

Efficient Cluster Detection by Ordered Neighborhoods

Emin Aksehirli¹(✉), Bart Goethals¹, and Emmanuel Müller^{1,2}

¹ University of Antwerp, Antwerp, Belgium

{emin.aksehirli,bart.goethals}@uantwerpen.be

² Karlsruhe Institute of Technology, Karlsruhe, Germany

emmanuel.mueller@kit.edu

Abstract. Detecting cluster structures seems to be a simple task, i.e. separating similar from dissimilar objects. However, given today’s complex data, (dis-)similarity measures and traditional clustering algorithms are not reliable in separating clusters from each other. For example, when too many dimensions are considered simultaneously, objects become unique and (dis-)similarity does not provide meaningful information to detect clusters anymore. While the (dis-)similarity measures might be meaningful for individual dimensions, algorithms fail to combine this information for cluster detection. In particular, it is the severe issue of a combinatorial search space that results in inefficient algorithms.

In this paper we propose a cluster detection method based on the *ordered neighborhoods*. By considering such ordered neighborhoods in each dimension individually, we derive properties that allow us to detect clustered objects in dimensions in linear time. Our algorithm exploits the ordered neighborhoods in order to find both the similar objects and the dimensions in which these objects show high similarity. Evaluation results show that our method is scalable with both database size and dimensionality and enhances cluster detection w.r.t. state-of-the-art clustering techniques.

1 Introduction

In the information era that we live in, there are a huge amount of data about almost anything. Institutions, both government and private, realized the importance of data and started to collect any kind of information on their business objects, hoping that they will derive useful knowledge out of it one day. Considering the amount, annotating and labeling the data is often infeasible. Therefore, unsupervised methods such as clustering are more suitable for knowledge extraction.

One of the challenging effects of this “data hoarding” is the increased number of attributes associated with each object. Unfortunately, due to the phenomenon of the so called *curse of dimensionality*, the similarity between objects becomes meaningless in high dimensional data spaces [5]. This means that the clusters cannot be separated from each other by (dis-)similarity assessment in high dimensional space, which in turn poses a serious challenge for traditional

clustering algorithms such as k -means or hierarchical clustering. Nevertheless, (dis-)similarity in the individual dimensions can still be exploited and clusters can be detected in combinations of these dimensions. The open challenge is how to exploit this (dis-)similarity information in individual dimensions for any combination of dimensions while not falling prey to the exponential nature of this combinatorial search space.

For example, nowadays a customer can be associated with credit ratings, shopping habits, travel patterns, entertainment choices, sport habits, etc. Even medical conditions or other private information might be available due to the data integration with a multitude of data sources. Considering all the aforementioned data, each customer becomes very unique and almost equally dissimilar to any other customer. This makes the notion of “similar customers” meaningless. Nevertheless, a meaningful customer segmentation can still be achieved by looking at a subset of dimensions, e.g. travel and sport habits only.

Although similarity between high dimensional objects is not meaningful, similarity according to subsets of attributes is still meaningful. To address this issue, subspace clustering [20] aims at cluster detection in any combination of the given attributes. Even though they enhance clustering quality compared to traditional clustering methods, they come with their own challenges. They suffer from noise sensitivity [1, 19], density estimation challenges [4, 9, 10], many complex parameters [4, 10, 12, 13], or inefficient search through the combinatorial search space [2, 9].

In contrast to these methods, we tackle the problem by exploiting local neighborhoods of objects in individual dimensions. These neighborhoods are used as means for cluster detection in combinations of dimensions. We show that clusters can be detected directly from these local neighborhoods. Moreover, exploiting their natural orders, we can avoid common scalability pitfalls in the algorithmic computation of clusters. In particular, we propose a linear-time algorithm for detecting cluster structures in dimensions. Key characteristics of our approach are its scalability w.r.t. both database size and dimensionality, the detection of clusters with variable scales, and its few user-friendly parameters. Further, it allows for individual (dis-)similarity measures for each dimension, which is important for data with different data types, complex distance functions, and the lack of joint (dis-)similarity measures in combinations of different data types and dimensions with complex distance functions.

In summary, our contributions are as follows: - A formal analysis of ordered neighborhoods in the context of cluster detection. - Prove completeness of cluster detection based on ordered neighborhoods. - A scalable algorithm exploiting ordered neighborhoods for cluster detection in data projections. - Exhaustive experiments showing that our method (1) works well with variable densities, (2) scales well with database size and dimensionality, (3) robust w.r.t. noise and irrelevant dimensions, (4) and suitable for exploratory data analysis in real world scenarios.

Please find the supplementary for the paper, in which we cover the topics in detail, along with the implementation and datasets on our website.¹

¹ <http://adrem.uantwerpen.be/clon>.

2 Formal Properties

We use the ordered local neighborhoods of objects in each individual dimension as indicators for clustered structures. This idea is based on the observation that similarity of objects is captured by the order of objects contained in the local neighborhood [17]. To the best of our knowledge, we are the first ones to exploit this theoretic principle for scalable cluster detection in high dimensional data.

For our solution, objects only need to be sorted once according to the given (dis-)similarity measure. This might even be given (for free) due to the index structures maintained in the relational database system that stores the data. Based on these ordered neighborhoods, we detect clusters by mining object sets that re-occur in a similar order in multiple dimensions. Starting with one dimensional clusters and iteratively refining them allows us to detect clusters in projections (i.e. dimensions which agree in the similarity of objects) of a high dimensional data space.

Although the general idea seems simple to implement, it has several open research questions: Can we guarantee to detect all hidden clusters? How to capture the neighborhood information, and how can it be exploited for clustering high dimensional data? In order to answer all of these questions one by one, we structure our main contributions as follows: We first show that ordered neighborhoods ensure complete cluster detection in Sect. 2.1. Second, we propose a data transformation from relational data to an ordered neighborhood space and investigate its properties in Sect. 2.2. We introduce a scalable algorithm that exploits these properties for cluster detection in Sect. 3.

2.1 Clustering in Neighborhoods

For a *relational database* \mathcal{D} we represent each data object \mathbf{o} in \mathcal{D} as a vector defined over attributes $\mathcal{A} = \{A_1, A_2, \dots, A_m\}$, with \mathbf{o}_i the value of the object for the attribute A_i . For each attribute A_i , the projected database is denoted as \mathcal{D}^{A_i} . Further, we use $\delta(\mathbf{o}, \mathbf{p})$ as a dissimilarity between the objects \mathbf{o} and \mathbf{p} . $|S|$ denotes the cardinality of set S .

As basic notion we use local neighborhoods based on k -nearest neighborhood (k -NN), which have already been exploited for lazy classification [7], dimensionality reduction [21], collaborative filtering [16] and many other techniques in various communities [6, 11]. For clustering, one has used shared nearest neighborhoods [8], neighborhoods in random projections [18], and frequent co-occurrence of neighborhoods [2]. In our work, we build the basic clustering principle of our method on the k -nearest neighborhood of an object, which is the set of objects that are the most similar to it:

Definition 1 (*k -Nearest Neighborhood (k -NN)*). *Let $\mathbf{o} \in \mathcal{D}$, δ a dissimilarity measure, and the $NN_k(\mathbf{o})$ be the k^{th} nearest object to \mathbf{o} according to δ , the k -nearest neighborhood (k -NN) of \mathbf{o} is defined as*

$$k\text{-NN}(\mathbf{o}) = \{\mathbf{p} \in \mathcal{D} \mid \delta(\mathbf{o}, \mathbf{p}) \leq \delta(\mathbf{o}, NN_k(\mathbf{o}))\}.$$

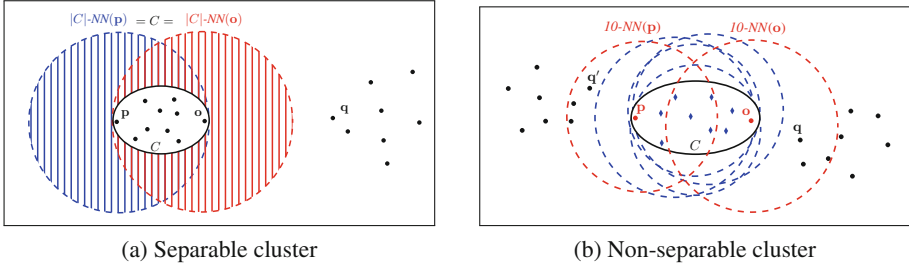


Fig. 1. Nearest neighborhoods of the objects in a cluster

k - $NN(\mathbf{o})$ captures the similar objects near object \mathbf{o} , which we can use for cluster detection. Figure 1 shows the nearest neighborhoods of the objects in clusters. Cluster sized nearest neighborhoods of objects in a *separable cluster* are equal to the cluster itself. In Fig. 1a, $|C|$ - NN of all 10 objects in C are equal to the C itself. Although the k - NN of some of the objects in a non-separable cluster include objects that are not in the cluster, the k - NN of the majority of objects in the cluster include the whole cluster. For example in Fig. 1b, the neighborhoods of the data points \mathbf{p} and \mathbf{q} are different from each other although they are in the same cluster. On the other hand, they are in the $|C|$ - NN of the 8 out of 10 objects in the cluster, shown as \blacklozenge . Therefore, they can be identified as being in the same cluster by checking their co-occurrences in neighborhoods of other objects instead of comparing their neighborhoods. This is one of the advantages of using co-occurrences in neighborhoods instead of shared nearest neighborhoods.

2.2 Ordered Neighborhoods

In the original space, computing the co-occurrence of object sets in neighborhoods of objects in \mathcal{D} requires expensive neighborhood computations. Instead, we compute the neighborhood of every object once and conduct co-occurrence search in these neighborhoods.

Definition 2 (Ordered Neighborhood). *The ordered neighborhood of $\mathbf{o} \in \mathcal{D}^{A_i}$ is the ordered list of the objects in k - $NN(\mathbf{o})$, in which the order reflects the order of objects in A_i .*

$$k\text{-}NN(\mathbf{o}) = (\dots, \mathbf{p}, \dots, \mathbf{q}, \dots) \iff \mathbf{p} < \mathbf{q}$$

Definition 3 (Ordered Neighborhood Database). *An ordered neighborhood database \mathcal{N} is the ordered list of ordered neighborhoods. Order of the neighborhoods is the order of objects in \mathcal{D}^{A_i} . If $\mathbf{o} < \mathbf{p}$, then*

$$\mathcal{N}^{A_i} = (\dots, k\text{-}NN(\mathbf{o}), \dots, k\text{-}NN(\mathbf{p}), \dots)$$

For the remaining of the paper, we use *neighborhood* and k - $NN(\mathbf{o})$ to refer to the ordered neighborhood, and *neighborhood database* to refer to the ordered

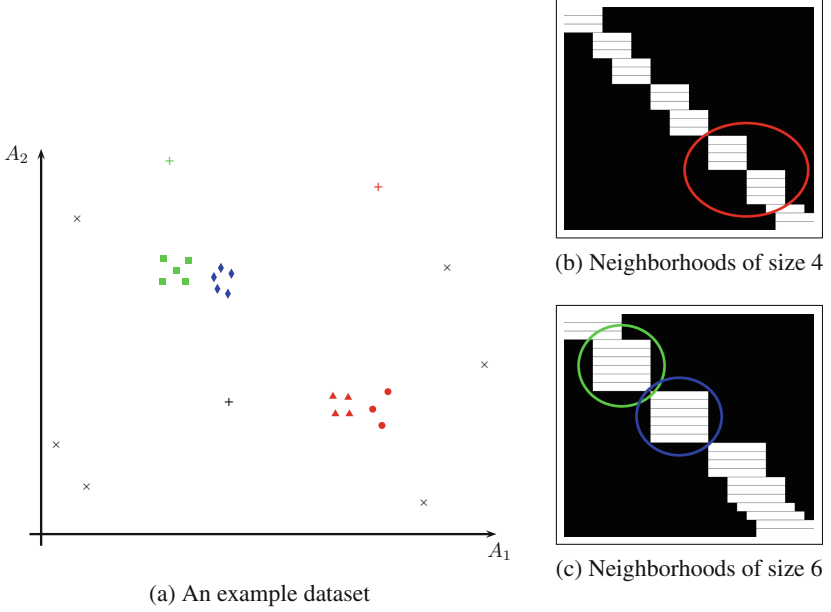


Fig. 2. An example dataset and its neighborhood databases for different sizes

neighborhood database unless it is noted otherwise. Further, if the projection attribute is clear from the context or it is irrelevant, we drop the attribute and use \mathcal{N} instead.

Let us look at the example dataset in Fig. 2a. Figures 2b and 2c show \mathcal{N}^{A_1} of the example dataset with $k = 4$ and $k = 6$. Each column in a neighborhood database represents an object while each row represents the neighborhood of an object. If an object occurs in a neighborhood, the corresponding cell is white. E.g., in Fig. 2b each horizontal white line, which denotes 4- NN , consists of exactly 4 cells.

With a formal analysis, the following properties hold for the neighborhood DBs, their proofs are provided in Sect. S.1 (in Supplementary):

Property 1 (Consecutive Objects). If two objects are in a neighborhood, all the objects in between them are also in that neighborhood. Let $\mathbf{o}, \mathbf{p}, \mathbf{q}, \mathbf{r} \in \mathcal{D}^{A_i}$ and $\mathbf{o} < \mathbf{p} < \mathbf{q}$. $\mathbf{o} \in k\text{-}NN(\mathbf{r}) \wedge \mathbf{q} \in k\text{-}NN(\mathbf{r}) \implies \mathbf{p} \in k\text{-}NN(\mathbf{r})$.

Property 2 (Consecutive Neighborhoods). If an object is in the neighborhood of two objects, then it is in the neighborhood of all the objects in-between them. Let $\mathbf{o}, \mathbf{p}, \mathbf{q}, \mathbf{r} \in \mathcal{D}^{A_i}$ and $\mathbf{o} < \mathbf{p} < \mathbf{q}$. If $\mathbf{r} \in k\text{-}NN(\mathbf{o}) \wedge \mathbf{r} \in k\text{-}NN(\mathbf{q}) \implies \mathbf{r} \in k\text{-}NN(\mathbf{p})$.

We can see Properties 1 and 2 in Fig. 2b. While the objects in a neighborhood are consecutive, an object only appears in consecutive neighborhoods, hence, they respectively form straight horizontal and vertical lines.

Theorem 1 (Recurrent Neighborhoods). *Clusters form square structures on the diagonal of the neighborhood databases.*

Proof. According to Properties 1 and 2, neighborhoods form continuous shapes. As we discuss in Sect. 2.1, a cluster is repeated as the neighborhoods of its objects, hence, the shape should be a square. Since each object is its own nearest neighbor, squares should be on the diagonal of the neighborhood DB. \square

The two squares in the lower right quarter of the Fig. 2b, which are circled in red, are the neighborhoods of the objects in the clusters that are denoted by \blacktriangle and \bullet in Fig. 2a. Similarly, squares on the upper left quarter of Fig. 2c, which are circled in green and blue, are respectively from the clusters \blacksquare , \blacklozenge .

Property 3 (Baseline). If there are no clusters in a segment of the dataset, i.e., the objects are distributed uniformly, all of the object sets for that segment have the same amount of recurrency in the neighborhood DB, which is a function of k , cf. Figure S.1.

Property 4 (Transitivity). If an object \mathbf{o} is not in the neighborhood of $\mathbf{p} > \mathbf{o}$, then it is not in the neighborhood of any $\mathbf{q} > \mathbf{p}$. Formally let $\mathbf{o}, \mathbf{p}, \mathbf{q} \in \mathcal{D}^{A_i}$ and $\mathbf{o} < \mathbf{p} < \mathbf{q}$, $\mathbf{o} \notin k\text{-NN}(\mathbf{p}) \implies \mathbf{o} \notin k\text{-NN}(\mathbf{q})$.

In Fig. 2b, we can see that the 6th object occurs in neighborhoods 4 to 9. It does not occur in any neighborhood before 4 and not in any neighborhood after 9. We use the transitivity to limit the search space of our Algorithm, cf. Sect. 3, while the baseline property helps us to discard the segments without a cluster structure and to eliminate the artifacts that are caused by the transformation.

3 Mining the Neighborhoods

In this section we introduce our proposed algorithm CLON, which efficiently finds all subspace clusters by using ordered neighborhoods. A detailed explanation and pseudo code of the algorithm can be found in Section S.3 (in the supplementary).

If a cluster exists only in a subset of dimensions, the irrelevant dimensions hinder the search and the quality of the cluster. As we mentioned before, this is often the case in high dimensional spaces. To eliminate the negative effects of irrelevant dimensions, CLON follows a bottom-up strategy and exploits the fact that clusters in lower dimensional projections are supersets of clusters in higher dimensional spaces [9]. Hence, CLON finds the one dimensional clusters first and then iteratively refines them by evaluating their higher dimensional subsets.

In its first phase, CLON converts the relational database into a neighborhood database. Values in each dimension have to be sorted only once, after which the neighborhoods can be computed very fast by using a sliding window. A neighborhood DB is created for each of the dimensions. CLON already benefits from the ordered neighborhoods while creating the neighborhood DBs: It exploits the consecutiveness properties to store the neighborhood information so

that instead of storing the whole neighborhoods, only the start and end of them are stored. This provides a dramatic memory saving, cf. Algorithm S.1.

In the second phase, the clusters in individual projections are found. As Theorem 1 states, clusters form recurrent object sets in the neighborhood DB. It has been proposed to find these recurrent object sets by using frequent itemset mining methods, however, since these methods are computationally expensive, only approximate methods are feasible for non-trivial datasets [2]. Our method CLON exploits the properties of ordered neighborhoods, so that it can mine these recurrent object sets in linear time instead of exponential time (both in the number of objects). Thanks to this speed improvement, *all* maximally large objects sets that are more recurrent than a certain threshold can be found in a reasonable time, cf. Sect. 4.3. Therefore in this phase, CLON can identify all of the recurrent objects sets in all individual projections, that is all one dimensional clusters, cf. Algorithm S.2.

In the third and the last phase, one dimensional clusters are used to find higher dimensional clusters. For any higher dimensional cluster, there are clusters that have more or equal number of objects, in the subsets of the dimensions of the original cluster [2, 9, 10]. For example, if a cluster exists in dimensions 1, 2, and 3, then a superset of its objects form clusters in dimension pairs (1,2), (2,3) and (1,3). Therefore, we can find higher dimensional clusters by refining lower dimensional ones, i.e. removing the objects that are not in the cluster. CLON refines each one dimensional cluster by checking whether any of its subsets are recurrent in other dimensions. If there are any, CLON continues with a depth-first search of higher dimensional clusters by traversing all the possible dimensions until none of the subsets are large or recurrent enough. Since CLON uses neighborhood DBs also for higher dimensional search, properties of ordered neighborhoods are exploited during this phase too. Lastly, the recurrent object sets are output as clusters along with the dimensions that they are recurrent in, cf. Algorithm S.3.

CLON requires two user-friendly parameters: (1) Minimum size of a cluster and (2) the neighborhood size. Selecting the minimum size should be straightforward, it depends on the size of the cluster that the user would like to find. The neighborhood size should be larger than the cluster size. Because of the recurrency oriented mining, CLON is robust to these parameters. A detailed discussion about parameter selection can be found in Section S.2.

3.1 Complexity Analysis

For a numeric \mathcal{D} that has n objects and d attributes, the computational complexity of transforming it to the neighborhood DB is $O(d \times n \times \log(n) \times k)$. Thanks to the consecutiveness properties, we do not have to store the entire database but only the start and end of the neighborhoods are enough. Therefore, the space requirement for the neighborhood DBs is $O(n \times d)$ which is equivalent to the storage of the original \mathcal{D} . Exploiting the properties of the neighborhood DB allows us to find one dimensional clusters extremely efficiently, in $O(n)$. Therefore, the complexity of finding all one dimensional clusters is $O(d \times n)$. Total time to find

all subspace clusters depends on the dimensionality of the clusters in the data. Although the computational complexity to search a cluster in a dimension is linear, the number of dimensions to search is combinatorial. Empirical results in Sect. 4.3 show that the scalability and efficiency of CLON is not hindered by the dimensionality of clusters.

4 Empirical Evaluation

4.1 Experimental Setup

We evaluate our method on heterogeneous datasets, i.e., datasets with different scales, both in terms of dimension scale and cluster distribution. We also conduct a set of experiments to evaluate the scalability of our method according to the number of objects, number of dimensions and the amount of noise. To see whether our method is a good fit for real world scenarios, we conduct experiments on very high dimensional real world datasets. We compare the quality of the clusters that are found by our method with the state of the art subspace clustering methods, such as CartiClus [2], FIRES [10], PROCLUS [1], STATPC [12], and SUBCLU [9], cf. Section S.4.

Object cluster finding capabilities of the methods are evaluated by the supervised F1 Measure. Where appropriate, we evaluate the subspace discovery capabilities by using the established E4SC score [14]. Runtime results are given for the performance versus scale experiments. More information about the experimental setup can be found in Section S.5. For the reproducibility of the experiments, a cross-platform implementation of the method and the information about the datasets are available on our website.

4.2 Heterogeneous Datasets

We measure the cluster discovery capabilities of our method on datasets that (1) have different scaling in different dimensions and (2) have clusters with different spreads. To evaluate these capabilities we generate two sets of datasets: One with different scales of dimensions and another one with clusters of different scales. Details about these datasets can be found in Section S.5.

Figures 3a and 4a show the F1 scores of the methods on a selected subset of these datasets. Different scale factors are indicated with different colors. On both set of experiments, CLON produces very stable and high quality clusters. CartiClus is comparable with CLON, which shows the effectiveness of nearest neighborhood based clustering. Radius based neighborhood evaluation of FIRES and distance sensitive projections of PROCLUS cannot cope with scaling. SUBCLU and STATPC produce high quality, although unstable, results thanks to their thorough search. Considering the execution times of CLON and its closest competitors, cf. Fig. 5, we can see the advantages of cluster refining through dimension search in ordered neighborhoods: CLON consistently produces better or comparable clusters in an order of magnitude less time.

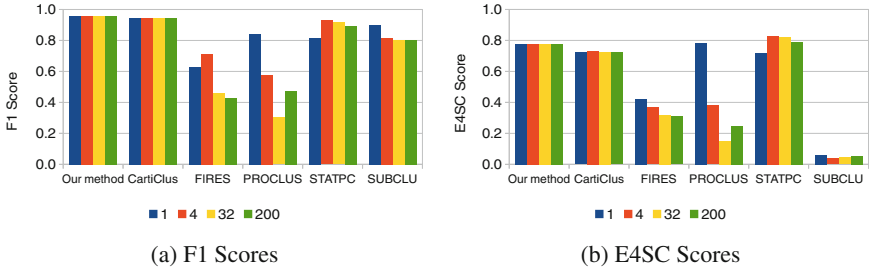


Fig. 3. Quality scores of the methods on datasets with dimension of various scales

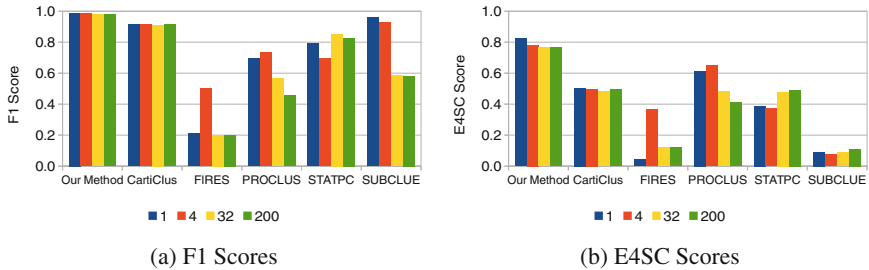


Fig. 4. Quality scores of methods on datasets that have clusters of various scales

Figures 3b and 4b show the E4SC scores for the same sets of datasets. While some of the methods perform poorly on dimension detection, dimension finding capabilities of the most of them are in parallel with the object cluster finding capabilities. CLON performs better than CartiClus in dimension detection although they both mine recurrent object sets, cf. Fig. 4b. This is because CLON is fast enough to mine all subspace clusters instead of a sample of them.

4.3 Scalability Results

We conduct experiments to understand the scalability of our algorithm according to the data size, dimensionality, and the noise ratio. For these experiments we use the datasets that are used in the literature for similar comparative studies [2, 14].

Figure 5a shows the run times of the methods on datasets with an increasing number of objects. This set of datasets include 5 different datasets which have approximately 1500, 2500, 3500, 4500, and 5500 objects and 20 dimensions. Our method scales well compared to the algorithms that search the combinations of dimensions.

Figure 5b shows the run times for different methods on datasets with an increasing number of dimensions. We conduct experiments on 6 different datasets that have approximately 1500 objects in 5, 10, 15, 25, 50, and 75 dimensions. In all of these datasets, each cluster exists in approximately 80% of the dimensions. Therefore, these experiments also evaluate the performance of the algorithms

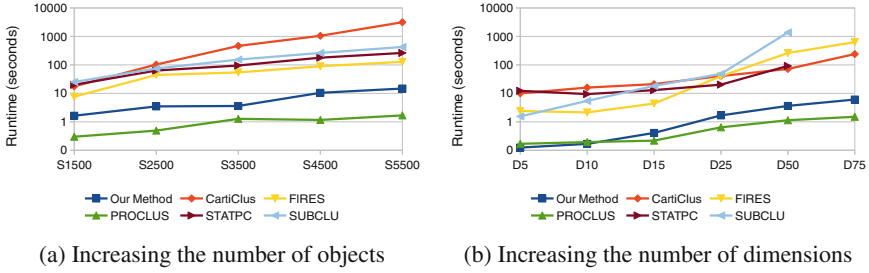


Fig. 5. Execution times

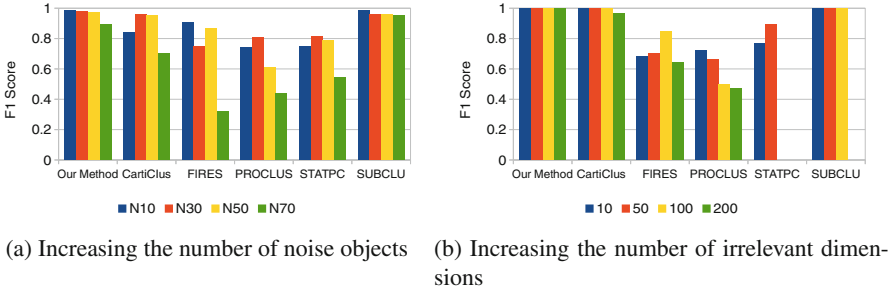


Fig. 6. Noise detection capability experiments

regarding dimensionality of clusters. Our method scales well with the increasing number of dimensions while utilizing the information in combinations of dimensions. Note that the missing values indicates that the method did not complete in a practical time or failed because of excessive memory requirements.

Figure 6a shows the quality results for the datasets that contain 10 %, 30 %, 50 %, and 70 % noise. Our method and CartiClus can effectively discard noise thanks to their recurrency and neighborhood based foundation. Only SUBCLU is on par with these methods because of its thorough search which also makes it slower around 100 times than our method, cf. Fig. 5.

Figure 6b shows the capability of irrelevant dimension detection. A dataset which has 10 clusters hidden in subsets of 10 dimensions is generated. Then, 10, 50, 100 and 200 uniformly randomly generated dimensions are added to the dataset. Here again, we see that our method, CartiClus and SUBCLU can effectively discard irrelevant dimensions even though the meaningful structures are hidden in noise that is 20 times of its size. Some of the methods did not complete in meaningful time for some of the datasets, including SUBCLU for more than 100 dimensions.

4.4 Real World Datasets

To evaluate our method, we used two gene expression datasets, *Nutt* and *Alon*, along with a movie rating dataset, *movies*. *Alon* is a dataset of 2000 gene expression across 62 tissues from a colon cancer research. The tissues are grouped

into 2 categories: 40 healthy tissues and 22 tissues with tumor [3]. *Nutt* dataset contains expressions of 1377 genes on 50 glioma tissue samples that are grouped into 4 different pathology categories [15]. High dimensional nature of both datasets make the clustering challenging even for the subspace clustering methods. Table 1 shows the F1 scores of the methods. “n/a” indicates that the method did not complete in a practical time or failed because of excessive memory requirements. These results show that, although our method searches the clusters in combinations of dimensions, it keeps being scalable even for very high dimensional datasets.

Table 1. F1 scores for gene expression datasets

	<i>Alon</i>	<i>Nutt</i>
Our method	0.78	0.75
PROCLUS	0.46	0.44
FIRES	0.52	0.55
SUBCLU	0.58	n/a
STATPC	n/a	n/a
CartiClus	n/a	n/a

We conduct an exploratory data analysis on a movie ratings dataset from GroupLens.² We use 10M movieLens dataset which includes 10000054 ratings applied to 10681 movies by 71567 users from the website <http://movielens.org>. We use movies as objects and users as attributes. We start with finding the most similar 3 movies according to the user ratings. The result is, not surprisingly, the original Star Wars trilogy. When we lower the similarity threshold, we start to see a cluster of the Lord of the Rings Trilogy along with some other high profile movies in clusters of science fiction movies, action movies and crime movies. Some of the selected clusters are given in Table S.4a.

We continue our analysis by increasing the number of similar movies to 6. Considering the similarity of the original Star Wars trilogy, we search for a cluster of all 6 of the released Star Wars movies, with no luck. Actually, lack of this 6 pack of Star Wars movies cluster does not contradict with the general opinion of the fans of the franchise because the last 3 movies are not as popular as the originals. Table S.4b shows some of the interesting clusters of size 6, which includes a cluster of two most popular trilogies, a cluster of distopian movies, a cluster of popular thrillers, and a cluster of very popular classic movies.

5 Conclusion

In this paper, we tackle the problem of finding clusters in high dimensional databases. We make a formal analysis of ordered neighborhoods and their intrinsic

² <http://grouplens.org/>.

properties. We show that co-occurrences in local neighborhoods of objects are indicators of cluster formations, and hence, clusters can be found by mining recurrent neighborhoods. We propose a scalable algorithm that exploits these intrinsic properties to find all of the clusters and their relevant dimensions. Along with its scalability, key properties of our algorithm include its intuitive and user friendly parameters, its noise detection capabilities, its adaptivity to datasets with various scales, and its capability to exploit multiple (dis-)similarity measures.

Besides conducting experiments on scalability, noise detection, and adaptivity; we tested our method on some very high dimensional gene expression datasets. Further, we did an exploratory analysis on a movie rating dataset to show that our method is a good fit for real world scenarios. Finally, we supply a software tool that can be used to interpret the data for further analysis.

Acknowledgements. Emmanuel Müller is supported by Post-Doctoral Fellowships of the Research Foundation – Flanders (FWO). Further, this work is supported by the Young Investigator Group program of KIT as part of the German Excellence Initiative.

References

1. Aggarwal, C.C., Wolf, J.L., Yu, P.S., Procopiuc, C., Park, J.S.: Fast algorithms for projected clustering. *SIGMOD Rec.* **28**(2), 61–72 (1999)
2. Aksehirli, E., Goethals, B., Müller, E., Vreeken, J.: Cartification: a neighborhood preserving transformation for mining high dimensional data. In: *ICDM*, pp. 937–942, December 2013
3. Alon, U., Barkai, N., Notterman, D.A., Gish, K., Ybarra, S., Mack, D., Levine, A.J.: Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays. *PNAS* **96**(12), 6745–6750 (1999)
4. Assent, I., Krieger, R., Müller, E., Seidl, T.: INSCY: indexing subspace clusters with in-process-removal of redundancy. In: *ICDM*, pp. 719–724, December 2008
5. Beyer, K., Goldstein, J., Ramakrishnan, R., Shaft, U.: When is nearest neighbor meaningful? In: Beerl, C., Bruneman, P. (eds.) *ICDT 1999*. LNCS, vol. 1540, pp. 217–235. Springer, Heidelberg (1998)
6. Enrich, T., Kriegel, H.-P., Mamoulis, N., Niedermayer, J., Renz, M., Züfle, A.: Reverse-nearest neighbor queries on uncertain moving object trajectories. In: Bhowmick, S.S., Dyreson, C.E., Jensen, C.S., Lee, M.L., Muliantara, A., Thalheim, B. (eds.) *DASFAA 2014, Part II*. LNCS, vol. 8422, pp. 92–107. Springer, Heidelberg (2014)
7. Goldstein, M.: k_n -nearest neighbor classification. *IEEE TIT* **18**(5), 627–630 (1972)
8. Jarvis, R., Patrick, E.: Clustering using a similarity measure based on shared near neighbors. *IEEE TC* **C-22**(11), 1025–1034 (1973)
9. Kailing, K., Kriegel, H.P., Kröger, P.: Density-connected subspace clustering for high-dimensional data. In: *SDM*, vol. 4. SIAM (2004)
10. Kriegel, H.P., Kröger, P., Renz, M., Wurst, S.: A generic framework for efficient subspace clustering of high-dimensional data. In: *ICDM*, pp. 8, November 2005
11. McCann, S., Lowe, D.: Local naive bayes nearest neighbor for image classification. In: *CVPR*, pp. 3650–3656, June 2012

12. Moise, G., Sander, J.: Finding non-redundant, statistically significant regions in high dimensional data: a novel approach to projected and subspace clustering. In: KDD, KDD 2008, pp. 533–541. ACM, New York (2008)
13. Müller, E., Assent, I., Günnemann, S., Krieger, R., Seidl, T.: Relevant subspace clustering: Mining the most interesting non-redundant concepts in high dimensional data. In: ICDM, pp. 377–386, December 2009
14. Müller, E., Günnemann, S., Assent, I., Seidl, T.: Evaluating clustering in subspace projections of high dimensional data. PVLDB **2**(1), 1270–1281 (2009)
15. Nutt, C.L., Mani, D.R., Betensky, R.A., Tamayo, P., Cairncross, J.G., Ladd, C., Pohl, U., Hartmann, C., McLaughlin, M.E., Batchelor, T.T., Black, P.M., Deimling, A.V., Pomeroy, S.L., Golub, T.R., Louis, D.N.: Gene expression-based classification of malignant gliomas correlates better with survival than histological classification. *Cancer Res.* **63**(7), 1602–1607 (2003)
16. Park, Y., Park, S., Jung, W., Lee, S.G.: Reversed CF: a fast collaborative filtering algorithm using a k-nearest neighbor graph. *Expert Syst. Appl.* **42**(8), 4022–4028 (2015)
17. Rodriguez, A., Laio, A.: Clustering by fast search and find of density peaks. *Sci.* **344**(6191), 1492–1496 (2014)
18. Schneider, J., Vlachos, M.: Fast parameterless density-based clustering via random projections. In: CIKM, CIKM 2013, pp. 861–866. ACM, New York (2013)
19. Sequeira, K., Zaki, M.: SCHISM: A new approach for interesting subspace mining. In: ICDM, vol. 0, pp. 186–193. IEEE Computer Society (2004)
20. Sim, K., Gopalkrishnan, V., Zimek, A., Cong, G.: A survey on enhanced subspace clustering. *Data Min. Knowl. Disc.* **26**(2), 332–397 (2013)
21. Weinberger, K., Saul, L.: Unsupervised learning of image manifolds by semidefinite programming. *Int. J. Comput. Vis.* **70**(1), 77–90 (2006)