

# MAFIA: An Adaptive Grid-Based Subspace Clustering Approach

Henrik Daniel Christensen<sup>[hench13@student.sdu.dk]</sup>

University of Southern Denmark, SDU  
Department of Mathematics and Computer Science

**Abstract.** This paper presents a comprehensive analysis of three bottom-up subspace clustering algorithms: *MAFIA*, *CLIQUE*, and *SUBCLU*. *MAFIA* extends the grid-based approach of *CLIQUE* by introducing adaptive grid sizes, offering improved scalability and clustering quality. *SUBCLU*, in contrast, utilizes a density-connectivity method that allows better identification of arbitrarily shaped clusters, overcoming some of the limitations inherent in grid-based algorithms. The paper explores the relationships and distinctions between these algorithms, evaluating their performance in terms of scalability of dataset size, data- and cluster-dimensionality, as well as their clustering quality. Through a series of experiments with both synthetic and real-world datasets, we demonstrate the strengths and limitations of each approach. The findings reveal that while *MAFIA* is the most scalable, *SUBCLU* excels at detecting clusters with irregular shapes. However, all algorithms require proper parameter tuning for optimal performance.

**Keywords:** High-Dimensional Subspace Clustering · Grid-Based- and Density-Connectivity approach · Comparative Study.

## 1 Introduction

*Clustering* is one of the main techniques within data mining. The main goal is to discover unknown patterns within a data set, by partitioning the data objects into *clusters*. Here, each object in a cluster is similar to one another, but different from objects in other clusters. Clustering is widely used in many applications, such as advertising, biology, web search and business intelligence [?, p. 444].

The task of clustering data is a challenging task, first, the data sets are typical large in size, which means that the clustering algorithm must be *scalable*. Additionally, data sets often contains numerous features (*attributes*), which introduces the problem of *curse of dimensionality*, which refers to a several challenges related to high-dimensional data spaces:

First, the issue of *concentration of distances*, where distances between objects in high-dimensional spaces become increasingly similar as dimensionality increases. This means that data points tend to become nearly equidistant from one another, making it difficult for traditional distance-based algorithms to discover clusters.

Secondly, there is the problem of *local feature relevance* and *local feature correlation*, where only a subset of features or different combinations of feature correlations may be relevant for clustering. Consequently, feature reduction techniques like *Principal Component Analysis* (PCA), which project the original space onto a lower-dimensional subspace, are inadequate because they typically identify only one global subspace that best fits the entire dataset. Also, algorithms that evaluate the entire feature space does not address this issue effectively. [?, p. 43–46]

Instead of relying on a global approach to feature selection, a local approach that addresses the issues of local feature relevance and local feature correlation is necessary. However, we then need to deal with two separate problems, which both needs to be solved simultaneously. First, is the problem of finding the relevant subspaces of each cluster. Second, is the problem of finding the clusters in each relevant subspace. To solve them efficiently, heuristics needs to be employed into the clustering algorithms. [?, p. 6–7]

For many applications, it is reasonable to focus only on clusters in axis-parallel subspaces, thus restricting the search space to  $O(2^d)$  dimensions. These algorithms are called *subspace clustering* (or *projected clustering*) algorithms. These can be further divided into: *top-down*- or *bottom-up* approach. In top-down approach, the relevant subspaces for the clusters are determined by gradually reducing the subspaces, starting from the entire space. In contrast, bottom-up approaches, finds the relevant subspaces for the clusters from the original space starting from 1-dimensional using the *monotonicity property* (or *downward closure property*), see Lemma ?? [?]. [?, p. 8, 11]

For the rest of this paper, the following notation will be adopted: Let  $\mathcal{A} = \{A_1, \dots, A_d\}$  represent the set of dimensions (attributes) in a  $d$ -dimensional numerical space,  $\mathcal{S} = A_1 \times A_2 \times \dots \times A_d$ . A subset of  $\mathcal{A}$  is called a subspace. Let  $\mathcal{D} = \{p_1, \dots, p_n\}$  be the dataset, where  $p_i \in \mathcal{S}$  represents a data point. A cluster  $\mathcal{C}$  is a subset of  $\mathcal{D}$ , where each point  $p_i \in \mathcal{C}$  is more similar to the other points in  $\mathcal{C}$  than to points outside of  $\mathcal{C}$ .

**Lemma 1.** *If  $\mathcal{C}$  is a cluster in a  $k$ -dimensional subspace, then  $\mathcal{C}$  must also form a cluster in any  $(k - 1)$ -dimensional projection of the same space.*

A proof is provided in [?].

## 1.1 Contributions

The primary focus of this paper is to analyze the grid-based bottom-up subspace clustering algorithm *MAFIA* [?], which builds upon the first grid-based approach *CLIQUE* [?]. This paper examines the relationship between these two, highlighting their similarities and differences. Additionally, the density-connectivity-based algorithm *SUBCLU* [?] is included, as it offers a contrasting approach to the grid-based approach.

The remainder of the paper is structured as follows: Section 2 describes and analyzes the three algorithms in detail. Section 3 evaluates their performance in

terms of scalability, considering dataset size, data- and cluster-dimensionality, as well as their clustering quality. Section 5 discusses the findings and explores the contributions of each algorithm to the field of subspace clustering. Finally, Section 6 concludes the main findings and suggests some future work.

## 2 Subspace Clustering

The idea of subspace clustering is to identify the subspaces that allows better clustering. Two different approaches of subspace clustering will be discussed in this section including three algorithms: CLIQUE, MAFIA and SUBCLU.

### 2.1 Grid-Based approach

In a grid-based approach,  $\mathcal{S}$  is partitioned into an axis-parallel grid structure. Each grid cell forms a hyper-rectangular *unit*, and the number of points within each unit is counted. This process is initially applied in the 1-dimensional space for each  $A_i$ . Units that contain a number of points exceeding a predefined threshold are retained as *dense units* (DUs). These DUs are then merged with adjacent ones in the next level (2-dimensional space) to form *candidate dense units* (CDUs), those of which are actually dense are kept. The final DUs are then clusters, for which the goal is to describe them using minimal representations in the form of *Disjunctive Normal Form* (DNF) expressions using intervals of the attributes. Two algorithms that follow this approach are CLIQUE and MAFIA.

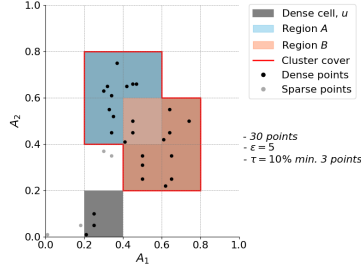
**CLIQUE** CLIQUE, one of the first proposed subspace clustering methods, uses the monotonicity property as its clustering criterion, see Lemma ?? . First,  $\mathcal{S}$  is partitioned into equal-sized *windows* (or *units*) of width  $\varepsilon$  (input parameter). Then, DUs are identified in for the 1-dimensional case, by counting the number of points in each window, for example, using a histogram, as shown in Figure ?? , where each window has a certain *bin* count. For example, if  $\varepsilon = 0.2$ , then 3 DUs are identified in  $A_1$  as they exceed the *density threshold*  $\tau$  (input parameter), while 4 DUs are identified in  $A_2$ . The total of 7 DUs from the 1-dimensional space, will then be used to form CDUs in the 2-dimensional space. Specially, CLIQUE uses the following definition to form CDUs in  $k$ -dimensional space:

**Definition 1.** *CDUs in  $k$ -dimensional space are formed by merging DUs from the  $(k - 1)$ -dimensional space that share the first  $(k - 2)$  attributes.*

Then, as metioned, only the CDUs that are actually dense are kept for further processing.

To optimize performance, CLIQUE applies a pruning technique called *coverage* (from [?]) to reduce computation time, as it may report many CDUs. This technique prunes subspaces with low coverage (i.e., subspaces with low interest are removed). However, this may risk eliminating subspaces that could potentially contain clusters.

Once the final DUs have been identified, CLIQUE seeks to group these DUs into clusters. A *region* is formed when adjacent DUs are connected within a subspace. A *maximal region* is the largest set of connected DUs that cannot be expanded further by adding more adjacent DUs. These maximal regions represent the extent of a cluster in the given subspace. For instance, in Figure ??, two distinct maximal regions, labeled *A* and *B*, are shown. The minimal cover for the clusters has the following DNF expression:  $((0.2 \leq A_1 < 0.6) \wedge (0.4 \leq A_2 < 0.8)) \vee ((0.4 \leq A_1 < 0.8) \wedge (0.2 \leq A_2 < 0.6))$ .



**Fig. 1.** 2-dimensional space containing 8 DUs and depicting two overlapping dense regions *A* and *B*.

**MAFIA** MAFIA can be seen as an extension to CLIQUE and follows mostly the same procedure except from the following:

1. Grid sizes are adaptive and automatically determined based on the data distribution.
2. CDUs are formed using all combinations of the dimensions, that the DUs contains.
3. Does not use the pruning technique as it could result in lost information, as already noted in [?].
4. Tries to parallelize the clustering process.

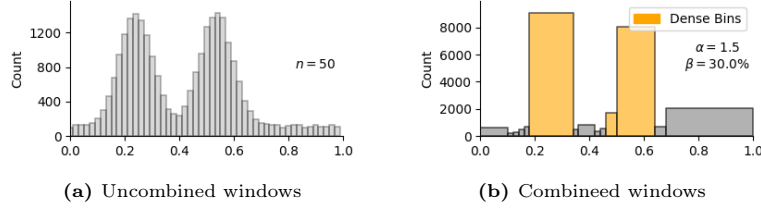
Point 1 and 2 will be explained in more detail next.

*Adaptive grids.* The *Adaptive* algorithm starts by divide each  $A_i$  into high  $n$ -number of equal-sized windows (default  $n = 1000$ ). Then, from left to right, two windows are merged together if their bin count are within a percentage of difference  $\beta$  (input parameter). That means, a high  $\beta$  value result in many merged windows, and vice-versa. If two windows are merged together, the highest bin count are assigned to the window, meaning the bin count that should be used to be compared with the next window. In other words, the algorithm does not use the total bin count of the two for the next comparison. Figure ?? shows an example before and after running the Adaptive algorithm. Having merged the windows, the algorithm finally determines a *dense-level* for each of the windows, that determines how many points a window must contain for a unit to be dense, thus determines whether it should be a candidate for higher dimensional spaces.

The dense-level is given by:

$$\text{dense-level} = \frac{\alpha \cdot a \cdot |\mathcal{D}|}{D_i} \quad (1)$$

where  $a$  is the size of the window,  $\alpha$  is the *cluster dominance factor* (input parameter) and  $D_i$  is the total range  $A_i$  covers. Thus, a higher  $\alpha$  results in a higher dense-level, and vice-versa.



**Fig. 2.** Illustration of the adaptive grids computation.

*CDU Generation.* MAFIA uses the following definition when forming CDUs in  $k$ -dimensional space:

**Definition 2.** *CDUs in  $k$ -dimensional space are formed by merging DUs from the  $(k - 1)$ -dimensional space that share any  $(k - 2)$  attributes.*

Indeed, this results in a high time complexity of  $O(Ndu^2)$ , where  $Ndu$  is the number of DUs. To control this complexity, CLIQUE restricts merging two DUs that share only the first  $(k - 2)$  dimensions. MAFIA, however, allows a broader merging strategy, combining DUs that share any  $(k - 2)$  dimensions. MAFIA can afford this approach because its adaptive grids reduce the number of DUs, effectively controlling the number of comparisons. This flexibility allows MAFIA to achieve a balance between efficiency and more comprehensive cluster detection.

An example illustrating the advantage of merging DUs that share any  $(k - 2)$  dimensions is shown in the following Table ?? . Here, 3 DUs are provided to the a 4-dimensional space. Both CLIQUE and MAFIA would form a 4-dimensional CDU based on DU1 and DU2, as they share the *first* two dimensions. However, only MAFIA would form a CDU based on DU1 and DU3, since it allows merging based on *any* dimensions.

DU1	DU2	DU3
{1, 2, 3}	{1, 2, 4}	{2, 3, 4}

**Table 1.** Illustration of forming a CDU in 4-dimensional space. The number given in each DU represents the dimensions it contains.

## 2.2 Density-Connected approach

A drawback of grid-based methods is that the quality of clustering depends on the positioning of the grids. In Figure ??, we see that the due to the rigid grid

structure we might miss cluster points, that is point not in a dense cell. This problem especially occurs for clusters that are not of rectangular shapes thus, allowing to detect clusters of arbitrary shapes and orientations. One way to address this issue is to a density-connected algorithm, which uses the neighbourhood of points to determine clusters. One such algorithm is SUBCLU.

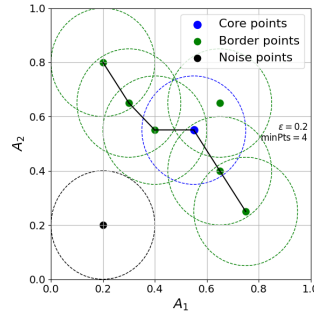
**SUBCLU** SUBCLU extends the principles of the well-known DBSCAN algorithm [?] to higher dimensions by leveraging Lemma ??, which guides the search for clusters across subspaces.

The algorithm begins by applying DBSCAN to each attribute  $A_i$  separately in 1-dimensional subspaces. It identifies dense regions using two input parameters:  $\varepsilon$  which defines the neighborhood radius for points to be considered neighbors, and  $m$ , the minimum number of points required for a point to qualify as a *core point*. If a point is not a core point but lies within  $\varepsilon$  of a core point, it is classified as a *border point*. Points that are neither core nor border points are treated as *noise*.

A *density-connected set* consists of core points that are linked through other core points, and border points that extend the cluster but do not meet the density criteria themselves. These border points are connected to the cluster via core points within their  $\varepsilon$ -neighborhood. A cluster is thus defined as the union of all density-connected sets, including both core and border points.

SUBCLU follows a bottom-up approach: after detecting clusters in 1-dimensional subspaces, it combines  $k$ -dimensional subspaces that share  $(k - 1)$  dimensions to generate candidate  $(k + 1)$ -dimensional subspaces. DBSCAN is reapplied to these candidate subspaces using the same  $\varepsilon$  and  $m$  values to verify whether the clusters found in lower-dimensional subspaces persist as new dimensions are added. Lemma ?? helps to prune the search space efficiently, allowing SUBCLU to prune subspaces where no clusters are detected, thus reducing unnecessary computations.

An example is shown in Figure ??, where SUBCLU identifies a cluster in a 2-dimensional subspace consisting of 6 points. Note that if  $minpts = 2$  instead of 3 it identifies two core points.



**Fig. 3.** Illustration of SUBCLU identifying a cluster in a 2-dimensional subspace.

### 3 Evaluation

The evaluation of the clustering algorithms was performed on a Intel i7 1.70 GHz processor (12th gen.) with 16 GB of RAM running Windows 11.

The evaluation of MAFIA was performed using *GPUMAFIA* [?], which was installed on a virtual machine running Ubuntu, configured with 4 CPUs and 4 GB of RAM.

CLIQUE and SUBCLU were evaluated on the main machine using *ELKI* [?]. Thus, one should be careful to compare the results of the three algorithms directly, as the execution environment and the implementation may affect the results. However, the growth rate and their clustering of the data sets can be compared.

Throughout the different experiments, a range of input parameters for the three algorithms were tested. The best found were selected. The complete evaluation project can be found in the GitHub repository: <https://github.com/henrikdchristensen/SDU-Data-Mining-Exam>. Here, additional experiments as well as detailed descriptions of how to generate the data sets and how to install and use GPUMAFIA is described.

#### 3.1 Data set generation

The aim for the synthetic data generation was to be able to produce similar data sets as discussed in both [?,?]. That is, hyper-rectangles (axis-parallel) non-overlapping clusters in different subspaces. This was achieved by using *MD-CGen* [?], where different sizes and different densities of the clusters could be determined, as well as which attributes for each cluster that are noise for which values at random is selected from a uniform distribution over the entire range of the attribute.

However, to see the one of the main advantages of SUBCLU compared to CLIQUE and MAFIA, a data set containing a Bezier-shaped cluster was created using *Artificial Cluster* [?]. In addition, a self-populated data set containing a plus-shaped cluster was created as discussed in [?].

Finally, two real-world data sets were selected to evaluate the algorithms in a more realistic setting.

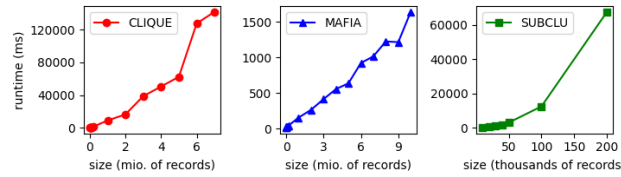
All the data sets were normalized such that each attribute are in the range  $[0, 1]$ .

#### 3.2 Experimental Results

**Scalability with Data Set Size** To evaluate the scalability of the algorithms a 20-dimensional data set containing 5 clusters in 5 different subspaces with 10% noise records was used. The data set size ranges from 10k to 1mio records. The results can be seen in Figure ???. As can be seen, MAFIA is the most scalable algorithm of the three. CLIQUE could handle up to 500k records, while SUBCLU was only able to handle up to 100k records. Similar results in terms of growth rate were reported in all three papers [?, ?, ?]. That is, linear growth for MAFIA and

CLIQUE, and quadratic growth for SUBCLU. The linear behaviour of CLIQUE and MAFIA comes from the fact that the number of passes over the data set depends only on the cluster dimensionality. An increase in the size of the data set just means that more data has to be scanned on every pass over the data set while finding DUs resulting in a linear growth rate. SUBCLU, on the other hand, has a quadratic growth rate, as it relies on the DBSCAN algorithm which needs partial range queries that are costly in terms of running time. Better running time of all three is achieved compared to the three papers, which is probably due to better hardware.

Note that, the `minpts` in SUBCLU were scaled linearly with the data set size, as the clusters in the dataset were of fixed sizes, which means, by increasing the data set size, the density of the clusters increases linearly. However, for the other two algorithms, there was no need to scale any of input parameters, as they relies on the density of the units for the total amount of records.



**Fig. 4.** Scalability with increasing data set size.

**Clustering Accuracy** Many different synthetic data sets were generated to test the clustering accuracy of the algorithms.

The first dataset was a 2-dimensional set containing a single cluster shaped like a skewed plus-sign, with 10% noise added. Results similar to those reported in [?] were achieved. MAFIA successfully identified a single cluster that accurately captured the cluster’s boundaries. In contrast, CLIQUE was unable to report a single cluster that correctly detected the borders. However, MAFIA’s accuracy comes at a cost, as it also identifies many lower-dimensional clusters.

The second data was a 10-dimensional set containing 2 clusters embedded in a different 4 dimensional subspace. 10% of the data was added as noise records. Similar data set as the one used in [?]. The results can be seen in Figure ???. Here, MAFIA was able to detect the clusters without any additional lower-dimensional clusters. In contrast, CLIQUE reports some overlapping clusters and some of the noise records as clusters. SUBCLU detects also the two clusters, but reports many lower-dimensional clusters and noise records as clusters as well.



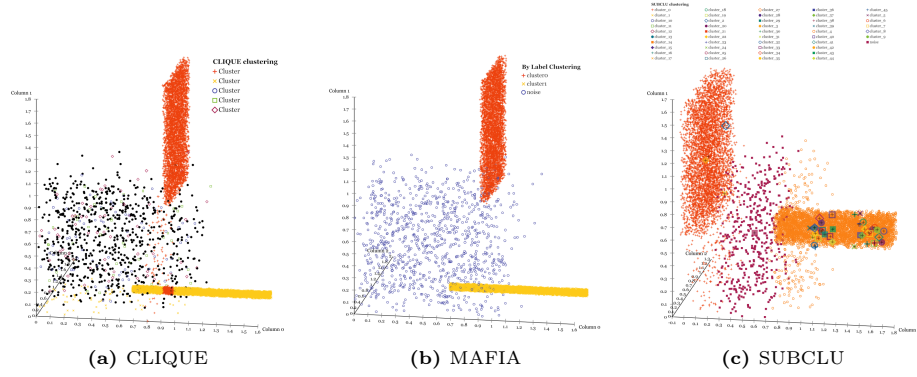


Fig. 5. Two clusters in 4 different subspaces.

The third data set were a 2-dimensional set containing a single cluster formed as a Bezier curve with 10% noise added to the data set. As expected, SUBCLU outperforms the other two algorithms as can be seen in Figure ?? . MAFIA and CLIQUE only partly detects the cluster and reports additional clusters and noise records as clusters. The result of SUBCLU demonstrates its use of DBSCAN.

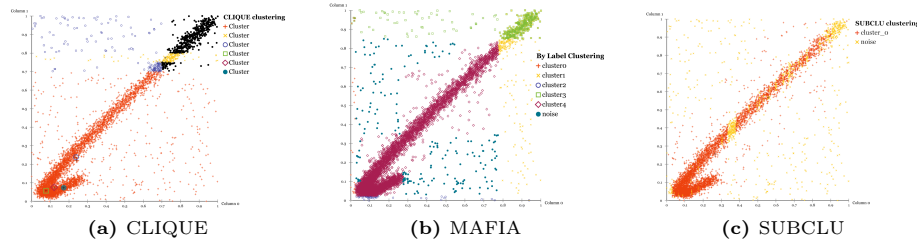
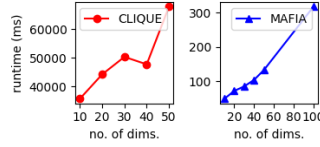


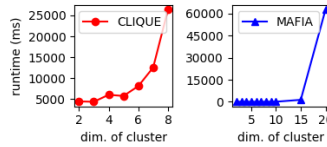
Fig. 6. A single bezier-shaped cluster.

**Scalability with Data Dimensionality and Cluster Dimensionality** Figure ?? shows the scalability of the CLIQUE and MAFIA with increasing data set dimensionality. The data set contains 1 mio. records with 3 clusters in 5 different subspaces and 10% noise records, similar to the data set in [?]. The data set dimensionality ranges from 10 to 100 dimensions. However, CLIQUE was only able to handle half of the dimensions, as the PC constantly freezes – probably due to the high memory consumption. In [?] it is noted that CLIQUE has a quadratic growth rate in terms of data set dimensionality. To investigate this even further, one could try to reduce the memory consumption by having a smaller amount of records in the data set. Nevertheless, MAFIA is the most scalable algorithm of the two.



**Fig. 7.** Scalability with increasing data set dimensionality.

Figure ?? shows the scalability of the CLIQUE and MAFIA with increasing cluster dimensionality. A 20-dimensional data set, containing 500k records with a single cluster were used. The cluster was embedded in increasing number of dimensions starting from 10 to 100 dimensions. 10% of the records was added as noise. The data set is similar to the one used in [?]. Both algorithms heavily suffers from the increasing cluster dimensionality, however, MAFIA is able to handle a higher dimensional cluster than CLIQUE. The reason why both algorithms depends on the cluster dimensionality is due to the fact that a higher cluster dimensionality results in a large subspace coverage and a large number of CDUs. In other words, each pass on the data needs to populate a large number of CDUs, and increase in cluster dimensionality also increases the number of passes over the data set.



**Fig. 8.** Scalability with increasing cluster dimensionality.

The results reported for SUBCLU in [?] were also tried to be replicated, however, similar results was not possible evaluate in different scales of data set dimensionality and cluster dimensionality, as the PC constantly freezes – probably due to the high memory consumption. In future work, it would be interesting to investigate this further.

**Real Data Sets** Two real-world data sets are used to evaluate the algorithms in a more realistic setting. Small data sets were selected for both performance and visualization reasons. Both data sets were normalized such that each attribute was in the range  $[0, 1]$ .

The first data set is the well-known *Iris* data set [?], which contains 150 records of 3 different iris types (e.g. Setosa, Versicolor, Virginica). The data set contains 4 features (sepal length, sepal width, petal length and petal width). All three algorithms produced some meaningful clusters, however, as can be seen in Figure ?. It is hard to determine which algorithm performs the best in this case.

The second data set is the *Date Fruit* data set [?], which contains 898 records of 7 different date fruit types (e.g. Barhee, Deglet Nour, Sukkary, Rotab). These were obtained via a computer vision, where 34 features (e.g shape and color) was

extracted. Only SUBCLU were able to produce clusters somewhat close to the true clusters but also finds many lower-dimensional clusters and small clusters. CLIQUE and MAFIA either produced way too many clusters or close to one single cluster.

## 4 Discussion

CLIQUE, MAFIA, and SUBCLU are three prominent subspace clustering algorithms, each offering unique strengths and limitations. CLIQUE, as one of the earliest subspace clustering approaches. Despite the rigid grid structure it is very simple and intuitive. However, its reliance on fixed grid boundaries can lead to issues when clusters do not align with these grids, resulting in missed or fragmented clusters. Moreover, the algorithm is sensitive to grid size and density thresholds, which may not be optimal for all datasets.

MAFIA builds upon CLIQUE by introducing adaptive grid sizes, allowing it to dynamically adjust grid boundaries based on data density. This modification enables MAFIA to detect clusters more precisely, providing better resolution than CLIQUE, especially in complex datasets. Additionally, experiments indicate that MAFIA is the most scalable algorithm when it comes to handling large data sizes. However, this scalability comes at a cost; MAFIA is highly sensitive to input parameters, such as grid size and density thresholds, and requires careful tuning to achieve optimal results. While it offers more precise boundary detection than CLIQUE, its dependence on parameters can lead to inconsistent performance across different datasets.

SUBCLU, on the other hand, departs from the grid-based approach and utilizes a density-connectivity principle similar to DBSCAN. This enables it to detect clusters of arbitrary shapes, which is a significant advantage when clusters are irregular or not aligned with the axes. This was demonstrated in experiments using real datasets, where SUBCLU outperformed grid-based methods in identifying clusters with complex boundaries. However, SUBCLU’s approach, while flexible, may not be as efficient as MAFIA in terms of data size scalability, particularly for very large datasets. Despite the advantages, one main issue is its scalability in high-dimensional spaces. Since it applies DBSCAN recursively to various subspaces, it incurs a high computational cost due to the large number of range queries required. As the number of dimensions increases, the number of subspaces to be explored grows exponentially, making it computationally expensive for large datasets with many attributes. Also, the input parameters, such as  $\varepsilon$  and  $m$ . These parameters are critical for determining the structure of clusters, but finding the right values can be difficult, especially in high-dimensional datasets where different subspaces may require different parameter settings. If the parameters are not well-tuned, SUBCLU may either miss important clusters or over-partition the data into smaller, less meaningful clusters.

A critical observation across these algorithms is that their performance can vary significantly depending on the data characteristics. Experiments using synthetic data show that clustering quality is highly use-case dependent, and allow-

ing or disallowing lower-dimensional clusters can affect outcomes. Furthermore, while these synthetic experiments highlight each algorithm's strengths, they also reveal their limitations. It is possible to generate datasets where none of these methods can identify the correct clusters, showing that subspace clustering remains a challenging problem.

All three algorithms struggle with closely clustered data points, often requiring fine-tuning of parameters such as grid size in CLIQUE and MAFIA, or  $\epsilon$  and  $m$  in SUBCLU, to effectively separate clusters. This need for parameter optimization underscores the complexity of subspace clustering and the necessity for tailored solutions depending on the specific data distribution and clustering objectives.

"The clustering result can be improved in a variety of ways, from mining significant subspace clusters to using parameter-insensitive clustering approaches. In parameter-insensitive clustering, the 'true' subspace clusters are discovered and they are not manifestations of skewed parameter settings." [?]

"Overcoming parameter-sensitive subspace clustering The current paradigm of subspace clustering requires the user to set parameters, and clusters that satisfy these parameters are returned. We can broadly classify the parameters into cluster parameters and algorithm parameters, where cluster parameters are used in defining the clusters and algorithm parameters are used in guiding the clustering process." [?]

"There are some weaknesses in using intervals. First, as the intervals are non-overlapping, and wrong positioning of the grids may lead to 'truth' subspace clusters being overlooked. Second, setting a fixed size on the intervals using  $\epsilon$  may result in poor clustering quality, as the distribution of the objects in each attribute is different." [?]

"Setting a fixed  $\tau$  may degrade the clustering quality, as a high threshold leads to a small number of subspace clusters, while a low threshold leads to a large number of clusters. To circumvent this problem" [?]

"The cluster is sensitive to the tuning parameters. If wrong  $\tau$  and  $\epsilon$  are set, actual dense units may be overlooked." [?]

"The parameters are difficult to set as they are non-meaningful and non-intuitive. A possible remedy is to try a range of parameter settings, and check for results which are stable in a particular range of parameter settings. Another possible option is to adjust the parameter setting until a suitable number of clusters is obtained." [?]

"Kailing et al. (2004) proposed density based subspace clustering, which overcomes the problems of grid based subspace clustering by dropping the usage of grids. Moreover, it is able to mine arbitrarily shaped subspace clusters in the hyperplane of the dataset." [?]

Both can have overlapping clusters meaning that the same data point can belong to multiple clusters.

Could be interesting to investigate other algorithms like top-down approaches, and other approaches which do not use the monotonicity property.

"It has been observed that a global density threshold, as used by SUBCLU and the grid-based approaches, leads to a bias towards a certain dimensionality: A tighter threshold which is able to separate clusters from noise well in low dimensions tends to lose clusters in higher dimensions, whereas a more relaxed threshold which is able to detect high-dimensional clusters will produce an excessive amount of low-dimensional clusters. Therefore, the dimensionality" [?, p.1:16]

"we find most bottom-up approaches free from the locality assumption. Instead, they pursue a complete enumeration approach facilitated by an APRIORI-like search. Thus, they remain in  $O(2^d)$  in the worst case." [?, p.1:41]

"In order to apply an efficient bottom-up subspace search approach similar to frequent itemset mining, the cluster criterion must implement the downward closure property. Existing bottom-up approaches usually rely on a density-based cluster criterion. A limitation of most of these approaches is that the cluster criterion must use a fixed density threshold for all subspaces in order to implement the downward closure property, although the use of a global density threshold is obviously contraindicated by the observations made before concerning the curse of dimensionality. As a consequence, the same globally defined density threshold applies for subspaces of considerably different dimensionality, although a significant cluster in a higher-dimensional subspace will most likely be less dense (in an absolute sense) than a significant cluster in a lower-dimensional subspace. In order to find higher-dimensional subspaces, the user has to define a less strict density threshold. This, however, would produce a lot of meaningless lower-dimensional clusters. On the other hand, choosing a more strict density threshold, the reported lower-dimensional clusters will probably be more meaningful but higher-dimensional subspace clusters will most likely be lost. The problem of applying a global density threshold obviously also affects the detection of meaningful subspace cluster hierarchies where lower-dimensional clusters are embedded in higher-dimensional ones. If the dimensionality of these" [?, p.1:48]

"The different heuristics and assumptions (refer to Section 5.1) and the different problems tackled (refer to Section 5.2) should always be kept in sight. In most cases, whether the trade-off between efficiency and effectiveness is tolerable will depend on the application." [?, p.1:50]

"The user must use their domain knowledge to help select and tune these settings." [?]

The main parameter for the density threshold can be hard to set across all dimensions. Adaptive grids try to solve this problem by adjusting the grid size based on the density of the data.

A big drawback of using SUBCLU is its use of DBSCAN, which takes logarithmic time for each range query if spatial index is used [?, p. 473]. SUBCLU runs DBSCAN on each  $d$  subspace initially, and then recursively run DBSCAN on each candidate subspace. This can be very time-consuming for high-dimensional data.

"application to synthetic data that has been designed with knowledge about the internal structure" [?]

Internal vs. External Evaluation - "internal" evaluation ( unsupervised) - how well does the result fit to the data, assuming certain properties of "good" clusters - cohesion/separation (e.g., TD2, Silhouette) - similarity matrix (correlation, visualization) - basic assumption: the clustering algorithm is suitable for the given problem - "external" evaluation ( supervised) - validation of the result independent of the algorithm, given some ground truth e.g., how well could the clustering algorithm identify known classes of objects - assumption: application to new data from the same domain and a similar problem set-up might perform similarly well - problem: the algorithm is "punished" if it actually detects "novel" patterns

Internal Evaluation - fundamental problem: is the algorithm suitable to tackle the given task? - decision (needs to precede the application) w.r.t. basic properties of the algorithm and expected characteristics of the data and clusters

Cohesion and Separation - cohesion: how strong are the cluster objects connected (how similar, pairwise, to each other)? - separation: how well is a cluster separated from other clusters? - validation index: suitable combination of cohesion and separation

Relationship Between Validity Index and Type of Clustering - cohesion and separation are suitable for convex cluster models, not for arbitrarily shaped clusters as possibly found, e.g., by density-based approaches

Problems of Internal Evaluation - suitability of the clustering model for the given dataset? - determinism? - determine k? - comparison of different cluster models? - connection between objective function of the clustering model and the validation index?

External Evaluation - based on a mapping between clustering-result and given clusters (classes), i.e., some ground truth or gold standard why not simply confusion matrix? - two basic approaches: - mapping of sets of objects - comparison of assignments for pairs of objects ("pair counting") - assessment of agreement between given and found partitions - mapping of sets: mainly information theoretic measures - pair counting: many measures available, e.g. also measures known from classification (such as F-measure etc.)

The advantage of entropy-based measures is that they do not require a one-to-one matching of sets (clusters/classes).

Problems of External Evaluation - It might establish a wrong bias on the development of clustering algorithms if they are designed with a bias to discover known classes again [Färber et al., 2010]. - Class structure and cluster structure in a dataset need not be the same [Färber et al., 2010]. - Data can contain different "truths", so different clustering results might be equally meaningful and interesting [Zimek and Vreeken, 2015, Bailey, 2013].

## 5 Conclusion

The analysis of CLIQUE, MAFIA, and SUBCLU highlights the strengths and limitations of each algorithm in subspace clustering. CLIQUE offers a straightforward, grid-based approach suitable for detecting axis-aligned clusters, but its reliance on fixed grids and sensitivity to parameters can limit performance. MAFIA improves upon CLIQUE by adapting grid sizes, providing better resolution and scalability for large datasets, though it still remains somewhat dependent on parameter tuning. SUBCLU, using a density-based method, excels in identifying arbitrarily shaped clusters, especially in real-world data, but it may not scale as efficiently as CLIQUE and MAFIA.

Overall, the choice of algorithm depends on the use case, data size, and the shape of clusters. The experiments demonstrated that these algorithms perform well under controlled conditions, but in highly complex data sets, none of the algorithms detects clusters accurately. Hence, while each method has its advantages, they all require careful adjustment of parameters and may need further refinement for optimal performance in diverse datasets.

## References

1. Adinetz, A., Kraus, J., Meinke, J., Pleiter, D.: Gpumafia: Efficient subspace clustering with mafia on gpus pp. 838–849 (08 2013). [https://doi.org/10.1007/978-3-642-40047-6\\_83](https://doi.org/10.1007/978-3-642-40047-6_83)
2. Agrawal, R., Gehrke, J., Gunopulos, D., Raghavan, P.: Automatic subspace clustering of high dimensional data for data mining applications p. 94–105 (1998). <https://doi.org/10.1145/276304.276314>, <https://doi.org/10.1145/276304.276314>
3. Ester, M., Kriegel, H.P., Sander, J., Xu, X.: A density-based algorithm for discovering clusters in large spatial databases with noise. In: Proceedings of the Second International Conference on Knowledge Discovery and Data Mining. p. 226–231. KDD’96, AAAI Press (1996)
4. Fisher, R.A.: Iris. UCI Machine Learning Repository (1936). <https://doi.org/https://doi.org/10.24432/C56C76>
5. Han, J., Kamber, M., Pei, J.: Data Mining: Concepts and techniques. Elsevier (6 2011)
6. Iglesias Vázquez, F., Zseby, T., Ferreira, D., Zimek, A.: Mdcgen: Multidimensional dataset generator for clustering. *Journal of Classification* **36** (04 2019). <https://doi.org/10.1007/s00357-019-9312-3>
7. Koklu, M., Kursun, R., Taspinar, Y.S., Cinar, I.: Classification of date fruits into genetic varieties using image analysis. *Mathematical Problems in Engineering* **2021**(1), 4793293 (2021). <https://doi.org/https://doi.org/10.1155/2021/4793293>, <https://onlinelibrary.wiley.com/doi/abs/10.1155/2021/4793293>
8. Kriegel, H.P., Kröger, P., Zimek, A.: Clustering high-dimensional data: A survey on subspace clustering, pattern-based clustering, and correlation clustering. *ACM Trans. Knowl. Discov. Data* **3**(1) (Mar 2009). <https://doi.org/10.1145/1497577.1497578>, <https://doi.org/10.1145/1497577.1497578>
9. Kröger, P., Kriegel, H.P., Kailing, K.: Density-connected subspace clustering for high-dimensional data **246-257** (04 2004). <https://doi.org/10.1137/1.9781611972740.23>

10. Li, W., Zhou, Z.: Ac: A data generator for evaluation of clustering (02 2022). <https://doi.org/10.36227/techrxiv.19091330>
11. Nagesh, H., Goil, S., Choud, A., Choudhary, P.: Adaptive grids for clustering massive data sets (01 2002). <https://doi.org/10.1137/1.9781611972719.7>
12. Parsons, L., Haque, E., Liu, H.: Subspace clustering for high dimensional data: A review. *SIGKDD Explor. Newsl.* **6**(1), 90–105 (Jun 2004). <https://doi.org/10.1145/1007730.1007731>, <https://doi.org/10.1145/1007730.1007731>
13. Schubert, E.: Automatic indexing for similarity search in ELKI **13590**, 205–213 (2022). [https://doi.org/10.1007/978-3-031-17849-8\\_16](https://doi.org/10.1007/978-3-031-17849-8_16), [https://doi.org/10.1007/978-3-031-17849-8\\_16](https://doi.org/10.1007/978-3-031-17849-8_16)
14. Sim, K., Gopalkrishnan, V., Zimek, A., Cong, G.: A survey on enhanced subspace clustering. *Data Mining and Knowledge Discovery* **26**(2), 332–397 (2 2012). <https://doi.org/10.1007/s10618-012-0258-x>, <https://doi.org/10.1007/s10618-012-0258-x>