\mathbf{z}$  denotes a point in  $\mathbb{R}^d$ ,  $\mathbf{x}_i$  is an observation in  $\mathbf{X}$ , while  $\hat{\Sigma}$  is the estimated covariance matrix of data set  $\mathbf{X}$ . Compared with DBSCAN, DBCA utilizes (1) to measure similarities among observations instead of Euclidean distance. The main advantages of DBCA are affine invariance and robustness to noises.

So far, the notion of depth serves an auxiliary role to the more traditional clustering methods, and is not employed directly in examining the multimodality of data. In this report, we propose an innovative clustering algorithm that relies on the center-outward ranks given by data depth, named Depth-based Local center Clustering (DLCC). DLCC determines first the neighbors for each point, and then computes depth values regarding the subset constructed by the point and its neighbors. Data depth can be used to locate centers of those subsets, which are designated as local centers. The method then detect inherent clusters by analyzing the local centers.

## 2 Algorithm

Data depth describes data from a global perspective. To include the concept of data depth into clustering, it is necessary to examine depth based on some subsets of the entire data set, where we refer to the depth calculated with regard to a subset as local depth (with respect to the subset). In DLCC,  $n$  subsets built by each observation and its neighbors are examined, where  $n$  is the total number of observations in a data set  $\mathbf{X}$ . Fix a size  $s$  of subset, neighbors of

each observation are defined as the top  $s - 1$  similar points. Possible approaches for building similarity matrix and defining neighbors are introduced in Section 2.1. Subset form by the observation  $\mathbf{x}_i$  and its neighbors is noted as  $\mathbf{X}_i$ . Data depth values are used to define the rank of any observation in  $\mathbf{X}_i$ .

**Definition 2.1.** Note  $\mathbf{x}_j$  is an observation in data set  $\mathbf{X}_i$ . The rank of  $\mathbf{x}_j$  in  $\mathbf{X}_i$  is defined as Equation (2).

$$r_{\mathbf{X}_i}(\mathbf{x}_j) = \# \{ \mathbf{x}_q : D(\mathbf{x}_j | \mathbf{X}_i) < D(\mathbf{x}_q | \mathbf{X}_i), \mathbf{x}_q \in \mathbf{X}_i \} + 1. \quad (2)$$

If  $r_{\mathbf{X}_i}(\mathbf{x}_j) = 1$ ,  $\mathbf{x}_j$  is the center of  $\mathbf{X}_i$ . Local centers are defined to be the centers of the subsets  $\mathbf{X}_i$ 's. If the center of  $\mathbf{X}_i$  is  $\mathbf{x}_i$ , then  $\mathbf{x}_i$  is said as the strict local center.

Local centers can be interpreted as representative points of  $\mathbf{X}$ , and the quality of them is crucial in following steps. Hence, before grouping local centers, it is necessary to filter them. For convenience, the local centers after filtering are called filtered centers. The algorithm of filtering local centers is explained in Section 2.2. After filtering, the local centers are grouped. One possible approach is based on the proportion of overlapped observations between neighbors of two local centers, and this proportion is called as the similarity between local centers. Here, a new parameter  $\delta$ , the similarity threshold, is introduced. Based on  $\delta$ , two strategies of DLCC named “min” and “max”, dealing with different scenarios of clustering problems, are summarized as follows:

**Definition 2.2.** Let  $\{c_1, \dots, c_P\}$  be a finite set of filtered centers. The grouping result under two strategies should satisfy following requirements.

- **Min strategy:**  $\forall c_i$  and  $c_j$  in the same group, there is  $\text{sim}(c_i, c_j) > \delta$ .
- **Max strategy:**  $\forall c_i$  in group  $\mathbf{g}$ ,  $\exists c_j$  in the same group s.t.  $\text{sim}(c_i, c_j) > \delta$ , and for  $\forall c_v \notin \mathbf{g}$ ,  $\text{sim}(c_i, c_v) \leq \delta$ .

Some insights into the strategy chosen for grouping local centers are significant in applications. Under min strategy, DLCC is a centroid-based clustering method, which attempts to find centroids of latent clusters from filtered centers. The min strategy assumes the sizes of all latent clusters are similar, and the parameter  $s$  under min strategy should be close to

the real size of latent clusters. On the other hand, the min strategy may potentially encounter the multiple assignment problem. For instance, there may be a collection of filtered centers  $\{c_p\}$  so that the only similarities larger than  $\delta$  are between  $c_1$  and each of the other centers, while the similarities between the other pairs in the collection are all less than  $\delta$ . Then  $c_1$  will be assigned into multiple groups by this strategy. In this case, we simply drop  $c_1$  from the grouping consideration due to the ambiguity.

The max strategy shows a resemblance with DBCA. It assumes each cluster can be built by connecting multiple finite convex shapes created by filtered centers. Filtered centers in the same group are connected with a chain of similarities. For example, if two filtered centers  $c_w$  and  $c_y$  are in the same group, then there must exist a chain of filtered centers,  $c_{p_1}, c_{p_2}, \dots, c_{p_h}$ , where  $c_{p_1} = c_w$ ,  $c_{p_h} = c_y$ , and the similarity between any two adjacent centers in that chain is larger than  $\delta$ . Generally, the parameters  $s$  and  $\delta$  are comparatively small under the max strategy. Unlike DBCA, which builds clusters with the chain of small ellipses with fixed shape, and the center of each ellipse is each observation, DLCC with the max strategy only considers chains of filtered centers. Moreover, the number of neighbors of filtered centers is fixed rather than defined in a fixed region. DLCC can avoid a major shortcoming of DBCA, which has difficulty in handling clusters with different densities, due to the fixed shape of the ellipses. Similar to density-based clustering methods, max strategy is preferable when the latent clusters are disjoint, and it can handle non-convex and varying size clustering problems.

The number of clusters, denoted  $K$ , can then be determined, following the grouping of the filtered centers. Let  $G$  be the number of groups obtained from the filtered centers using one of the grouping strategies, then  $K \leq G$ . A unique neighbor of a group is defined as an observation that can only be a neighbor for centers in this particular group. If a certain group has no unique neighbor, then it will be dropped and leads to  $K < G$ . We obtain the temporary clusters from the unique neighbors of filtered centers in each of the groups. Then, we can define the scores of observations relative to each cluster, and establish cutoff points to determine whether to assign the observations to the cluster they obtain highest score with. In Section 2.3 the clustering procedure is presented in detail.

At this stage, the number of clusters has been determined, and we expect that the majority of observations have been assigned to one of the clusters. There may still be points remaining

outside of the clusters. In this situation, the clustering problem now can be regarded as a classification problem. Intuitively, we can utilize classification methods to assign the remaining points to one of the clusters. By default, DDLC adopts the most elementary classifier based on data depth, named maximum depth classifier, i.e., assign the observation to the cluster which it obtains the largest depth value. It should be noted that the choice of classification methods is very flexible here. In fact, any classification method may be considered for this step.

Lastly, self-improvements procedures of DLCC are discussed. Although the clustering problem is converted to a classification problem in the last stage, it has some special features, because the observations that are already labelled may still be updated. Here we introduce a logical parameter “maxdepth”, which has the potential to revise clustering results.

**Definition 2.3. maxdepth:** *A logical parameter controls whether the algorithm loops until all observations in the allocated cluster have the highest depth values.*

Let the  $K$  clusters constructed so far be denoted by  $\mathcal{X}_k$ , where  $k = 1, 2, \dots, K$ . For any observation  $x_j$ , the depth value of it with respect to each cluster is then given by  $D(x_j | \mathcal{X}_k)$ . If the depth value of  $x_j$  in current cluster does not equal to  $\max_k D(x_j | \mathcal{X}_k)$ , and if maxdepth is TRUE, then  $x_j$  will be relocated. This procedure loops until all observations are assigned to the cluster where they obtain the largest depth value. For min strategy, there is one more step for self-improvements. Generally, for a centroid-based clustering algorithm, each cluster should have just one centroid, while the number of filtered centers can be larger than the number of clusters. Therefore, it is natural to find the deepest filtered center for each cluster, and use them to re-produce clusters until the clustering result remains unchanged.

The entire DLCC algorithm is summarized below:

Step 1: Build a  $n \times n$  similarity matrix, and locate  $s - 1$  neighbors for each observation.

Step 2: Generate  $n$  subsets based on step 1, and find all local centers.

Step 3: Obtain filtered centers.

Step 4: Compute the similarities among filtered centers, and divide them into  $G$  groups based on min/max strategy.

Step 5: Determine the number of clusters  $K$  ( $K \leq G$ ) based on groupings of filtered centers and scores of observations with regard to clusters. Then build temporary clusters which contain majority observations.

Step 6: Classify the remaining observations using classification methods, such as maximum depth classifier, mixture-model based classification and so on.

Step 7: If “maxdepth” is TRUE, loop until all observations in the allocated cluster have the highest depth values.

Step 8: Under min strategy, re-define centroids  $\{c_1, \dots, c_K\}$  which satisfy  $D(c_k | \mathcal{X}_k) = \sup D(C | \mathcal{X}_k)$ , where  $C$  is the set of filtered centers. Loop from step 5 until the result does not change.

## 2.1 Similarity matrix building

To build the similarity matrix for searching neighbors of each observation, we use the Mahalanobis depth matrix in DBCA. Sample Mahalanobis depth function is defined as Equation (3) in  $\mathbb{R}^d$  space with regard to data set  $\mathbf{X}$ .

$$D_M(\mathbf{z} | \mathbf{X}) = \left[ 1 + (\mathbf{z} - \bar{\mathbf{X}})' \hat{\Sigma}^{-1} (\mathbf{z} - \bar{\mathbf{X}}) \right]^{-1}, D_M(\mathbf{z} | \mathbf{X}) \in [0, 1], \quad (3)$$

where  $\mathbf{z} \in \mathbb{R}^d$  and,  $\bar{\mathbf{X}}$  and  $\hat{\Sigma}$  represent the sample mean vector and estimated covariance matrix respectively. In order to obtain a robust result about neighbors of each observation, the covariance matrix here is estimated with the Minimum Covariance Determinant (MCD) method [10]. It can be seen that the modified equation in DBCA [8]

$$D_M(\mathbf{z} | \mathbf{x}_i) = \left[ 1 + (\mathbf{z} - \mathbf{x}_i)' \hat{\Sigma}^{-1} (\mathbf{z} - \mathbf{x}_i) \right]^{-1}$$

uses observation  $\mathbf{x}_i$ ,  $i = 1, 2, \dots, n$  rather than the mean vector in calculating depth value. That is, the center for calculating Mahalanobis depth is translated to each observation and the depth values of other observations are recorded. Hence, the modified Mahalanobis depth value shows the similarity to the observation in the center. Obviously, the right hand side of (1) can

be written as  $[1 + d_{Mah}^2(\mathbf{z}, \mathbf{x}_i)]^{-1}$ , where  $d_{Mah}(\mathbf{z}, \mathbf{x}_i)$  is the Mahalanobis distance between  $\mathbf{z}$  and  $\mathbf{x}_i$ , and the similarity matrix is equivalent to the distance matrix of Mahalanobis distance (larger similarity value means smaller Mahalanobis distance).

### 2.1.1 Min strategy

Under min strategy, a proper similarity matrix is critical; However, the estimated global covariance matrix may not be suitable enough for observations in different latent clusters. From the strong relationship between Mahalanobis distance and Principal Component Analysis (PCA), the squared Mahalanobis distance equals the sum of squares of standardised principal component scores [11]. If PCA can use the top few principle components (PCs) to explain the bulk of the variances, then the distance determined in the space constructed by the top PCs will resemble the Mahalanobis distance, and the PCA result can indicate strong correlations exist between variables. In contrast, there is no solid evidence that the Mahalanobis distance is superior to the standard Euclidean distance when there are no apparent correlations between variables. A simple test based on PCA is presented below, to determine whether the estimated global covariance matrix should be utilized to generate the similarity matrix.

Let  $\mathbf{v} = \{v_1, v_2, \dots, v_U, \dots, v_d\}$ , which contains the proportion of variance explained by each PC from high to low, where  $U$  indicates the number of PCs that explain over 5% variance. If  $U = d$ , then  $U$  redefined to be  $d - 1$ . If either Equation (4) or Equation (5) is satisfied, the data set passes the test.

$$v_1 > 0.6, \quad (4)$$

$$\sum_{u=1}^U v_u > 0.95. \quad (5)$$

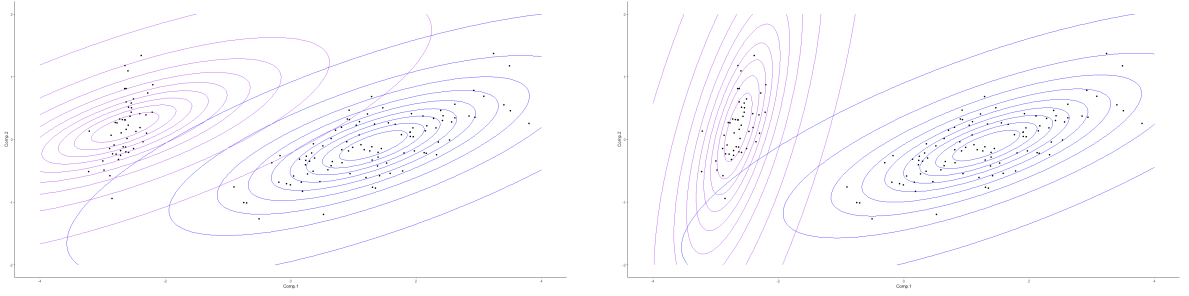
Even if a data set passes the test, a fixed covariance matrix may not be sufficient. Figure 1a shows an example from the iris data. With the fixed covariance matrix, it appears that the purple contour does not correspond to the orientation of the left side observations. Therefore, we investigate rotating the global covariance matrix based on the data distribution and utilize the Gaussian Mixture Model (GMM) concept. The estimated global covariance matrix can be



decomposed as Equation (6).

$$\hat{\Sigma} = \lambda \Gamma \Delta \Gamma', \quad (6)$$

where  $\lambda = |\hat{\Sigma}|^{1/d}$ ,  $\Gamma$  is an eigenvector matrix and  $\Delta$  is a diagonal matrix containing the eigenvalues, and  $|\Delta| = 1$  [12]. They represent the size, the orientation and the shape of the covariance matrix. To rotate the estimated covariance matrix, we keep  $\lambda$  and  $\Delta$ , and only adjust  $\Gamma$ . Then it is similar with a EEV model in Gaussian parsimonious clustering models (GPCM) family. Because  $\lambda$  and  $\Delta$  are given, the number of free covariance parameters is  $d$  less than in an EEV model, and we call it an Alternative EEV (AEEV) model. AEEV model gives two estimated covariance matrices for iris data. Figure 1b shows two similarity contours with the same shape and size, but different orientations for two selected data points.



(a) The similarity contour plots in two selected points based on the global covariance matrix estimated by MCD method.

(b) The similarity contour plots in two selected points based on the corresponding covariance matrices given by AEEV GMM.

Figure 1: Contour plots of the similarity measurement in the top two PCs of iris data set based on Mahalanobis depth.

If the data set does not pass the test, there are two scenarios. One is that there are significant correlations among variables within each latent cluster, but these correlations cannot be shown globally; the other is that there is no strong correlation among variables even within each latent cluster. Now, we try EEV/EEI model from GPCM family. The advantage of the EEV/EEI model is that it estimates clusters with identical eigenvalues. Therefore, we can conduct the same test using the eigenvalues provided by the EEV/EEI model. If it still fails, we assume the second scenario is correct and simply use Euclidean distance to construct the similarity matrix. To keep the value between 0 and 1 when using the Euclidean distance, we define

the  $(z, x_i)$ th entry of the similarity matrix as  $[1 + d_{Eu}(z, x_i)]^{-1}$ , where  $d_{Eu}$  is the Euclidean distance. Otherwise, we accept the first scenario and construct the similarity matrix using the covariance matrices provided by the EEV/EEI model.

Those steps can be summarized as trade-offs between GMM and DLCC. If GMM explains too much data, the size of covariance matrices given by GMM will be too small for measuring similarities among observations. Hence, if there is no significant difference in model performance (such as BIC) between GMM results with different numbers of components, we prefer the model with fewer components. Furthermore, if the similarity matrix is constructed from multiple covariance matrices provided by AEEV/EEV/EEI GMM, then similarities between observations are not mutual if they do not belong to the same cluster according to the GMM outcome, i.e.,  $D_M(x_i | x_j) \neq D_M(x_j | x_i)$ . Because neighbors of each observation are determined by order statistics, the non-mutual relationship in similarity is acceptable.

### 2.1.2 Max strategy

In terms of max strategy,  $s$  is small, and the number of filtered centers is large. Its clustering outcome is less sensitive than that of the min strategy for the similarity matrix. Furthermore, the mutual relationship for similarity between observations is useful in subsequent analyses of clustering performance under the max strategy. Thus, for the construction of the similarity matrix, we just employ the estimated global covariance matrix.

## 2.2 Local center selections

Locating and filtering local centers are the core part of DLCC. The chosen depth function is the Mahalanobis depth, where the covariance matrix is estimated using MCD. The point with the greatest depth value in a subset represents the center of that subset. For subsets  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ ,  $n$  local centers can be found, noted as  $\mathbf{a} = \{a_1, a_2, \dots, a_n\}$ . Owing to the robustness of data depth, if certain subsets are comparable, their deepest point may be same. Then we can define the frequency of local center as  $f$ , which stands for the times of an observation appears in  $\mathbf{a}$ .

Some local centers cannot contribute to DLCC. For example, an outlier which locates

between two latent clusters can be the deepest point of a subset containing observations from both clusters. To filter local centers, two benchmarks are proposed.

- A local center should locate at the central position of its own neighbors.
- A local center with higher  $f$  is more reliable.

The first step of filtering local centers is to pick up local centers satisfying  $r_{\mathbf{X}_i}(\mathbf{x}_i) \leq 2$  (Definition 2.1), which aims to meet the first benchmark. Due to the restricted number of strict local centers, the clustering result will be significantly impacted if the outcome of strict local centers is influenced by randomness from calculations. Therefore, the requirement of a central position is relaxed from  $r_{\mathbf{X}_i}(\mathbf{x}_i) = 1$  to  $r_{\mathbf{X}_i}(\mathbf{x}_i) \leq 2$ . For the remaining steps in filtering local centers, there are some differences between min and max strategies.

### 2.2.1 Min strategy

After the first step filtering, there are now  $t$  local centers remaining. Index these local centers in the decreasing order of their frequencies, i.e.

$$f_1 \geq f_2, \dots, \geq f_t,$$

Then the final result of the filtering step for min strategy will consists of the first  $P \leq t$  of these centers  $\{c_1, \dots, c_P\} \subseteq \mathbf{a}$ , respectively with frequencies  $f_1, \dots, f_P$ . To determine  $P$ , we consider the cumulative proportion of neighbors for the local centers down the list, as in (7).

$$\mathcal{P}_R = \frac{\#\{\bigcup_{p=1}^R \mathbf{X}_{c_p}\}}{n}. \quad (7)$$

Figure 2 plots the example for a simulated data set. Intuitively, if there is a significant jump in the cumulative proportion of a local center and all local centers after that jump contribute little to the cumulative proportion, then the particular local center is a potential cut-off point.

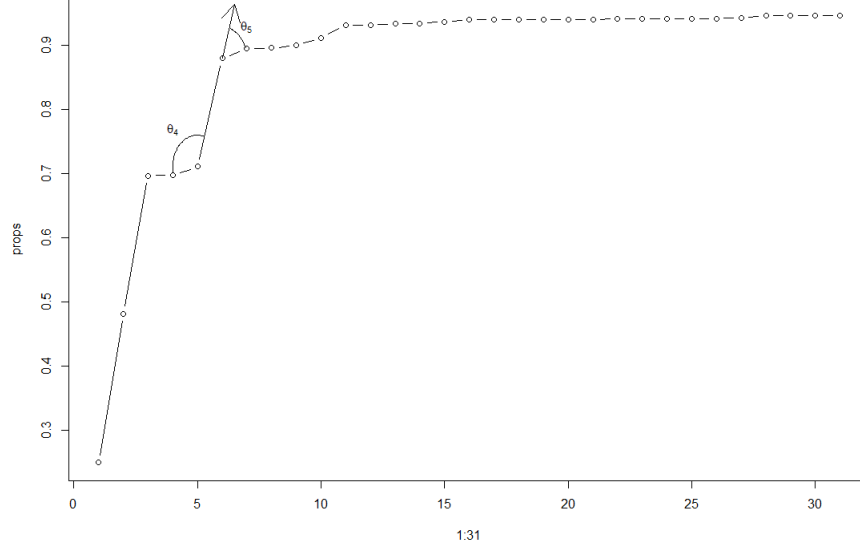


Figure 2: Cumulative proportion for 31 local centers from a simulated data set with 1000 observations.

The steps of searching the last significant jump when  $t > 3$  are summarized below:

- Let  $\Theta = \{\theta_1, \theta_2, \dots, \theta_{t-2}\}$  be the set consisting of intersection angles between lines connecting adjacent points in cumulative proportion plot. Calculate  $\tan \Theta$ , where  $\tan \theta_R = \frac{(\mathcal{P}_{R+1} - \mathcal{P}_R) - (\mathcal{P}_{R+2} - \mathcal{P}_{R+1})}{1 + (\mathcal{P}_{R+2} - \mathcal{P}_{R+1})(\mathcal{P}_{R+1} - \mathcal{P}_R)}$ .
- If there are  $\theta$ s satisfying  $\tan(\theta) < 0$  and  $|\tan(\theta)| > \text{Med}(|\tan \Theta|)$ , sort those  $|\tan(\theta)|$  in decreasing order. Otherwise, sort all  $|\tan(\theta)|$  in decreasing order. The permutation of  $|\tan(\theta)|$  defined here is written as  $(R_1, R_2, \dots, R_\tau)$ , where  $\tau = t - 2$  in the latter case.
- Set

$$\mathbb{L} = \{|\tan(\theta)|_{R_1} - |\tan(\theta)|_{R_2}, \dots, |\tan(\theta)|_{R_{\tau-1}} - |\tan(\theta)|_{R_\tau}\}.$$

For  $|\tan(\theta)|_{R_\tau} - |\tan(\theta)|_{R_{\tau+1}} = \max(\mathbb{L})$ ,  $P = R_\tau + 2$  if  $\tan \theta_{R_\tau} < 0$ ; Else,  $P = R_\tau + 1$ .

In the example of Figure 2,  $\tan(\theta_4) < 0$ , which indicates the slope of the line connecting  $c_4$  and  $c_5$  is smaller than the slope of the line connecting  $c_5$  and  $c_6$ . It means the jump right after  $c_5$  is more significant than the jump right before  $c_4$ . Moreover, as the number of local centers increases, the jump that may discover a latent cluster will become progressively smaller. If all

clusters have been found, then adding other local centers will contribute little to the cumulative proportion. Hence, we locate the cut-off point  $c_P$  by finding the largest gap in order statistics of  $|\tan(\theta)|$ . When  $t \leq 3$ , the number of latent clusters are assumed to be less or equal to  $t$ ,  $P$  is simply determined by satisfying  $\mathcal{P}_P = \min\{\mathcal{P}_R | \mathcal{P}_R > 0.7, R = 1, \dots, t\}$ .

### 2.2.2 Max strategy

With the same notions of Section 2.2.1, for max strategy, all remaining local centers satisfying  $f > 1$  included. Moreover, for any local center  $a_y$  with  $f = 1$ , if  $\bigcup_{p=1}^R \mathbf{X}_{c_p} \cap \mathbf{X}_{a_y} = \emptyset$  where  $R$  is the current number of filtered centers, it will be selected as well. The reason for this step is that if a cluster is constructed by subsets with non-convex shapes (such as ring and spiral), which contradicts the assumption of the max strategy, then it is possible that no local center in that cluster will remain after filtering because Mahalanobis depth is a convex depth function [1].

## 2.3 Clustering in DLCC

After unique neighbors for each group of local centers are assigned to temporary clusters, the score of observations to clusters are computed.

### 2.3.1 Min strategy

For the min strategy, the score accounts for both clustering and assessing if each observation's clustering result is acceptable. Let  $\mathbf{C}$  denote the set of filtered centers, then

$$\text{score}_{i|k} = \frac{\sup D_M(\mathbf{x}_i | \mathbf{C}_{g=k}) - \sup D_M(\mathbf{x}_i | \mathbf{C}_{g \neq k})}{\max\{\sup D_M(\mathbf{x}_i | \mathbf{C}_{g=k}), \sup D_M(\mathbf{x}_i | \mathbf{C}_{g \neq k})\}}. \quad (8)$$

Equation (8) is the score function for observation  $\mathbf{x}_i$  with respect to temporary cluster  $k$  under the min strategy. The score value is between  $-1$  and  $1$ . A positive score illustrates that compared with the nearest filtered center in other groups, the similarity between  $\mathbf{x}_i$  and the nearest filtered center in group  $k$  is larger.

For each cluster, two score pools are generated. The first is denoted  $\mathcal{S}$ , which consists of the unique neighbors in this cluster and their scores, while the other, denoted  $\hat{\mathcal{S}}_k$ , contains the

observations whose scores are positive with respect to this cluster, and their scores. Although unique neighbors are relatively trustworthy, some of them may not be reliable enough. A unique neighbor, for instance, can be the  $s - 1$ th similar point of a filtered center in group  $k$ , while also being the  $s$ th similar point of a filtered center in another group. Observations in  $\mathcal{S}_k$  whose scores are smaller than the median score in  $\hat{\mathcal{S}}_k$  are shifted to  $\hat{\mathcal{S}}_k$  to filter out less dependable points. Then, for scores in  $\hat{\mathcal{S}}_k$ , a cut-off point will be determined. All observations with scores above the cutoff point will be moved to  $\mathcal{S}_k$ . To determine the cutoff point, scores in  $\mathcal{S}_k$  are sorted in decreasing order as  $\{\gamma_1, \gamma_2, \dots, \gamma_E\}$ , where  $E$  represents the number of scores in  $\mathcal{S}_k$ . One candidate of the cut-off point is  $\gamma_e$  which satisfies Equation (9), where  $\lfloor * \rfloor$  is the floor function.

$$\gamma_e - \gamma_{e+1} = \max\{\gamma_{\lfloor E/2 \rfloor} - \gamma_{\lfloor E/2 \rfloor + 1}, \dots, \gamma_{E-1} - \gamma_E\}. \quad (9)$$

Another candidate for the cut-off point can be used when  $E < s/2$ . It is defined as the  $1 - \frac{s/2-E}{\#\{\hat{\mathcal{S}}_k\}}$  quantile of scores in  $\hat{\mathcal{S}}_k$ . This prevents a temporary cluster from having an insufficient amount of observations, and it can ensure that a temporary cluster contains at least  $s/2$  observations. It is necessary to note that, if  $E + \#\{\hat{\mathcal{S}}_k\} < s/2$ , DLCC under min strategy will fail because it violates the premise that all clusters are of comparable size. The cutoff point is specified as the lowest candidate value. After transferring observations with scores larger than the cutoff point in  $\hat{\mathcal{S}}_k$ , the temporary cluster  $k$  is updated as current observations in  $\mathcal{S}_k$ .

### 2.3.2 Max strategy

The score in this section is only for overlapping neighbors between filtered centers in different groups. Let  $\mathbf{C}_k$  be the set containing filtered centers in group  $k$ , then the score is defined as follows:

$$\text{score}_{i|k} = \frac{\#\{c \mid \mathbf{x}_i \in \mathbf{X}_c, c \in \mathbf{C}_k\}}{\#\{\mathbf{C}_k\}}. \quad (10)$$

All observations that are both neighbors of filtered centers in different groups are allocated to the cluster with the highest score. Observations that remain unlabeled are not neighbours of

any filtered center.

### 2.3.3 Classification and “maxdepth”

In terms of classification, if the default technique is used, the Mahalanobis depth function will continue to be used in both the maximum depth classifier and the algorithm for maximizing depth values in order to maintain consistency (if “maxdepth” is True). Unlike finding local centers, the covariance matrix now is based on the traditional moment estimation because in theory, temporary clusters are subsets of latent clusters, and all observations in the temporary clusters are selected with confidence after steps in Sections 2.3.1 and 2.3.2. It should not be necessary to further utilize MCD estimation for the covariance matrix. Besides the default algorithm, model-based classification method is also tested in applications.

## 3 Applications and Results

In this section, real data sets such as Iris, Seed and Wine from UCI machine learning repository [13] and synthetic data sets are used for investigating both the performance and the logic behind DLCC.

### 3.1 Min strategy

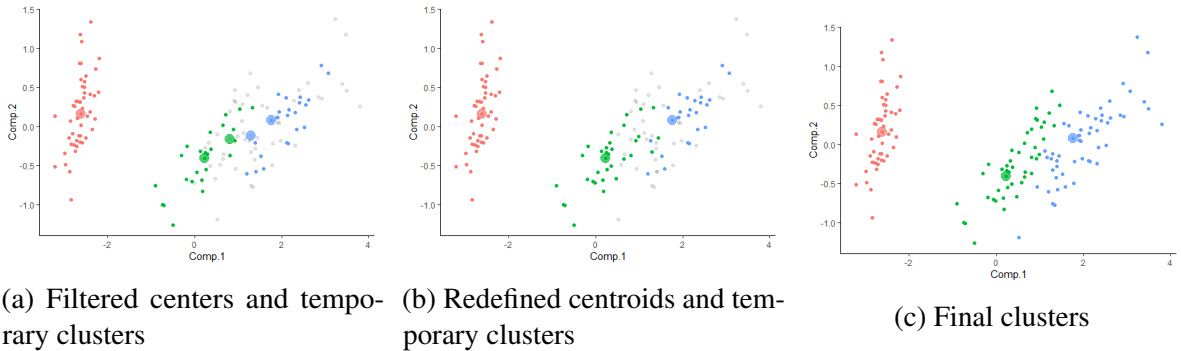


Figure 3: Visualizations for DLCC (min strategy) on the iris data set, where the  $x$  and  $y$  axes are the top two principle components and the color represents the clustering result. Grey points in temporary clusters are unlabelled observations.

Figure 3 shows there are 5 filtered centers on the Iris data set, and they are divided into 3 groups under min strategy. Notice that the number of filtered centers exceeds the number of clusters, the three deepest filtered centers from the current clustering result are selected as the re-defined centroids. In this step, two filtered centers locating at the border of two clusters are dropped in generating temporary clusters. Figure 3 also illustrates that, despite the presence of unlabeled observations in temporary clusters, the classification procedure yields a satisfactory clustering outcome.

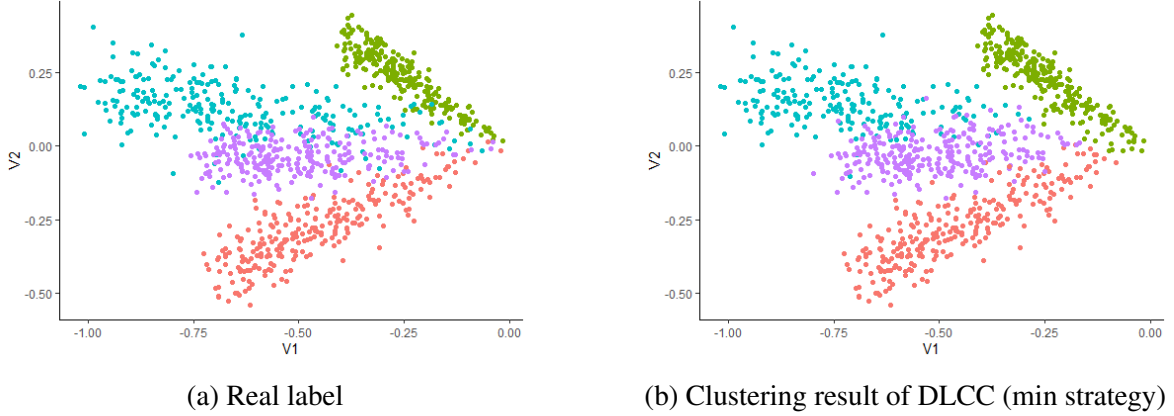


Figure 4: Graph of simulated data with ray shapes [14].  $x$  and  $y$  axes are the first two variables of ray data.

Figure 4 shows the first two dimensional visualization of ray shapes data simulated by Passino [14]. The orientation of each cluster is obviously different, which means that a global covariance matrix cannot provide a reasonable similarity matrix. In Figure 4, the clustering result given by DLCC closely resembles the real label, because the covariance matrices for constructing the similarity matrix here are provided by EEV GMM.

Table 1: Information regarding some chosen data sets, as well as settings of DLCC (min strategy) in them, where SM represents the approach for building similarity matrix, EU denotes the Euclidean distance.

| Data set | $n$  | $d$ | $K$ | SM   | $\delta$ | $s$ | Maxdepth | Classification method |
|----------|------|-----|-----|------|----------|-----|----------|-----------------------|
| Iris     | 150  | 4   | 3   | AEEV | 0.70     | 50  | TRUE     | Default               |
| Seed     | 210  | 7   | 3   | AEEV | 0.80     | 70  | TRUE     | Default               |
| Wine     | 178  | 13  | 3   | EU   | 0.70     | 50  | TRUE     | Default               |
| Ray      | 1000 | 5   | 4   | EEV  | 0.70     | 250 | FALSE    | Default               |

Table 1 provides some basic information of 4 data sets, as well as estimated number of



clusters  $K$  and parameters in DLCC. The clustering performance of DLCC in these data sets is compared to that of GMM and PAM, using the external metric Adjusted Rand Index (ARI) [15] to evaluate the clustering quality. For convenience, the number of clusters in GMM and PAM are both set to  $K$  in Table 1. As shown in Table 2, DLCC generally gives better performances with suitable similarity matrix and parameters.

Table 2: Clustering performances comparisons in chosen data sets.

| DATA SET | ARI   |       |       |
|----------|-------|-------|-------|
|          | DLCC  | GMM   | PAM   |
| Iris     | 0.904 | 0.904 | 0.642 |
| Seed     | 0.787 | 0.737 | 0.747 |
| Wine     | 0.982 | 0.930 | 0.741 |
| Ray      | 0.823 | 0.780 | 0.484 |

### 3.2 Max strategy

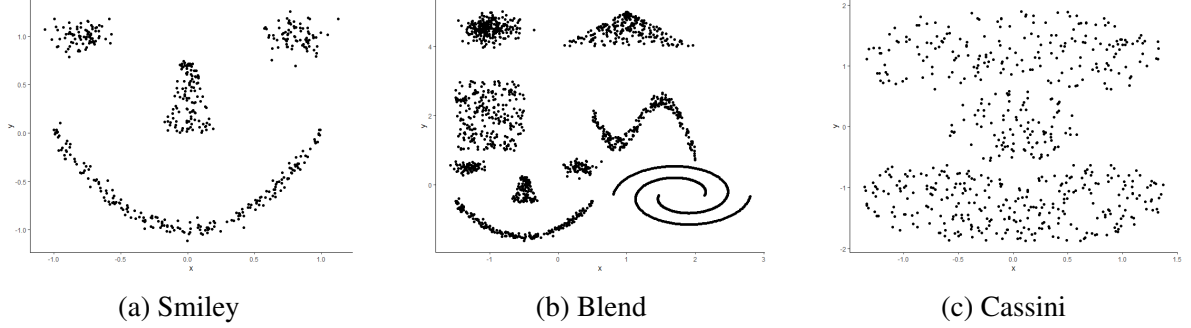


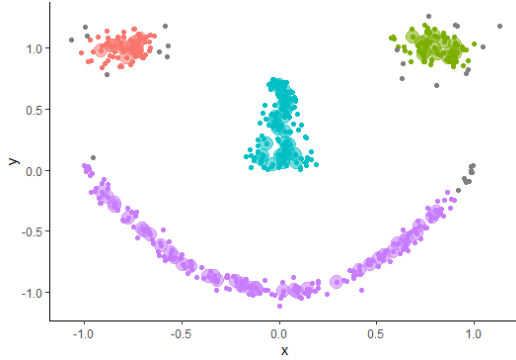
Figure 5: Graphs of three simulated data sets.

Figure 5 shows three synthetic data sets based on the mlbench package in R [16]. The third graph titled Cassini has three clusters, where two banana-shaped clusters gripping an oval-shaped cluster. From top to bottom, the number of points for three clusters are 200, 100 and 300 respectively, which results in varying densities among the three clusters.

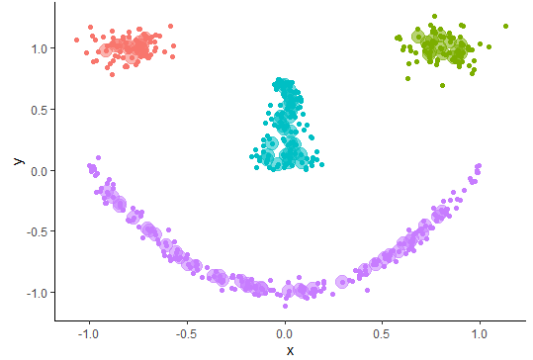
Figure 6 reveals that the max strategy tends to define many filtered centers. With the chains of filtered centers, eyes, nose and mouth are correctly clustered.

Table 3: Information regarding synthetic data sets, as well as the settings and performances of DLCC (max strategy) in them, where Mclass in Classification method means the mixture-model based classification.

| Data set | $n$  | $d$ | $K$ | $\delta$ | $s$ | Maxdepth | Classification method | ARI   |
|----------|------|-----|-----|----------|-----|----------|-----------------------|-------|
| Smiley   | 500  | 2   | 4   | 0.43     | 25  | FALSE    | Default               | 1     |
| Blend    | 2000 | 2   | 9   | 0.27     | 25  | FALSE    | Mclass                | 0.858 |
| Cassini  | 600  | 2   | 3   | 0.40     | 25  | TRUE     | Default               | 1     |

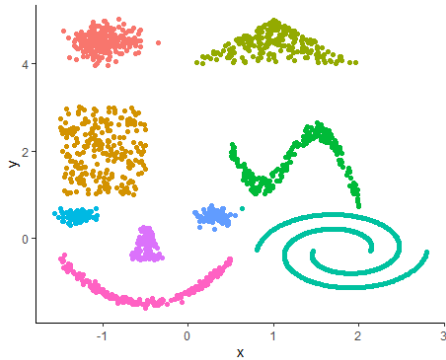


(a) Filtered centers and temporary clusters

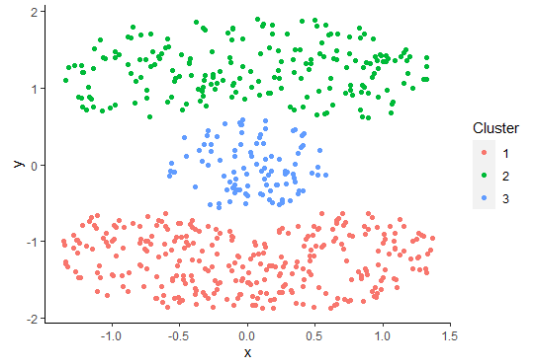


(b) Final clusters

Figure 6: Visualizations about for DLCC (max strategy) on the simulated smiley data set, where the color represents the clustering result. Grey points in temporary clusters are unlabelled observations



(a) Clustering result in Blend



(b) Clustering result in Cassini

Figure 7: Clustering results of DLCC (max strategy) in Blend and Cassini

Figures 6 and 7, and Table 3 show DLCC achieves perfect clustering in both smiley and Cassini, demonstrating its ability to tackle clustering problems with varying cluster densities. Although DLCC's performance in Blend is acceptable, it is unable to differentiate the spiral

into two clusters. In fact all filtered centers in the spiral with  $f = 1$ , which indicates that the geometry of any subset in the spiral used to locate local centers is non-convex. Denote this type of cluster as a local non-convex cluster, and DLCC seems to be not performing well with local non-convex clusters. Besides, a miss clustered point between the eye and the wave may be corrected by using simple distance-based classification method like K-nearest neighbors instead of depth or model based classification methods.

From the examples above using both min and max strategies, the following rule-of-thumb on “maxdepth” seems reasonable. If the shapes of all latent clusters are convex, and no clusters overlap, then “maxdepth” should be set to TRUE. For clusters with small overlapping regions, if the sizes of the covariance matrices of them are similar, a true “maxdepth” is also reasonable.

In conclusion, DLCC algorithm can obtain impressive clustering results under different scenarios of clustering problems, and its framework has the potential to be extended to other depth functions.

## 4 Parameter selection (max strategy)

In this section, we introduce a parameter selection algorithm under the max strategy, as well as a novel internal clustering criterion which can work for non-convex clustering problems. To avoid confusion, DLCC is defaulted to the max strategy in this section.

In DLCC, the choice of parameters  $s$  (size) and  $\delta$  (similarity threshold) affect the resulting clustering differently. For a particular dataset, the size  $s$  determines the set of filtered centers  $\mathbf{C} = \{c_1, \dots, c_P\}$ . With any given  $s$ , the value of  $\delta$  determines how the filtered centers are grouped. Basically, a smaller value of  $\delta$  results in fewer, and larger, groups of filtered centers. For example, start with  $\delta_0$  and suppose that  $g_1$  and  $g_2$  are two groups of the filtered centers, and  $c_j \in g_j$ , for  $j = 1, 2$  are the most similar centers among the two groups

$$\text{sim}(c_1, c_2) \geq \text{sim}(c_i, c_v) \text{ for any } c_i \in g_1, c_v \in g_2.$$

Set  $\tilde{\delta} = \text{sim}(c_1, c_2)$ , then  $\tilde{\delta} \leq \delta_0$ . The two groups  $g_1$  and  $g_2$  will combine into a single group if the value of  $\delta$  is updated to  $0.99\tilde{\delta}$ . This provides a useful tool to show the hierarchical structure in the clustering of data.

### 4.1 Hierarchical structure of DLCC

Heuristically, the initial choice for  $\delta$  can be based on the following value:

$$\delta_U := \min_{c_1} \max_{c_2} \{\text{sim}(c_1, c_2) : c_2 \neq c_1 \in \mathbf{C}\} \quad (11)$$

For any  $\delta \geq \delta_U$ , there will be at least one group with a single filtered center. It is thus desirable to start with a smaller value, say  $\delta_0 = 0.99\delta_U$ , so that the multiple smaller subsets are connected to form clusters. Suppose  $\delta_0$  is chosen and  $\mathbf{C}$  is now partitioned into groups  $g_1, g_2, \dots$ . For any two groups of filtered centers,  $g_i$  and  $g_j$ , their threshold similarity is the maximal similarity value between the centers in them:

$$\sigma(g_i, g_j) := \max\{\text{sim}(c_i, c_j) : c_i \in g_i, c_j \in g_j\},$$

which define the similarity matrix for  $\delta_0$ . Obviously, different values of  $\delta$  may produce different similarity matrices. Following the changes of the similarity matrices by varying  $\delta$  starting from  $\delta = \delta_0$ , it makes evident the hierarchical structure in the data.

We use a toy example to illustrate the process, which is presented in (12). In each step, the similarity matrix is obtained from the previous one by merging columns / rows that contain the largest non-trivial (not 0 or 1) similarity value. The merging process retains the larger of the values in the corresponding entries of the columns / rows. Eventually, there are left with completely disjoint groups, represented by the identity matrix as the similarity matrix. The process produces two corresponding lists of numbers:

- $G$ : the number of groups in each stage
- $\delta$ : the corresponding threshold values of  $\delta$

$$\begin{pmatrix} 1 & 0.6 & 0 & 0 & 0 & 0 \\ 0.6 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0.5 & 0.3 & 0.2 \\ 0 & 0 & 0.5 & 1 & 0.4 & 0.1 \\ 0 & 0 & 0.3 & 0.4 & 1 & 0.6 \\ 0 & 0 & 0.2 & 0.1 & 0.6 & 1 \end{pmatrix} \rightarrow \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0.5 & 0.3 \\ 0 & 0.5 & 1 & 0.4 \\ 0 & 0.3 & 0.4 & 1 \end{pmatrix} \rightarrow \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0.4 \\ 0 & 0.4 & 1 \end{pmatrix} \rightarrow \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \quad (12)$$

In this toy example, it is easy to see the list of  $G$  is  $\{6, 4, 3, 2\}$  and correspondingly  $\delta = \{0.6, 0.5, 0.4, 0\}$ . More precisely, the changes in groupings described in the matrices correspond to  $\delta$  values varying as below

$$\delta_U > \delta \geq 0.6 \longrightarrow 0.6 > \delta \geq 0.5 \longrightarrow 0.5 > \delta \geq 0.4 \longrightarrow 0.4 > \delta \geq 0$$

Hence, with fixed  $s$  and filtered centers, the potential values of  $G$  are fixed, and the minimal number of groups can be greater than one if there are no neighbors that overlap across groups.

Understanding this hierarchical structure provides us with a suitable  $\delta$  with respect to a chosen  $s$ , if the number of clusters is known for any reason. In general, without the knowledge of the number of clusters, we propose below an algorithm to heuristically select a suitable

value of  $\delta$ .

## 4.2 Heuristic guessing algorithm for $\delta$

Fix the parameter  $s$  and the set of filtered centers  $\mathbf{C} = \{c_1, c_2, \dots, c_P\}$ . Based on the discussion above and  $\delta$  obtained from  $\delta_0 = 0.99\delta_U$ , we describe a heuristic guessing algorithm for selecting the most suitable  $\delta$ , where the output of the algorithm is noted as  $\hat{\delta}$ . It depends on a parameter  $q$  (in Step 3), which has a default value set to be 10%.

Step 1 If  $\delta = \{0\}$ , then set  $\hat{\delta} = 0.99\delta_U$ . If the minimum number in the list of  $\mathbf{G}$  is larger than 1,  $\hat{\delta} = \min\{\delta \setminus \{0\}\}$ . Else, run left steps.

Step 2 For each  $c_j \in \mathbf{C}$  define  $N_j \subset \mathbf{C}$  as the collection of top  $\lambda = \min(2d, P/2)$  filtered centers in terms of similarities with  $c_j$  and compute the average similarity

$$m_j = \frac{1}{\lambda} \sum_{c \in N_j} \text{sim}(c_j, c).$$

Step 3 Let  $\tilde{\mathbf{C}}_{\text{edge}}$  consist of the filtered centers in the lowest  $q$  quantile (default to 10%) of the average similarities  $\mathbf{m} = \{m_1, \dots, m_P\}$ . The edge centers  $\mathbf{C}_{\text{edge}}$  is a maximal subset of  $\tilde{\mathbf{C}}_{\text{edge}}$  such that  $\max\{\text{sim}(c_i, \mathbf{C}_{\text{edge}} \setminus \{c_i\})\} > 0$  for all  $c_i \in \mathbf{C}_{\text{edge}}$ . It can be simply obtained by deleting filtered centers  $c$  with  $\max\{\text{sim}(c, \tilde{\mathbf{C}}_{\text{edge}} \setminus \{c\})\} = 0$  because the similarity between filtered centers is mutual.

Step 4 For any  $c_i$  in  $\mathbf{C}_{\text{edge}}$ , compute  $l_i = \max\{\text{sim}(c_i, \mathbf{C}_{\text{edge}} \setminus \{c_i\})\}$  and the gap value

$$\text{gap}_i = \max\{\text{sim}(c_i, \mathbf{C} \setminus \{c_i\})\} - l_i.$$

Step 5 Find the  $l_i$ 's that are also in  $\delta$ , among which, the output  $\hat{\delta}$  of the algorithm is the one corresponding to the largest gap value. If  $\delta$  does not contain any  $l_i$ 's,  $\hat{\delta}$  is set to 0.

The process is based on the intuitive observation that, in a reasonable grouping, if a group  $g$  contains a filtered center  $c_j$  and another filtered center, then  $g$  should also contain the filtered center(s) having the highest similarity with  $c_j$ . Thus  $\max\{\text{sim}(c_i, \mathbf{C} \setminus \{c_i\})\}$  is an indication

of the largest with-in group similarity. On the other hand, the filtered centers in  $C_{\text{edge}}$  are comparatively dissimilar with other filtered centers. The algorithm guesses two groups are cut between filtered centers in  $C_{\text{edge}}$ . Based on this, the  $l_i$  can be thought of as an indication of the largest between group similarity if  $c_i$  and its most similar filtered center in  $C_{\text{edge}}$  are assumed to be in two groups. The algorithm works by identifying the largest gap in these indications. As groupings only changes when  $\delta$  crosses the values in  $\delta$ , only these values are considered in the Step 5. If  $q$  is too small, it is possible that the algorithm cannot find a suitable  $\hat{\delta}$  in  $\delta$ , and it will return  $\hat{\delta} = 0$ , which indicates that there is no suitable partition among filtered centers in  $C_{\text{edge}}$ .

### 4.3 Density-based Clustering score

In sections 4.1 and 4.2, we provided an algorithm to determine the optimal  $\delta$  for each value of  $s$ . It is then important to have an internal metric for evaluating the performance of clustering determined by the pair of parameters  $(s, \delta)$ , because the ground truth is typically unavailable in real-world clustering applications. Well-known internal metrics such as the silhouette width (sw) [17] and the Calinski-Harabasz (CH) score [18] do not perform well in non-convex clustering problems or clustering problems with extremely nearby clusters.

Their limitations are illustrated using two synthetic data sets with different clustering results from various algorithms. In Figure 8, it concerns the Cassini data, and in Figure 9, the SmileyF data, which adds a circular “face” to the Smiley data. The average sw and the CH scores for those clustering results in the two data sets are presented in Table 4.

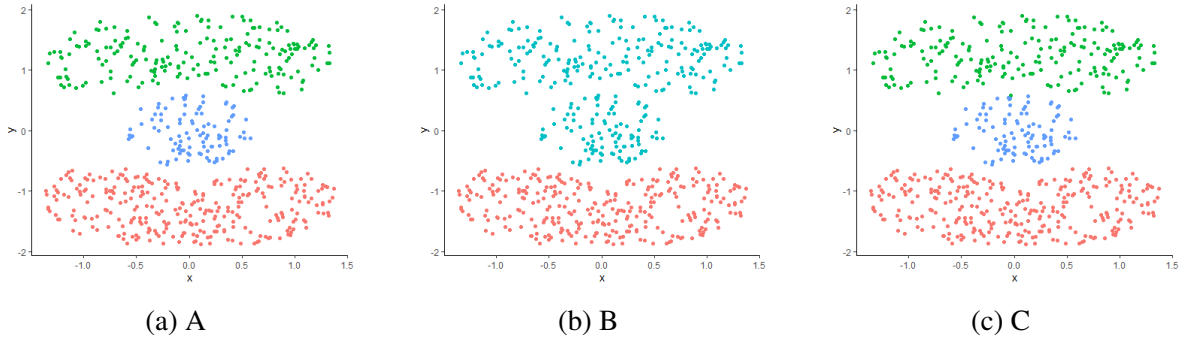


Figure 8: A-C are 3 clustering results in the synthetic data set Cassini.

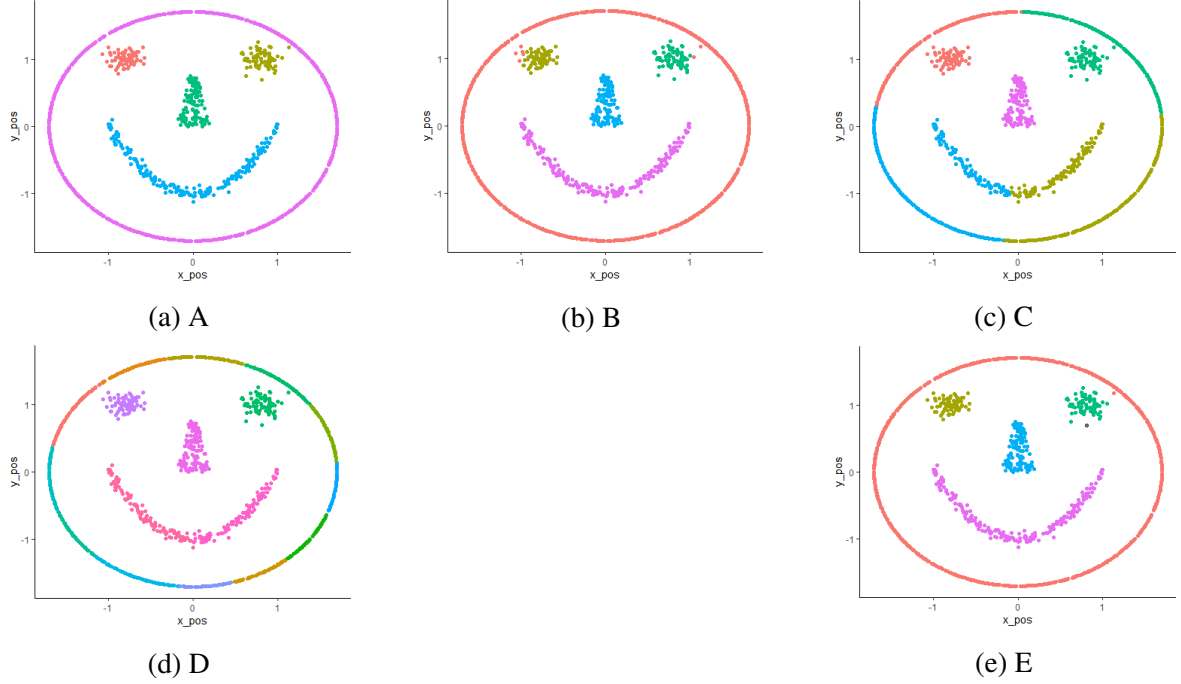


Figure 9: A-E are 5 clustering results in the synthetic data set SmileyF. One black point in E represents an unlabeled outlier.

Table 4: The average silhouette widths (asw) and the Calinski-Harabasz scores (CH) for clustering results in Figures 8 and 9.

| Clustering Result | A      | B      | C       | D       | E      |
|-------------------|--------|--------|---------|---------|--------|
| DATA SET          |        |        |         |         |        |
| Cassini: ASW      | 0.341  | 0.491  | 0.339   |         |        |
| Cassini: CH       | 642.13 | 842.52 | 641.99  |         |        |
| SmileyF: ASW      | -0.146 | -0.152 | 0.435   | 0.514   | -0.147 |
| SmileyF: CH       | 45.25  | 43.18  | 1470.62 | 2687.72 | 44.74  |

In both figures, the clustering results in A are the actual labels for the data sets, for which both asw and CH yield low scores. Both criteria favour the clustering result B for the Cassini data set, which contains only two clusters. Similarly, they awarded high scores to the clustering results C and D in the SmileyF data set, while based on human assessments, C and D perform the poorest. We propose the Density-based Clustering (DC) score for DLCC, a new internal metric for evaluating clustering performances. Instead of computing similarities globally, the idea is to calculate within and between cluster similarities in a small, fixed region, which works well with non-convex clusters.



Using within cluster / between cluster similarities in place of within cluster / between cluster distances, we recall some relevant notions below, analogous to those in DBCA [8]:

**Definition 4.1** (Core neighbors). *Given a similarity threshold  $\eta$ , a point  $\mathbf{x}_j$  is a core neighbor of  $\mathbf{x}_i$  (with respect to  $\eta$ ) if  $D_M(\mathbf{x}_j \mid \mathbf{x}_i) \geq \eta$ .*

**Definition 4.2** (Depth-connection). *A point  $\mathbf{x}_v$  is depth-connected to  $\mathbf{x}_i$  if there is a chain of points  $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_y$  where each two adjacent points are core neighbors for each other, and  $\tilde{\mathbf{x}}_1 = \mathbf{x}_v$  and  $\tilde{\mathbf{x}}_p = \mathbf{x}_i$ .*

**Definition 4.3** (DBCA-Cluster). *A chain  $\tilde{\mathbf{C}} = \{\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_y\}$  is called a maximal depth-connected chain if there is no  $\mathbf{x}_j \notin \tilde{\mathbf{C}}$  satisfying  $\sup_{\mathbf{x}_i \in \tilde{\mathbf{C}}} D_M(\mathbf{x}_i \mid \mathbf{x}_j) \geq \eta$ . Let  $\tilde{n}$  be the minimum number of points required to form a cluster in DBCA, then each cluster is formed by observations that are in the same maximal depth-connected chain with length at least  $\tilde{n}$ .*

**Definition 4.4** (Outliers). *All observations remain unlabelled in DBCA are regarded as outliers.*

To prevent misunderstanding, similarity in this context refers to the similarity between points, not groups or filtered centres, as defined in Equation (1). Additionally, Definition 4.3 is mentioned for defining outliers. It is important to note that the clustering results that the DC score examines can be produced by a variety of algorithms and are not necessarily DBCA, so the clustering results do not need to satisfy Definition 4.3.

In Definition 4.1, larger  $\eta$  corresponds to fewer core neighbors, which corresponds to a smaller fixed region. Let  $\eta_X$  be the largest value of  $\eta$  such that all observations in the dataset are depth-connected. Let  $\zeta$  be a cluster in the clustering result which is examined by DC score, and  $\eta_\zeta$  be the largest value of  $\eta$  such that  $\zeta$  is decomposed into a single DBCA-cluster and outliers (without considering observations in other clusters). Let  $n_\zeta$  be the number of points in  $\zeta$  and  $\mathbf{o}_\zeta$  be the set of DBCA-outliers. Let  $\mathcal{X}_\zeta$  be the set of observations in cluster  $\zeta$ , then we

define the within cluster similarity  $J_\zeta$  of  $\zeta$  by (13).

$$\begin{aligned}
 J_\zeta &= \sum_{i=1}^{n_\zeta} w_i \mu_i, \quad \text{where} \\
 \mu_i &= \begin{cases} \frac{\sum \{D_M(\mathbf{x}_j|\mathbf{x}_i) | D_M(\mathbf{x}_j|\mathbf{x}_i) \geq \eta_\zeta\}}{\#\{\mathbf{x}_j | D_M(\mathbf{x}_j|\mathbf{x}_i) \geq \eta_\zeta\}}, & \text{if } \mathbf{x}_i \notin \mathbf{o}_\zeta \\ \max\{D_M(\mathbf{x}_j | \mathbf{x}_i) | \mathbf{x}_j \notin \mathbf{o}_\zeta\}, & \text{if } \mathbf{x}_i \in \mathbf{o}_\zeta \end{cases} \\
 w_i &= \begin{cases} \frac{\#\{\mathbf{x}_j | D_M(\mathbf{x}_j|\mathbf{x}_i) \geq \eta_\zeta\}}{\sum_{\mathbf{x}_v \notin \mathbf{o}_\zeta} \#\{\mathbf{x}_j | D_M(\mathbf{x}_j|\mathbf{x}_v) \geq \eta_\zeta\} + \#\{\mathbf{o}_\zeta\}}, & \text{if } \mathbf{x}_i \notin \mathbf{o}_\zeta \\ \frac{1}{\sum_{\mathbf{x}_v \notin \mathbf{o}_\zeta} \#\{\mathbf{x}_j | D_M(\mathbf{x}_j|\mathbf{x}_v) \geq \eta_\zeta\} + \#\{\mathbf{o}_\zeta\}}, & \text{if } \mathbf{x}_i \in \mathbf{o}_\zeta \end{cases}
 \end{aligned} \tag{13}$$

where  $\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_v \in \mathcal{X}_\zeta$ . By considering DBCA-outliers, the DC score can be made more robust by avoiding the significant adverse impact of an individual outlier.

In the cluster  $\zeta$ , the set of inner observation of  $\zeta$  is

$$\mathbf{I}_\zeta := \{\mathbf{x}_i \in \mathcal{X}_\zeta : D_M(\mathbf{x}_j | \mathbf{x}_i) < \eta_X \text{ for all } \mathbf{x}_j \in \mathbf{X} \setminus \mathcal{X}_\zeta\},$$

where  $\mathbf{X}$  is the total data set. If  $\mathbf{x}_i \in \mathbf{I}_\zeta$  then  $\mathbf{x}_i$  is said as an inner observation of cluster  $\zeta$ .

The between cluster similarity  $H_\zeta$  for cluster  $\zeta$  is defined as (14).

$$\begin{aligned}
 H_\zeta &= \sum_{i=1}^{n_\zeta} \tilde{w}_i \tilde{\mu}_i, \quad \text{where} \\
 \tilde{\mu}_i &= \begin{cases} \frac{\sum \{D_M(\mathbf{x}_h|\mathbf{x}_i) | D_M(\mathbf{x}_h|\mathbf{x}_i) \geq \eta_X\}}{\#\{\mathbf{x}_h | D_M(\mathbf{x}_h|\mathbf{x}_i) \geq \eta_X\}}, & \text{if } \mathbf{x}_i \notin \mathbf{I}_\zeta \\ \eta_X, & \text{if } \mathbf{x}_i \in \mathbf{I}_\zeta \end{cases} \\
 \tilde{w}_i &= \begin{cases} \frac{\#\{\mathbf{x}_h | D_M(\mathbf{x}_h|\mathbf{x}_i) \geq \eta_X\}}{\sum_{\mathbf{x}_u \notin \mathbf{I}_\zeta} \#\{\mathbf{x}_u | D_M(\mathbf{x}_h|\mathbf{x}_u) \geq \eta_X\} + \#\{\mathbf{I}_\zeta\}}, & \text{if } \mathbf{x}_i \notin \mathbf{I}_\zeta \\ \frac{1}{\sum_{\mathbf{x}_u \notin \mathbf{I}_\zeta} \#\{\mathbf{x}_u | D_M(\mathbf{x}_h|\mathbf{x}_u) \geq \eta_X\} + \#\{\mathbf{I}_\zeta\}}, & \text{if } \mathbf{x}_i \in \mathbf{I}_\zeta \end{cases}
 \end{aligned} \tag{14}$$

where  $\mathbf{x}_i \in \mathcal{X}_\zeta$  and  $\mathbf{x}_h, \mathbf{x}_u \notin \mathcal{X}_\zeta$ .

The total DC score for a clustering result is then given by

$$\text{DC} = \sum_{\zeta} \frac{n_\zeta J_\zeta}{H_\zeta}, \tag{15}$$

where  $n_\zeta$  weighs more the clustering performance for the cluster with more observations. Weight function  $w_i$  in  $J_\zeta$  gives observations with more core neighbors more weight and less weights for outliers. In  $H_\zeta$ , weight function  $\tilde{w}_i$  gives observations which are close to other clusters more weights, but less weights for inner observations. Some properties of DC score are listed below:

Property 1: For a cluster  $\zeta$ , if  $\eta'_\zeta < \eta_\zeta$ , then  $\mu'_i \leq \mu_i$  for any  $\mathbf{x}_i \in \mathcal{X}_\zeta$ .

Property 2: For a cluster  $\zeta$ , if  $\eta'_\zeta < \eta_\zeta$ , then  $\tilde{\mu}'_i \leq \tilde{\mu}_i$  for any  $\mathbf{x}_i \in \mathcal{X}_\zeta$ .

Property 3: For any cluster  $\zeta$ , the  $\mu_i$  given by any outlier is less than  $\eta_\zeta$ .

Property 4: For any cluster  $\zeta$ ,  $H_\zeta \geq \eta_X$ . Higher proportion of inner observations in  $\zeta$  leads to smaller  $H_\zeta$ .

The proofs of Properties 1 and 2 are shown in Appendix 6.

Table 5 presents the DC scores for above examples of clustering results of Cassini and SmileyF, which seems that the ordering of the scores now better matches human assessments.

Table 5: The DC scores for clustering results in Figures 8 and 9, where  $\tilde{n} = 3$ .

|         | A       | B       | C       | D       | E       |
|---------|---------|---------|---------|---------|---------|
| Cassini | 629.73  | 628.57  | 629.05  |         |         |
| SmileyF | 1822.27 | 1814.40 | 1666.02 | 1752.26 | 1821.83 |

We now turn our attention to the final phase of selecting the best pair of DLCC parameters using the proposed internal metric. All derived parameter pairs are tested using the DC score for measuring clustering performances. If the parameters  $\hat{\delta}$  used are obtained from the heuristic guessing algorithm in Section 4.2, the clustering will be also done using  $0.99\hat{\delta}$  for comparison. The parameter pair that produces the results with the greatest DC score is selected.

## 5 Projection Depth with a fixed center

The projection depth proposed by Zuo and Serfling [19] has some favourable properties, such as a relatively high breakdown point for the sample projection median. For a point  $z$  with respect to a data set  $\mathbf{X}$ , it is defined by

$$PD(z|\mathbf{X}) = \left( 1 + \sup_{\|\mathbf{u}\|=1} \frac{|\mathbf{u}^T z - \text{Med}(\mathbf{u}^T \mathbf{X})|}{\text{MAD}(\mathbf{u}^T \mathbf{X})} \right)^{-1}, \quad (16)$$

where Med stands for the median and MAD is the median absolute deviation defined as

$$\text{MAD}(U) = \text{Med}(|U - \text{Med}(U)|)$$

for a univariate random variable  $U$ , and

$$\mathbf{u}^T \mathbf{X} = \{u^T \mathbf{x}_1, u^T \mathbf{x}_2, \dots, u^T \mathbf{x}_n\}.$$

The second term in (16) represents the Stahel-Donoho outlyingness with respect to  $\mathbf{X}$ :

$$O_z = \sup_{\|\mathbf{u}\|=1} \frac{|\mathbf{u}^T z - \text{Med}(\mathbf{u}^T \mathbf{X})|}{\text{MAD}(\mathbf{u}^T \mathbf{X})}$$

hence (16) can also be written as  $(1 + O_z)^{-1}$ .

Inspired by the idea of using the Mahalanobis depth definition to measure similarities among observations, we propose a fixed center version of the projection depth:

$$\widetilde{PD}(z|\mathbf{X}, \mathbf{x}_j) = \left( 1 + \sup_{\|\mathbf{u}\|=1} \frac{|\mathbf{u}^T z - (\mathbf{u}^T \mathbf{x}_j)|}{\text{Med}(|\mathbf{u}^T \mathbf{X} - \mathbf{u}^T \mathbf{x}_j|)} \right)^{-1}, \quad (17)$$

where  $\widetilde{PD}(\mathbf{x}_j|\mathbf{X}, \mathbf{x}_j) = 1$  and  $\text{Med}(\mathbf{u}^T \mathbf{X})$  in (16) is replaced by  $\mathbf{u}^T \mathbf{x}_j$ . The depth value of (17) can be explained as the similarity of other observations to a single point.

We describe below the overall concept for computing the projection depth with a fixed center, which is based on Liu and Zuo's exact algorithm for calculating projection depth, under the assumption that the data are in general position [20, 21]. For convenience, the number of

observations in  $\mathbf{X}$  is assumed to be odd in latter explanations (The even case is similar, see the detailed algorithm for computing the projection depth in [20, 21]). Firstly we consider

$$\tilde{O}_j = \sup_{\|\mathbf{u}\|} \frac{|\mathbf{u}^T \mathbf{z} - (\mathbf{u}^T \mathbf{x}_j)|}{\text{Med}(|\mathbf{u}^T \mathbf{X} - \mathbf{u}^T \mathbf{x}_j|)}. \quad (18)$$

For simplicity, we say

$$Q_j = \frac{\mathbf{u}^T \mathbf{z} - (\mathbf{u}^T \mathbf{x}_j)}{\text{Med}(|\mathbf{u}^T \mathbf{X} - \mathbf{u}^T \mathbf{x}_j|)}.$$

$Q_j$  is odd with respect to  $\mathbf{u}$ , which represents that finding the superior of  $|Q_j|$  is equivalent with finding the extreme values of  $Q_j$ . For any  $\mathbf{u}$ , there is a permutation in  $\mathbf{u}^T \mathbf{X}$ , and the indices can be rearranged to  $\{b_1, b_2, \dots, b_n\}$ , such that

$$\mathbf{u}^T \mathbf{x}_{b_1} \leq \mathbf{u}^T \mathbf{x}_{b_2}, \dots, \leq \mathbf{u}^T \mathbf{x}_{b_n}.$$

Note the current index of  $\mathbf{u}^T \mathbf{x}_j$  is  $b_m$ , then a set of  $\mathbf{u}$  can be defined as  $\mathbb{U} = \{\mathbf{u} | \mathbb{A}^T \mathbf{u} \leq \mathbf{0}\}$ , where  $\mathbb{A} = \{\mathbf{x}_{b_1} - \mathbf{x}_{b_m}, \dots, \mathbf{x}_{b_{m-1}} - \mathbf{x}_{b_m}, \mathbf{x}_{b_m} - \mathbf{x}_{b_{m+1}}, \dots, \mathbf{x}_{b_m} - \mathbf{x}_{b_n}\}$ . Simply speaking, the position of  $\mathbf{x}_j$  in the permutation keeps the same when  $\mathbf{u} \in \mathbb{U}$ . Furthermore, the whole space of  $\mathbf{u}$  can be divided into finite sets  $\mathbb{U}_1, \mathbb{U}_2, \dots, \mathbb{U}_U$ . Based on the definition, when  $\mathbf{u}$  moves from  $\mathbb{U}_1$  to another region  $\mathbb{U}_2$ ,  $b_m$  will change. Now reverting to a single set  $\mathbb{U}$ , for any  $\mathbf{u} \in \mathbb{U}$ , there is

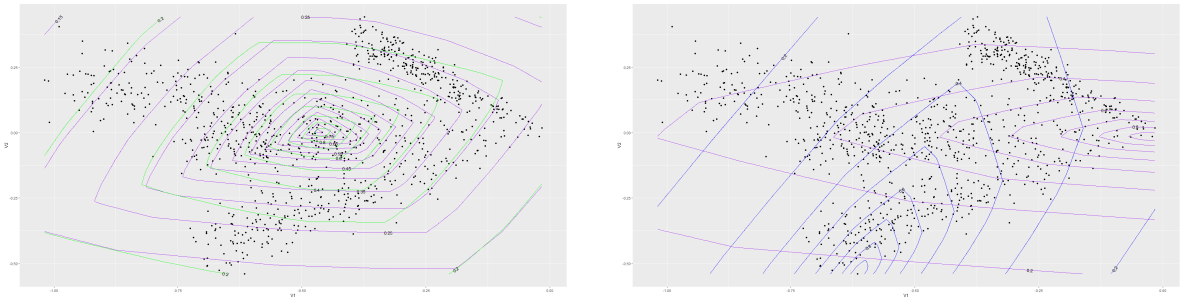
$$\begin{aligned} & |\mathbf{u}^T \mathbf{x}_l - \mathbf{u}^T \mathbf{x}_j| \\ &= \begin{cases} -(\mathbf{u}^T \mathbf{x}_l - \mathbf{u}^T \mathbf{x}_j), & \text{if } l \in \{b_1, b_2, \dots, b_{m-1}\}, \\ \mathbf{u}^T \mathbf{x}_l - \mathbf{u}^T \mathbf{x}_j, & \text{if } l \in \{b_m, b_{m+1}, \dots, b_n\}. \end{cases} \end{aligned}$$

Similarly, for any  $\mathbf{u}$  there is a permutation for  $|\mathbf{u}^T \mathbf{X} - \mathbf{u}^T \mathbf{x}_j|$ . Moreover, a region  $\mathbb{D} \subset \mathbb{U}$  can be defined, which ensures the index of the observation (for example the index is  $v$ ) which satisfies  $|\mathbf{u}^T \mathbf{x}_v - \mathbf{u}^T \mathbf{x}_j| = \text{Med}(|\mathbf{u}^T \mathbf{X} - \mathbf{u}^T \mathbf{x}_j|)$ , is fixed. The median of  $|\mathbf{u}^T \mathbf{X} - \mathbf{u}^T \mathbf{x}_j|$  in  $\mathbb{D}$  now can be noted as  $\mathbf{u}^T V$ , where  $V = \mathbf{x}_v - \mathbf{x}_j$  if  $v \in \{b_m, b_{m+1}, \dots, b_n\}$ , otherwise,  $V = \mathbf{x}_j - \mathbf{x}_v$ . Then  $Q_j = \frac{\mathbf{u}^T (\mathbf{z} - \mathbf{x}_j)}{\mathbf{u}^T V}$  over  $\mathbb{D}$ , which illustrates that  $Q_j$  is a piecewise linear fractional function over a finite number of pieces like  $\mathbb{D}$ , and for each piece the maximum

value of  $Q_j$  only happens at the edge [22, cited in [21]].

In the two dimensional case, unit vectors form a circle, and any unit vector can be represented by  $\mathbf{u} = (\cos \alpha, \sin \alpha)$ , where  $\alpha \in [0, 2\pi]$ . In fact, due to the symmetry, we only need to consider  $\alpha \in [0, \pi]$ . The circle is split into many pieces to find local maximum values of  $Q_j$ . We start with  $\alpha = 0$  and make  $\mathbf{u}$  pass through all pieces counter-clockwise, and when  $\mathbf{u}$  passes through the edge of each piece,  $\mathbf{x}_{\text{med}}$ , which represents the observation in the median position of  $|\mathbf{u}^T \mathbf{X} - \mathbf{u}^T \mathbf{x}_j|$  will change. For each time the position of  $\mathbf{x}_j$  changes in the permutation of  $\mathbf{u}^T \mathbf{X}$ ,  $b_m$  is either plus one or minus one depending on which observation switches position with it. Those unit vectors located at the edges can be used to calculate local maximum values, and as a result, we can find the global maximum value of  $Q_j$  from all local maxima to calculate the depth value. Lai and Zuo (2011) provided a detailed explanation for calculating the projection depth in two dimensions [23].

Figure 10 reveals contour plots of (16) and (17) in the Ray data mentioned in Section 3, where the depth values are calculated by the first two variables of the Ray data. From Figure 10a, the contour plot of the original projection depth (PD) and the projection depth with the fixed center (PDfix) exhibit similar shapes, while they are slightly different in orientation. Moreover, PDfix's contour plot is also convex and polylateral. Figure 10b illustrates how the contour plots of PDfix can have completely various orientations depending on the center points chosen, indicating that the orientation of the contour will be automatically modified to account for the positions of other points.



(a) Original projection depth (Purple) VS Projection depth with the fixed center (The chosen center is the deepest point given by (16), Green)

(b) Projection depth with the fixed center (small  $y$  value, Blue) VS Projection depth with the fixed center (large  $x$  value, Purple)

Figure 10: Projection depth's contour plots in Ray data set with first two variables.

The majority of the properties of PD, like Affine invariance, Null at infinity, Monotone on rays, and so forth, are inherited by PDfix. The idea of fixed center version of data depth may have the potential to expand the applications of data depth, which is one of our future directions.

## 6 Appendix

**Lemma 6.1.** *For a cluster  $\zeta$ , if  $\eta_{\zeta'} < \eta_\zeta$ , then  $\mu_{i'} \leq \mu_i$  for any  $\mathbf{x}_i \in \mathcal{X}_\zeta$ . Further, if there is no outlier, the average weighted within cluster similarity for observations  $\mathbf{x}_i \in \mathcal{X}_\zeta$  will remain the same or become smaller for smaller  $\eta$  value, i.e.,  $\tilde{J}_{\zeta'} \leq J_\zeta$  (Note that  $\tilde{J}_{\zeta'} \neq J_{\zeta'}$ , because different  $\eta$  value means different clustering results).*

*Proof.* Let  $\eta_{\zeta'} < \eta_\zeta$ , for any observation  $\mathbf{x}_i \in \mathcal{X}_\zeta$  and  $\notin \mathbf{o}_\zeta$ , there are

$$\mu_i = \frac{Q_i}{n_{Q_i}}, \quad (19)$$

$$\mu_{i'} = \frac{Q_i + W_i}{n_{Q_i} + n_{W_i}}, \quad (20)$$

where  $Q_i = \sum \{D_M(\mathbf{x}_j | \mathbf{x}_i) \mid D_M(\mathbf{x}_j | \mathbf{x}_i) \geq \eta_\zeta\}$ ,  $n_{Q_i} = \#\{\mathbf{x}_j \mid D_M(\mathbf{x}_j | \mathbf{x}_i) \geq \eta_\zeta\}$ ,  $W_i = \sum \{D_M(\mathbf{x}_j | \mathbf{x}_i) \mid \eta_{\zeta'} \leq D_M(\mathbf{x}_j | \mathbf{x}_i) < \eta_\zeta\}$  and  $n_{W_i} = \#\{\mathbf{x}_j \mid D_M(\mathbf{x}_j | \mathbf{x}_i) \leq \eta_{\zeta'} \leq D_M(\mathbf{x}_j | \mathbf{x}_i) < \eta_\zeta\}$ . If  $n_{W_i} = 0$ ,  $\mu_i = \mu_{i'}$ . For the case  $n_{W_i} > 0$ ,  $W_i, Q_i, n_{W_i}, n_{Q_i}$  are all positive, there are

$$\eta_{\zeta'} \leq W_i/n_{W_i} < \eta_\zeta \leq Q_i/n_{Q_i}, \quad (21)$$

$$\begin{aligned} W_i n_{Q_i} &< Q_i n_{W_i}. \\ \mu_i - \mu_{i'} &= \frac{Q_i n_{W_i} - W_i n_{Q_i}}{n_{Q_i}(n_{Q_i} + n_{W_i})} > 0. \end{aligned} \quad (22)$$

Therefore,  $\mu_i > \mu_{i'}$  when  $\eta_{\zeta'} < \eta_\zeta$ .

Without considering outliers, for the within cluster similarity, there are

$$J_\zeta = \frac{\sum Q_i}{\sum n_{Q_i}}, \quad (23)$$

$$\tilde{J}_{\zeta'} = \frac{\sum (Q_i + W_i)}{\sum (n_{Q_i} + n_{W_i})}. \quad (24)$$

To examine if  $J_\zeta > \tilde{J}_{\zeta'}$ , we can begin with a simple case with only two observations in a



cluster  $\zeta$ , which is shown in (25).

$$\frac{Q_1 + Q_2}{n_{Q_1} + n_{Q_2}} - \frac{Q_1 + W_1 + Q_2 + W_2}{n_{Q_1} + n_{W_1} + n_{Q_2} + n_{W_2}} = \frac{n_{W_1}Q_1 - n_{Q_1}W_1 + n_{W_2}Q_2 - n_{Q_2}W_2 + n_{W_2}Q_1 - n_{Q_1}W_2 + n_{W_1}Q_2 - n_{Q_2}W_1}{(n_{Q_1} + n_{Q_2})(n_{Q_1} + n_{W_1} + n_{Q_2} + n_{W_2})}, \quad (25)$$

where  $n_{W_1}$  and  $n_{W_2}$  are non-zero. Following (21),  $n_{W_1}Q_1 - n_{Q_1}W_1 > 0$  and  $n_{W_2}Q_2 - n_{Q_2}W_2 > 0$ . Moreover, all  $Q_i/n_{Q_i} \geq \eta_\zeta$  and all  $W_j/n_{W_j} < \eta_\zeta$  with  $n_{W_j} \neq 0$ ; Thereby, even when  $i \neq j$ ,  $n_{Q_i}W_j < n_{W_j}Q_i$  holds. Then there are  $n_{W_2}Q_1 - n_{Q_1}W_2 > 0$  and  $n_{W_1}Q_2 - n_{Q_2}W_1 > 0$ . Hence,

$$\frac{Q_1 + Q_2}{n_{Q_1} + n_{Q_2}} - \frac{Q_1 + W_1 + Q_2 + W_2}{n_{Q_1} + n_{W_1} + n_{Q_2} + n_{W_2}} > 0. \quad (26)$$

Similarly, if  $Q_1 + Q_2$ ,  $n_{Q_1} + n_{Q_2}$ ,  $W_1 + W_2$  and  $n_{W_1} + n_{W_2}$  are written as  $Q_{1 \sim 2}$ ,  $n_{Q_{1 \sim 2}}$ ,  $W_{1 \sim 2}$  and  $n_{W_{1 \sim 2}}$ , it is easy to show  $J_\zeta - \tilde{J}_{\zeta'} > 0$  with more points. Therefore,  $J_\zeta > \tilde{J}_{\zeta'}$ . Besides,  $J_\zeta = \tilde{J}_{\zeta'}$  if and only if  $\sum n_{W_i} = 0$ .  $\square$

**Lemma 6.2.** *For a cluster  $\zeta$ , if  $\eta_{\zeta'} < \eta_\zeta$  then  $\tilde{\mu}_{i'} \leq \tilde{\mu}_i$  for any  $\mathbf{x}_i \in \mathcal{X}_\zeta$ . Further, if core neighbors of all  $\mathbf{x}_i \in \mathcal{X}_{\zeta'}$  are in  $\mathcal{X}_{\zeta'}$ , then the average between cluster similarity for observations  $\mathbf{x}_i \in \mathcal{X}_\zeta$  will remain the same or become smaller for smaller  $\eta$  value, i.e.,  $\tilde{H}_{\zeta'} \leq H_\zeta$  (Note that  $\tilde{H}_{\zeta'} \neq H_{\zeta'}$ , because different  $\eta$  value means different clustering results).*

*Proof.* Let  $\eta_\zeta > \eta_{\zeta'} \geq \eta_X$  and assume for every  $\mathbf{x}_j$  satisfying  $D_M(\mathbf{x}_j | \mathbf{x}_i) \geq \eta_{\zeta'}$ , there is  $\mathbf{x}_j \in \mathcal{X}_{\zeta'}$ . Then for any  $\mathbf{x}_i \in \mathcal{X}_\zeta$ , there is,

$$\begin{aligned} \{\mathbf{x}_j | D_M(\mathbf{x}_j | \mathbf{x}_i) \geq \eta_{\zeta'}\} &= \{\mathbf{x}_v | D_M(\mathbf{x}_v | \mathbf{x}_i) \geq \eta_\zeta\} + \\ &\quad \{\mathbf{x}_h | \eta_{\zeta'} \leq D_M(\mathbf{x}_h | \mathbf{x}_i) < \eta_\zeta\}, \end{aligned}$$

where  $\mathbf{x}_j \in \mathcal{X}_{\zeta'}$ ,  $\mathbf{x}_v \in \mathcal{X}_\zeta$  and  $\mathbf{x}_h \notin \mathcal{X}_\zeta$ , which means with a higher  $\eta$ , there may exist some

observations being assigned to other clusters.

$$\tilde{\mu}_{i'} = \frac{V_i}{n_{V_i}}, \quad (27)$$

$$\tilde{\mu}_i = \frac{V_i + L_i}{n_{V_i} + n_{L_i}}, \quad (28)$$

$$\tilde{H}_{\zeta'} = \frac{\sum V_i + n_{I_1} \eta_X}{\sum n_{V_i} + n_{I_1}}, \quad (29)$$

$$H_{\zeta} = \frac{\sum V_i + \sum L_i + n_{I_2} \eta_X}{\sum n_{V_i} + \sum n_{L_i} + n_{I_2}}, \quad (30)$$

where  $V_i = \sum \{D_M(\mathbf{x}_h \mid \mathbf{x}_i) \mid \eta_{\zeta'} > D_M(\mathbf{x}_h \mid \mathbf{x}_i) \geq \eta_X\}$ ,  $n_{V_i} = \#\{\mathbf{x}_h \mid \eta_{\zeta'} > D_M(\mathbf{x}_h \mid \mathbf{x}_i) \geq \eta_X\}$ ,  $n_{I_1} = \#\{\mathbf{I}_{\zeta'}\}$ ,  $L_i = \sum \{D_M(\mathbf{x}_h \mid \mathbf{x}_i) \mid \eta_{\zeta'} \leq D_M(\mathbf{x}_h \mid \mathbf{x}_i) < \eta_{\zeta}\}$ ,  $n_{L_i} = \#\{\mathbf{x}_h \mid \eta_{\zeta'} \leq D_M(\mathbf{x}_h \mid \mathbf{x}_i) < \eta_{\zeta}\}$  and  $n_{I_2} = \#\{\mathbf{I}_{\zeta}\}$ . Similar with the proof of Property 1, it is easy to show  $\tilde{\mu}_{i'} - \tilde{\mu}_i < 0$  when  $n_{L_i} > 0$  and  $\tilde{\mu}_{i'} = \tilde{\mu}_i$  when  $n_{L_i} = 0$ . For the case  $\sum n_{L_i} > 0$ , consider

$$\begin{aligned} & \frac{\sum V_i + n_{I_1} \eta_X}{\sum n_{V_i} + n_{I_1}} - \frac{\sum V_i + \sum L_i + n_{I_2} \eta_X}{\sum n_{V_i} + \sum n_{L_i} + n_{I_2}} = \\ & \frac{\sum V_i \sum n_{L_i} - \sum L_i \sum n_{V_i} + n_{I_1}(\sum n_{L_i} \eta_X - \sum L_i) + (n_{I_1} - n_{I_2})(\sum n_{V_i} \eta_X - \sum V_i)}{(\sum n_{V_i} + n_{I_1})(\sum n_{V_i} + \sum n_{L_i} + n_{I_2})}. \end{aligned}$$

Because  $\eta_X \leq \sum V_i / \sum n_{V_i} < \eta_{\zeta'} \leq \sum L_i / \sum n_{L_i}$ , there are  $\sum V_i \sum n_{L_i} - \sum L_i \sum n_{V_i} < 0$  and  $n_{I_1}(\sum n_{L_i} \eta_X - \sum L_i) < 0$ . Since  $\eta_{\zeta'} < \eta_{\zeta}$  (With larger  $\eta$ , more similar points can be assigned to other clusters, and hence less number of inner observations.),  $n_{I_1} \geq n_{I_2}$ , and we know  $\sum \eta_X \leq \sum L_i / n_{L_i}$ , then there is  $(n_{I_1} - n_{I_2})(\sum n_{V_i} \eta_X - \sum V_i) \leq 0$ . Therefore,  $\tilde{H}_{\zeta'} < H_{\zeta}$  when  $\sum n_{L_i} > 0$ . If  $\sum n_{L_i} = 0$ ,  $H_{\zeta} = \tilde{H}_{\zeta'}$ .  $\square$

## References

- [1] Karl Mosler. Depth statistics. In Robustness and complex data structures, pages 17–34. Springer, 2013.
- [2] Ricardo Fraiman, Regina Y Liu, and Jean Meloche. Multivariate density estimation by probing depth. Lecture Notes-Monograph Series, pages 415–430, 1997.
- [3] Regina Y Liu and Kesar Singh. Rank tests for multivariate scale difference based on data depth. DIMACS series in discrete mathematics and theoretical computer science, 72:17, 2006.
- [4] David L Donoho and Miriam Gasko. Breakdown properties of location estimates based on halfspace depth and projected outlyingness. The Annals of Statistics, pages 1803–1827, 1992.
- [5] Rebecka Jörnsten. Clustering and classification based on the l1 data depth. Journal of Multivariate Analysis, 90(1):67–89, 2004.
- [6] Karl Mosler and Richard Hoberg. Data analysis and classification with the zonoid depth. DIMACS Series in Discrete Mathematics and Theoretical Computer Science, 72:49, 2006.
- [7] Oleksii Pokotylo, Pavlo Mozharovskyi, and Rainer Dyckerhoff. Depth and depth-based classification with r-package ddalpha. arXiv preprint arXiv:1608.04109, 2016.
- [8] Myeong-Hun Jeong, Yaping Cai, Clair J Sullivan, and Shaowen Wang. Data depth based clustering analysis. In Proceedings of the 24th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems, pages 1–10, 2016.
- [9] Martin Ester, Hans-Peter Kriegel, Jörg Sander, Xiaowei Xu, et al. A density-based algorithm for discovering clusters in large spatial databases with noise. In kdd, volume 96, pages 226–231, 1996.
- [10] Mia Hubert and Michiel Debruyne. Minimum covariance determinant. Wiley interdisciplinary reviews: Computational statistics, 2(1):36–43, 2010.

- [11] Richard G Brereton. The mahalanobis distance and its relationship to principal component scores. Journal of Chemometrics, 29(3):143–145, 2015.
- [12] Paul D McNicholas. Mixture model-based classification. Chapman and Hall/CRC, 2016.
- [13] Dheeru Dua and Casey Graff. UCI machine learning repository, 2017.
- [14] Francesco Sanna Passino, Nicholas A Heard, and Patrick Rubin-Delanchy. Spectral clustering on spherical coordinates under the degree-corrected stochastic blockmodel. Technometrics, pages 1–12, 2022.
- [15] Lawrence Hubert and Phipps Arabie. Comparing partitions. Journal of classification, 2(1):193–218, 1985.
- [16] Friedrich Leisch and Evgenia Dimitriadou. mlbench: Machine Learning Benchmark Problems, 2021. R package version 2.1-3.
- [17] Peter J Rousseeuw. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. Journal of computational and applied mathematics, 20:53–65, 1987.
- [18] Tadeusz Caliński and Jerzy Harabasz. A dendrite method for cluster analysis. Communications in Statistics-theory and Methods, 3(1):1–27, 1974.
- [19] Yijun Zuo. Projection-based depth functions and associated medians. The Annals of Statistics, 31(5):1460–1490, 2003.
- [20] Xiaohui Liu, Yijun Zuo, and Zhizhong Wang. Exactly computing bivariate projection depth contours and median. Computational Statistics & Data Analysis, 60:1–11, 2013.
- [21] Xiaohui Liu and Yijun Zuo. Computing projection depth and its associated estimators. Statistics and Computing, 24(1):51–63, 2014.
- [22] Kanti Swarup. Linear fractional functionals programming. Operations Research, 13(6):1029–1036, 1965.

- 
- [23] Yijun Zuo and Shaoyong Lai. Exact computation of bivariate projection depth and the stahel–donoho estimator. Computational Statistics & Data Analysis, 55(3):1173–1179, 2011.
- [24] Yijun Zuo and Robert Serfling. General notions of statistical depth function. Annals of statistics, pages 461–482, 2000.