

# iVIBRATE: Interactive Visualization Based Framework for Clustering Large Datasets

Keke Chen      Ling Liu

College of Computing, Georgia Institute of Technology

{kekechen, lingliu}@cc.gatech.edu

## Abstract

With continued advances in communication network technology and sensing technology, there is an astounding growth in the amount of data produced and made available through the cyberspace. Efficient and high quality clustering of large datasets continues to be one of the most important problems in large-scale data analysis. A commonly used methodology for cluster analysis on large datasets is the three-phase framework of “sampling/summarization – iterative cluster analysis – disk-labeling”. There are three known problems with this framework, which demand effective solutions. The first problem is how to effectively define and validate irregularly shaped clusters, especially in large datasets. Automated algorithms and statistical methods are typically not effective in handling such particular clusters. The second problem is how to effectively label the entire data on disk (disk-labeling) without introducing additional errors, including the solutions for dealing with outliers, irregular clusters, and cluster boundary extension. The third problem is the lack of research about the issues for effectively integrating the three phases. In this paper, we describe iVIBRATE – an interactive-visualization based three-phase framework for clustering large datasets. The two main components of iVIBRATE are its VISTA visual cluster rendering subsystem, which invites human into the large-scale iterative clustering process through interactive visualization, and its Adaptive ClusterMap Labeling subsystem, which offers visualization-guided disk-labeling solutions that are effective in dealing with outliers, irregular clusters, and cluster boundary extension. Another important contribution of iVIBRATE development is the identification of special issues presented in integrating the two components and the sampling approach into a coherent framework, and the solutions to improve the reliability of the framework and to minimize the amount of errors generated throughout the cluster analysis process. We study the effectiveness of the iVIBRATE framework through a walkthrough example dataset of a million records and experimentally evaluate the iVIBRATE approach using both real-life datasets and synthetic datasets. Our results show that iVIBRATE can efficiently involve the user into the clustering process and generate high-quality clustering results for large datasets.

## 1 Introduction

Cluster analysis is a critical component in large-scale data analysis. Over the past decade, large datasets have been collected and analyzed in many application domains, varying from bioinformatics, information retrieval,

physics, geology, to marketing and business trend prediction. Many have reached the level of terabytes to petabytes [26, 19]. There is a growing demand for efficient and flexible clustering techniques that can adapt to the large datasets with complex cluster structure.

A dataset used in clustering is typically represented as a table  $D$  consisting of  $d$  dimensions (columns) and  $N$  records (rows). A record can represent an event, an observation or some meaningful entity in practice, while a dimension could be an attribute/aspect of the entity. The clustering algorithms try to partition the records into groups with similarity measure [29]. A dataset can be large in terms of the number of dimensions (dimensionality), the number of records, or both. The problem of high dimensionality (hundreds or thousands of dimensions) is typically addressed by feature selection and dimensionality reduction techniques [11, 38, 56, 3]. In this paper, we will focus on cluster analysis for the *numerical* datasets with a very large number of records ( $> 1$  million records) and a medium number of dimensions (usually  $< 50$  dimensions), assuming that the high dimensionality has been reduced before datasets entering the iVIBRATE framework for cluster analysis.

## 1.1 General Problems with Clustering Large Datasets

Several clustering algorithms have aimed at processing the *entire* dataset in linear or near linear time, such as WaveCluster [47], DBSCAN [15], and DENCLUE [23]. However, there are some drawbacks with these approaches.

**(1)Time Complexity of Iterative Cluster Analysis.** Typically, cluster analysis continues after the clustering algorithm finishes in a run, unless the user has evaluated, understood and accepted the clustering patterns or results. Therefore, the user needs to be really involved in the iterative process of “clustering and analysis/evaluation”. In this process, multiple clustering algorithms, and multiple runs of the same algorithm with different parameter settings can be tested and evaluated. Even for a clustering algorithm with linear computational complexity, running such an algorithm on a very large dataset in multiple times could become intolerable. Moreover, most cluster validation methods cost non-linear time [22, 30, 14]. When performed on the entire large dataset, the validation of clusters hinders the performance improvement for the entire iterative process.

**(2)Cluster Analysis on Representative Dataset vs. on Entire Dataset.** Bearing the above problems in mind, a number of approaches were proposed to perform clustering algorithms on the sample datasets or data summarization instead of the entire large dataset. For example, CURE [20] applies random sampling to get the sample data and then runs a hierarchical clustering algorithm on the sample data. BIRCH [58] summarizes the entire dataset into a CF-tree and then runs a hierarchical clustering algorithm on the CF-tree. This “sampling/summarization – iterative cluster analysis” framework has been commonly recognized as a practical solution to large-scale cluster analysis. Since the size of dataset is reduced with the sampling/summarization techniques, any typical clustering algorithms and cluster validation techniques that have acceptable non-linear computational complexity can be applied in cluster analysis.

However, the above two-phase framework misses the bridge between the clustering result of the representative dataset and the requirement of retrieving cluster labels for the entire large dataset. There are some typical

questions frequently asked about the entire large datasets, by the applications requiring cluster analysis: 1) what is the cluster label for a particular data record? 2) what are the data records that belong to a particular cluster? Therefore, we also need to extend the clustering result to the entire dataset, which introduces in the third phase - labeling data on disk with the intermediate clustering result. Previous research on clustering with the three-phase framework has been primarily focused on the first two phases. Surprisingly, very few studies have considered the efficiency and quality of disk-labeling phase.

Disk-labeling also provides the opportunity to review and correct the errors generated by the first two phases. For example, sampling/summarization tends to cause the loss of small clusters in the representative dataset. It is easy to understand that when the sample size is small, some small clusters might be lost. This is also true when summarization is done in a high granularity. For example, when a CF-tree in BIRCH is relatively small compared to the number of records, a leaf node will possibly cover a large spatial locality and we have to consider the points in the leaf as one cluster. However, in real-life applications, such a leaf node may contain several small clusters, which are not observable in the summarized clusters. Although many applications only consider the large clusters, small clusters may become very significant for some applications. Thus, there is a need for monitoring the small clusters, which are missed by applying sampling/summarization techniques in the first phase.

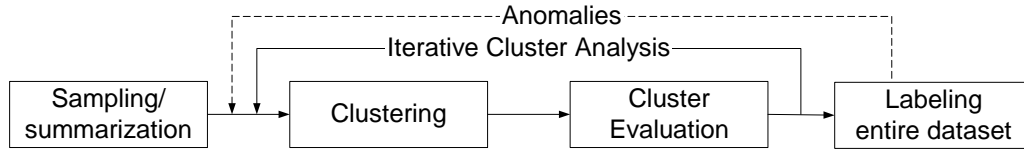


Figure 1: Three phases for cluster analysis of large datasets, (Sampling/summarization – Cluster Analysis – Disk Labeling)

**(3) Problems with Irregularly Shaped Clusters.** Many automated clustering algorithms work effectively in finding clusters in spherical or elongated shapes but they cannot handle arbitrarily shaped clusters very well [20]. Accordingly, the traditional validation methods, primarily based on statistical methods, cannot validate arbitrarily shaped clusters effectively [22, 46].

In some applications, irregularly shaped clusters may be formed by combining two regular clusters or by splitting one large cluster based on domain knowledge. Most of the existing clustering algorithms do not allow the user to participate in the clustering process until the clustering algorithm is completed. It is inconvenient to incorporate domain knowledge into the cluster analysis, or to allow the users to steer the clustering process for automated algorithms.

We argue that visualization techniques can play an important role in solving the problem of irregularly shaped clusters in large datasets. Some visualization-based algorithms, such as OPTICS [1], tried to find the arbitrarily shaped clusters, but they are applicable only to the datasets of a small manageable size, which is greatly restricted by the computer system capability, i.e., the memory capacity, the CPU speed and the screen size. Although visualization has shown unique advantages in handling arbitrarily shaped clusters, there is no re-

search extending the visualization techniques to clustering very large datasets, where sampling/summarization approaches have to be applied as the preprocessing stage of visualization systems.

**(4)Problems with Disk Labeling.** When disk labeling is used as the last phase of large-dataset clustering, it uses the intermediate clustering result to assign a cluster label to each data record on disk. Without an effective labeling phase, a large amount of errors can be generated when the intermediate clustering result is extended to the entire dataset.

The quality of labeling depends primarily on the precise description of cluster boundaries. All existing labeling algorithms are based on very rough cluster descriptions [30], such as a centroid or a set of representative cluster boundary points for a cluster. A typical labeling algorithm assigns each data record on disk to a cluster that has centroid or the representative boundary points closest to this data record. Centroid-based labeling (CBL) uses cluster center (centroid) only to represent a cluster; Representative-point-based labeling (RPBL) uses a set of representative points on cluster boundary to describe the cluster. The latter is better since it provides more information about the cluster boundary. With RPBL, the quality of boundary description mainly depends on the number of representative points, which could be very large for some irregular cluster shapes or for large clusters. However, it is not convenient for the user to determine how many representative points are sufficient for a particular dataset.

Moreover, in many cases, the cluster boundary cannot be defined precisely with the sample dataset. More importantly, the cluster boundaries continue to evolve as we incorporate more the labeled records during the disk labeling phase. We call it the “cluster boundary extension” problem and will describe it in detail with our solutions later in the paper.

## 1.2 The Scope and Contribution of the Paper

We have summarized four key problems in clustering large datasets: 1) the three-phase framework is necessary for reducing the time complexity of iterative cluster analysis; 2) extending clustering result on the representative dataset to the entire large dataset can raise problems; 3) clustering and validating irregularly shaped clusters in large datasets is important but difficult; and 4) existing disk-labeling algorithms may result in large errors when used to extend the intermediate clustering result to the entire dataset.

We also explicitly identify the boundary extension problem as an important challenge in the labeling phase. The point density over the boundary could increase significantly as the number of labeled records increases, which naturally leads to the boundary extension problem. When the sample size is much smaller than the size of the large dataset, e.g.  $< 1\%$  of the original data records, boundary extension can cause large error too if the labeling algorithms do not consider this phenomenon. Figure 2 shows the clusters evolving from the small ones in the representative dataset to the larger ones in the entire dataset, where boundary extension results in significant difference in cluster definition.

Boundary extension can cause another two important problems. 1) For the regular spherical clusters as shown in Figure 2, existing labeling algorithms usually mistakenly assign all outliers to the nearby clusters, or treat

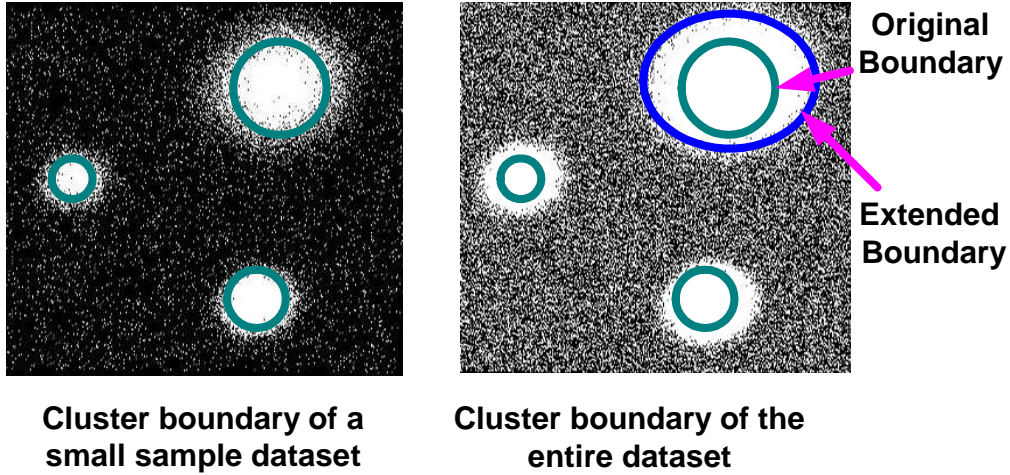


Figure 2: Comparing the cluster boundary of small and large dataset (data points are white)

the cluster members out of the original boundary as outliers. As a result, none of them can deal with outliers and boundary extension effectively. For irregular cluster boundary, the situation becomes worse. 2) Boundary extension might also result in the overlapping of different clusters, which are originally separated in the representative set. This may require us to adjust the cluster definition.

To address all of the above problems, we propose the iVIBRATE framework — an interactive visualization based three-phase framework for clustering large datasets. The iVIBRATE framework includes the three phases “Sampling — Visual Cluster Analysis — Visualization-based Adaptive Disk-labeling”. In this paper, we introduce the two main components: visual cluster analysis and visualization-based adaptive disk-labeling, while focusing on the important issues in integrating the three components in the iVIBRATE framework, and demonstrating how to analyze very large datasets with the iVIBRATE framework.

In the clustering phase, we use visual cluster analysis, including visual clustering and visual cluster validation, to allow the user to be more effectively involved in the iterative cluster analysis, reducing both the length of single iteration and the number of iterations. We develop a visual cluster rendering system, VISTA, to perform “visual cluster analysis”. The VISTA system interactively visualizes the multi-dimensional datasets (usually  $<50$  dimensions). It allows the user to interactively observe potential clusters in a series of continuously changing visualizations. More importantly, it can incorporate the algorithmic clustering results, and serve as an effective validation and refinement tool for identifying irregular clusters.

In the disk-labeling phase, we develop the Adaptive ClusterMap Labeling. ClusterMap encodes the irregular cluster boundaries defined on the visualization generated by the VISTA subsystem and clearly distinguishes the outliers. It extends the cluster boundary by clearly distinguishing the outliers, continuously detecting irregular shaped clusters, and visually monitoring the anomalies in labeling process. As a result, the Adaptive ClusterMap labeling generates fewer errors than the existing disk-labeling algorithms, and thus effectively reduces the error introduced by the first two phases.

When the three phases are integrated in the iVIBRATE framework, errors caused by the improper operations

in the prior phases may propagate to the later phases. Visualization helps to detect and reduce such errors. We identify and analyze the issues related to the integration, and develop the theory and tools to monitor the possible errors. The iVIBRATE framework is evaluated with real and synthetic datasets and results show that iVIBRATE can efficiently incorporate the user into the clustering process to generate high-quality clustering results for large datasets.

The rest of the paper is organized as follows. Section 2 reviews some related work. Section 3 presents the iVIBRATE framework. Section 4 and Section 5 introduce the two main components of the framework: the VISTA visual cluster rendering system and the ClusterMap labeling algorithms. Section 6 describes the problems and the solutions for integrating the three phases into the iVIBRATE framework. Section 7 reports the experimental results, demonstrating that the iVIBRATE approach can effectively discover and validate irregular clusters, and the intermediate clustering result can be effectively extended to the entire large dataset through adaptive ClusterMap labeling. We also present an example using the steps defined in iVIBRATE framework to explore the very large real dataset: census dataset, in section 8.

## 2 Related Work

**Clustering Large Data.** We have described four challenges related to clustering large datasets: time complexity (scalability), sampling/summarization based clustering, irregular clusters and disk-labeling. Although each of these issues has been studied in the literature, there is surprisingly little study on how they impact on the cluster quality when sampling/summarization, iterative cluster analysis, and disk labeling are integrated into a unifying framework for large datasets with complex cluster structures.

Concretely, time complexity of clustering algorithm has been addressed from early on. K-means algorithm [29] is the most popular algorithm with linear time complexity. However, most studies related to K-means assume that the clusters are in spherical shapes. Recently, there are some other algorithms [47, 15, 23] started looking at the problem of clustering irregular clusters in linear/near-linear time.

A general cluster analysis framework is described in a review paper by Jain [30], showing that cluster analysis is usually an iterative process. One approach to address the scalability of iterative clustering analysis is the use of the “sampling(summarization) – clustering–labeling” framework, represented by CURE [20] and BIRCH [58]. However, the labeling phase and interactions between the phases are not sufficiently addressed.

Regarding summarization/sampling phase, similar to BIRCH summarization, Bradley et al. [4] suggest using sufficient statistics to model the data summary. In comparison to summarization, sampling is used more extensively for data analysis: commercial vendors of statistical packages (e.g., SAS) typically use uniform sampling to handle large datasets. Vitter’s reservoir sampling [51] presents another efficient uniform sampling technique. The main problem with uniform sampling is the loss of small clusters. CURE proposed a method to estimate the minimum sample size given that the size of dataset and the smallest size of clusters are known. The minimum sample size shall increase dramatically with the increase of dataset size. Thus, for a fixed sample size, the small clusters should be monitored using a different approach to maintain the consistency in cluster analysis

quality, especially in the disk labeling process. Another popular sampling approach is called density-biased sampling proposed by Palmer et al. [41]. A density-biased sampling preserves the small clusters in sampling process. However, this technique also reduces the actual size of large clusters, introducing too much inconsistency between the clusters in the sample set and the actual clusters in the entire large dataset. These problems call for effective solutions that can improve the quality of both the iterative cluster analysis phase and the disk labeling phase in order to reduce the impact of the errors introduced by sampling on the overall quality of the data clustering process.

Another important issue in clustering large datasets is the problem of arbitrarily shaped clusters or so called irregular shaped clusters. Dealing with arbitrarily shaped clusters is well-recognized as a challenging problem in clustering research community. A number of clustering algorithms have aimed at this particular problem, such as CURE [20], CHAMELEON [32], DBSCAN [15], DBCLASD [54], WaveCluster [47], DENCLUE [23] and so on. The semi-automatic algorithm OPTICS [1], derived from the DBSCAN algorithm, shows that visualization can be very helpful in cluster analysis. However, all these algorithms are known to be effective only in low dimensional (typically,  $<10D$ ) datasets or in small/medium datasets with hundreds or thousands of records.

The disk labeling solutions so far heavily depend on the concrete cluster representations generated at the iterative cluster analysis phase. Existing cluster representations can be classified into four categories [30] : centroid-based, boundary-point-based (representative-point-based), classification-tree-based and rule-based representations. Centroid is the most popular representation. Many clustering algorithms only generate centroids, for example, K-means and most hierarchical algorithms. Good representative boundary points are often difficult to extract. The most typical algorithm for generating the representative points is CURE [20]. Classification-tree-based and rule-based representations are equivalent (each path in the classification tree can be represented as a rule) and not convenient to describe high-dimensional data or complicated clustering structure.

**Document Clustering in IR and Linkage-based clustering in Network Analysis.** Most of the work on clustering large datasets can be roughly categorized into three areas: scientific/business data clustering [29, 30], document clustering [53, 12, 44, 28, 45, 57, 49], and linkage-based clustering for large scale network analysis [30, 21, 40].

In scientific/business data clustering, the data is already formalized as a set of multi-dimensional vectors (i.e., a table). However, in document clustering, the original data is text data. Most document clustering techniques are focused on the two steps before applying clustering algorithms: extracting keywords/constructing the numerical features [56], and defining similarity measures [2]. Given the vector representation of documents and the similarity measure, document clustering is to some extent similar to scientific/business data clustering. Since document collection has become larger and larger with the wide spread of Internet-based applications, we expect that the iVIBRATE framework described in this paper can also be applied to clustering of large sets of documents.

Graph mining or linkage-based clustering [30] has received a growing interest in the recent years due to increased interests in analyzing large-scale networks, such as peer to peer online communities [42], and social

network [40]. Linkage-based clustering is also used in clustering categorical datasets [21]. In comparison, most of the business/scientific data clustering algorithms utilize the distances between multi-dimensional data points (records) to compute and derive data clusters, while most of the linkage-based clustering algorithms utilize the node connectivity as a main measurement to understand and derive interesting clustering structure of the network. Thus, the linkage-based clustering algorithms aim at finding communities in networks – groups of vertices within which connections are dense but between which connections are sparser. A commonality of data clustering, document clustering, and node clustering is the fact that they all emphasize on efficient algorithms to speed up the clustering process of large datasets. However, the subtle differences between distance based measure and connectivity-based measure may influence how the clustering algorithms are devised and what factors are critical to the performance of the algorithms.

Due to the scope of our paper, we will confine our discussion to the general clustering problem to the business and scientific datasets.

**Visualization of Multidimensional Data.** Information visualization has demonstrated great advantages in multi-dimensional data analysis. Here we only discuss the scatter-plot-based techniques because they are the most intuitive techniques for cluster visualization. The early research on general 2D-plot-based data visualization is Grand Tour and Projection Pursuit [10]. Since there are numerous projections from a multidimensional data space to a 2D space, Grand Tour and the Project Pursuit automatically interpolate the intermediate projections to connect several pre-selected projections. Yang [55] utilizes the Grand Tour technique to show projections in an animation. Dhillon, et al. [13] aimed at precise visualization of the clusters, but the technique is only effective for 3 clusters. When more than 3 clusters exist, this method needs the help of Grand Tour techniques. In fact, the mathematical interpolation used in these models is computationally complex. since we can fully utilize the human visual ability and the property of exploration tasks to simplify the visualization model and get the satisfactory visualization more efficiently. The iVIBRATE visual rendering development (VISTA) has shown that the special properties of the data exploration task (i.e., the clustering task in the context of this paper) can be utilized in user interaction to find the satisfactory visualization more efficiently.

In addition, there are static multidimensional visualization techniques such as Scatterplot Matrices, coplots, Parallel Coordinates [27], Glyphs [37], multidimensional stacking [35], prosection [8] and FastMap based visualization [16]. A nice tool, XvmdvTool [52], implements some of the above static visualization techniques. However, these techniques are generally not designed for effectively discovering the clusters in very large datasets. Many of them are limited by the dimensionality (10-20 dimensions at maximum), for example, Scatterplot Matrices and Parallel Coordinates. None of them are designed for very large number of records – neither can they be easily extended to the three-phase cluster analysis framework for large datasets.

Star Coordinates [31] is a visualization system designed to interactively visualize and analyze clusters. We utilize the form of Star Coordinates and build the  $\alpha$ -mapping model in our system.  $\alpha$ -mapping model extends the ability of the original mapping [6]. RadViz [25] utilizes the same coordinates system with a non-linear mapping function, but it is not suitable for cluster exploration. HD-Eye [24] is another interesting interactive visual clustering system. HD-Eye visualizes the density-plot of the interesting projection of any two of the



$k$  dimensions. It uses icons to represent the clusters and the relationship between the clusters. However, it is difficult for the user to synthesize all of the interesting 2D projections to find the general pattern of the clusters. In fact, visually determining the cluster distribution solely through user interaction is not necessary. A more practical approach is to incorporate all available clustering information, such as the algorithmic clustering results and the domain knowledge, into the visual cluster exploration.

### 3 The iVIBRATE Framework

In this section, we first briefly describe the motivation and the design ideas of the iVIBRATE development, and then present the components and working mechanism of iVIBRATE.

**Motivation.** In the three-phase framework, the cluster analysis involves the "clustering - analysis/evaluation" iteration, which can be concretely described in the following steps:

1. Run the clustering algorithms with the initial setting of parameters.
2. Analyze the clustering results with statistical measures and domain knowledge.
3. If the result obtained in Step 2 is not satisfactory, adjust the parameters and re-run the clustering algorithms, then go to Step 2 to evaluate the clustering result again until the satisfactory result is obtained.
4. If the result is satisfactory, then perform post-processing, which may include labeling the data records on disk with the cluster labels.

**Open problems.** We first discuss a number of open problems in the above steps, and then we describe how iVIBRATE addresses these problems. Traditional statistical methods, such as variance and intra/inter cluster similarity, are typically used in Step 2, which assume that the shape of cluster structure is hyper-sphere or hyper-ellipse. As a result, these traditional statistical methods have difficulty in effectively validating the irregular cluster shapes [14, 22]. Moreover, with automated algorithms, it is almost impossible to incorporate the domain knowledge. The critical task in step 3 is to learn and determine the appropriate parameter settings. In the absence of user experience, it is extremely difficult to determine the appropriate setting of parameters for a new run. For example, CURE [20] requires the number of representative points and the shrink factor. DBSCAN [15] needs proper  $\varepsilon$  and *MinPts* to get satisfactory clusters. DENCLUE [23] needs to define the smoothness level and the significance level. These parameter settings are different from dataset to dataset and depend primarily on the users to try different parameters and find the "best" set of parameters by hand. Therefore, there is a need for facilitating the users to find the appropriate parameter setting when automated algorithms are applied. Finally, a coarse labeling algorithm tends to deteriorate the intermediate clustering result.

Bearing these problems in mind, we observed that, if the step 2 and 3 can be carefully combined together, which means that the user can perform evaluation in the course of clustering and be able to refine the clusters interactively, the length of an iteration would also be greatly reduced. In addition, the user would understand

more about the dataset and thus be more confident in their judgment of the clustering results. This motivates us to develop and promote the interactive visual cluster rendering approach.

**Cluster Visualization and Visual Validation.** Cluster visualization can improve the understanding of clusters. Former studies [34] in the area of visual data exploration support the notion that visual exploration can help in cognition. Visual representations can be very powerful in revealing trends, highlighting outliers, showing clusters, and exposing gaps. According to the paper [48], with the right coding, human pre-attentive perceptual skills enable users to recognize patterns, spot outliers, identify gaps and find clusters in a few hundred milliseconds. In addition, it does not require the knowledge of complex mathematical/statistical algorithms or parameters [33].

Visualization is known to be the most intuitive method for validating clusters, especially clusters in irregular shape. Since the geometry and density features of clusters, which are derived from the distance (similarity) relationship, determine the validity of the clustering results, many clustering algorithms in literature use 2D-plot of clustering results to visually validate their effectiveness on 2D experimental datasets.

**Static vs. Dynamic Data Visualization.** In general, multi-dimensional data visualization can be categorized into static visualization or dynamic visualization. Static visualization displays the data with a fixed set of parameters, while dynamic cluster visualization allows the user to adjust a set of parameters, resulting in a series of continuously changed visualizations. It is commonly believed that static visualization is not sufficient in visualizing multi-dimensional data [33, 48], and it has been shown that clusters can hardly be satisfactorily preserved in a static visualization [10, 13]. Therefore, we consider using interactive dynamic cluster visualization in the iVIBRATE framework. We observe that a cluster can always be preserved as a point-cloud in visual space through linear mappings. The only problem is that, the point-clouds may overlap with one another and to find certain mapping that can satisfactorily separate the point clouds is mathematically complex. In the iVIBRATE framework, we incorporate the combination of visual cluster clues and interactive rendering into the iterative clustering analysis and refine the algorithmic clustering results with visual cluster clues and interactive rendering, which enables us to identify these point-cloud overlaps quickly and intuitively.

**Visualization-based Disk-labeling.** Another unique characteristic of iVIBRATE is its visualization-based disk-labeling algorithms. We argue that a fine cluster visualization of a dataset can serve as the visual clustering pattern of this dataset, where the cluster boundary can be precisely described and most outliers can be clearly distinguished. We develop the basic ClusterMap labeling algorithm for obtaining better description of cluster boundary and higher quality of disk-labeling. We extend the basic ClusterMap algorithm to adapt to the boundary extension, and refer to this enhanced algorithm the Adaptive ClusterMap labeling algorithm.

**Components and Working Mechanism.** Figure 3 sketches the main components of iVIBRATE. We briefly describe each of the main components and the method used to integrating them as follows.

- **Visual Cluster Rendering** The VISTA system can be independently used to render the clusters in a dataset without incorporating any external information. It can also visualize the result of an automated clustering algorithm or use the result to guide the interactive rendering. By interactively adjusting the

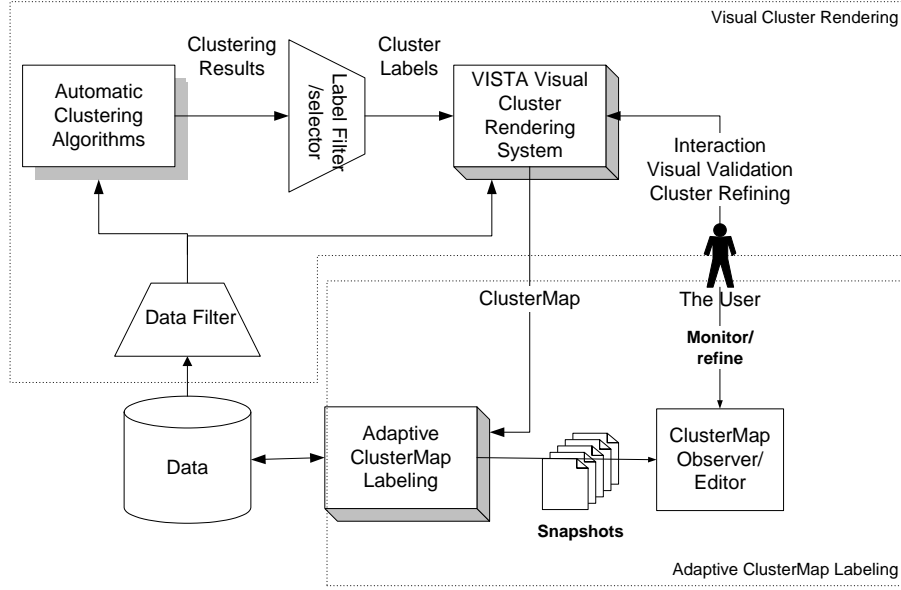


Figure 3: The iVIBRATE framework

visual parameters, the user can visually validate the algorithmic clustering results through multiple continuously changing visualizations. Using a couple of rendering rules, which are easy to learn, the user can quickly identify cluster overlaps and improve the cluster visualization. In addition, it also allows the user to conveniently incorporate domain knowledge into cluster definition through visual rendering operations. Semi-automated rendering method is also provided for larger number of dimensions ( $> 50D$  and  $< 100D$ ).

Data Filter prepares the data for visualization. It handles missing values and normalizes the data. If the dimensionality is too high, dimensionality reduction techniques or feature selection might be applied to get a manageable number of dimensions. When the datasets grow past a million items and cannot be easily visualized on a computer display, Data Filter also uniformly samples data to create a manageable representative dataset, or extracts certain relevant subsets, for example, one specific cluster.

Label Selector selects the clustering result that will be used in visualization. While a clustering algorithm finishes, it usually assigns a label to each data record. Label Filter extracts part of the labels corresponding to the data records extracted by Data Filter.

- **Adaptive ClusterMap Labeling** In the iVIBRATE framework, we introduce ClusterMap — a new cluster presentation, and the associated post-processing methods. ClusterMap makes the labeling result most consistent with the defined cluster distribution. Adaptive ClusterMap also automatically adjusts the cluster boundary according to the accumulated labeled data records.

ClusterMap Observer is an interactive monitoring tool. It observes the snapshots, i.e. the changing ClusterMaps during labeling. The snapshots may provide information about the bias of the representative sample set, for example, the missing small clusters, and the anomalies in labeling, such as inappropriate initial cluster boundaries 6.2.

A user of iVIBRATE will perform the cluster analysis on a large dataset in the following seven steps. 1) The large dataset is sampled to get a subset in manageable size (such as thousands or tens of thousands records sampled from over one million of records). 2) The sample set is used as an input to the selected automatic clustering algorithms and to the VISTA visual rendering subsystem. The algorithmic clustering result provides helpful information in the visual cluster rendering process. 3) The user interacts with the VISTA visual cluster rendering system to find the satisfactory visualization, which visualizes the clusters in well-separated areas. Since human vision is very sensitive to the gap between point clouds, which implies the boundary of clusters, the interactive rendering works very well in refining vague boundaries or irregular cluster shapes. 4) A ClusterMap is then defined on the satisfactory cluster visualization and used as the initial pattern in ClusterMap labeling. 5) The labeling process will adapt the boundary extension and refine the cluster definition in one pass through the entire dataset. An additional pass might be needed to reorganize the entire dataset for fast processing of queries. During the labeling process, the snapshots are saved periodically, which are then used to monitor the abnormality during the labeling process. 6) The user can use the ClusterMap Observer to check the snapshots and refine the extended ClusterMap. 7) To further observe the small clusters that may be omitted in the sampling process, the data filtering component is used to filter out the labeled outliers and performs sampling/visual rendering on the sampled outliers again (for details, see section 6.3).

In the following sections, we will introduce the two visualization-based subsystems with our focus on integrating the components into the framework so that the errors are effectively controlled.

## 4 VISTA Visual Cluster Rendering System

A main challenge in cluster visualization is cluster preservation, i.e., visualizing multi-dimensional datasets in 2D/3D visual space, while preserving the cluster structure. Previous studies have shown that preserving cluster structure precisely in static visualization is very difficult, if not impossible, and computationally expensive [33, 10, 55, 31, 48]. An emerging and more practical mechanism to address this problem is to allow the user to interactively explore the dataset [33] and to distinguish the visibly inaccurate cluster structure, such as cluster overlapping, broken clusters and false clusters (the situation where the outliers in the original space are mapped to the same visual area) through visual rendering operations.

The iVIBRATE visual cluster rendering subsystem (VISTA) is designed to be a dynamic visual cluster exploration system. It uses a visualization model, characterized by the max-min normalization and the  $\alpha$ -mapping to produce a linear transformation that maps a multi-dimensional data point onto a data point in a 2D visual space, each 2D point is annotated with a star-coordinate. This mapping model provides a set of visually adjustable parameters for each dimension through the  $\alpha$  parameters. By continuously changing one of the parameters, the user can see the dataset from different perspectives. Since the linear mapping does not break the clusters, the clusters in multi-dimensional space are still visualized as dense point clouds (the “visual clusters”) in 2D space. And the visible “gaps” between the visual clusters in 2D visual space indicate the *real gaps* between point clouds in the original high dimensional space. However, visible overlapping between the visual clusters

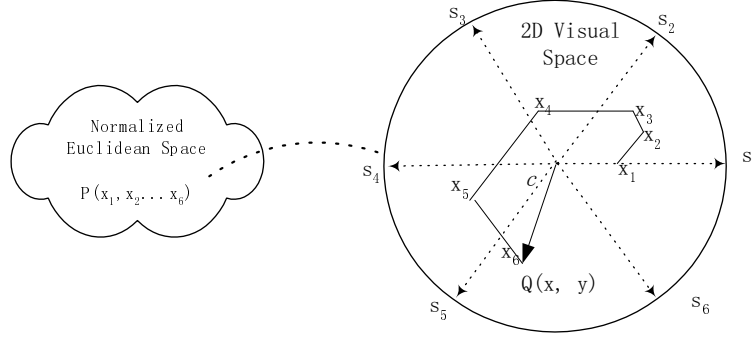


Figure 4: Illustration of  $\alpha$ -mapping with  $k = 6$

in the 2D space, i.e. the point clouds, may happen with certain parameter settings. We have developed a set of interactive operations and designed several heuristic rendering rules in order to efficiently distinguish the visual cluster overlaps. These developments have shown to be quite effective in achieving desired rendering efficiency.

In the current prototype of iVIBRATE VISTA subsystem, the system processes only the datasets that support the distance/similarity function defined by Euclidean distance, which is the most commonly used distance measure in the applications. For presentation convenience, in the rest of the paper we refer to the datasets we consider as the Euclidean datasets.

#### 4.1 The Visualization Model

The VISTA visualization model consists of two linear mappings – max-min normalization followed by  $\alpha$ -mapping. The columns and rows of the datasets referred in the following discussion represent the dimensions and data records, respectively.

**Max-min normalization** is used to normalize the columns in the datasets in order to eliminate the dominating effect of large-valued columns. For a column with value bounds  $[\min, \max]$ , max-min normalization scales a value  $v$  in the column into  $[-1, 1]$  as follows:

$$v' = \frac{2(v - \min)}{\max - \min} - 1 \quad (1)$$

where  $v$  is the original value and  $v'$  is the normalized value.  **$\alpha$ -mapping** maps  $k$ -D points onto the 2D visual space while providing the convenience of visual parameter tuning. We describe  $\alpha$ -mapping as follows. Let a 2D point  $Q(x, y)$  represent the image of a  $k$ -dimensional ( $k$ -D) max-min normalized data point  $P(x_1, \dots, x_i, \dots, x_k)$ ,  $x_i \in [-1, 1]$  in 2D space.  $Q(x, y)$  is determined by the average of the vector sum of the  $k$  vectors  $\vec{s}_i x_i$ , where  $\vec{s}_i = (\cos(\theta_i), \sin(\theta_i))$ ,  $i = 1 \dots k$  and  $\theta_i \in [0, 2\pi]$  are the star coordinates [31] that represent the  $k$  dimensions on the 2D visual space. Formula 2 defines  $\alpha$ -mapping.

$$A_{(\theta_1, \dots, \theta_k)}(x_1, \dots, x_k, \alpha_1, \dots, \alpha_k) = (c/k) \sum_{i=1}^k \alpha_i x_i \vec{s}_i - \vec{o} \quad (2)$$

i.e. a 2D point  $Q(x, y)$  is determined by

$$\{x, y\} = \{(c/k) \sum_{i=1}^k \alpha_i x_i \cos(\theta_i) - x_0, (c/k) \sum_{i=1}^k \alpha_i x_i \sin(\theta_i) - y_0\} \quad (3)$$

Here,  $\alpha_i (i = 1 \dots k, \alpha_i \in [-1, 1])$  provides the visually adjustable parameters, one for each of the  $k$  dimensions.  $\alpha_i \in [-1, 1]$  covers a considerable range of mapping. Experimental results show that this range combined with the scaling factor  $c$  is effective enough for finding satisfactory visualization.  $\theta_i$  is set to  $2i\pi/k$  initially and can be adjusted afterwards. We also proved that adjusting  $\theta$  values is often equivalent to a pair of  $\alpha$  adjustment plus zooming [6]. Therefore,  $\theta_i$  are always considered as a set of fixed numbers.  $\vec{o} = (x_0, y_0)$  is the center of the display area.

$\alpha$ -mapping is a linear mapping, with any fixed set of  $\alpha$  values. Without loss of generality, we set the center translation  $(x_0, y_0)$  as  $(0, 0)$ . The mapping  $A_{\alpha_1=a_1, \dots, \alpha_k=a_k}(x_1, \dots, x_k)$  can be represented as the following transformation.

$$A_{\alpha_1=a_1, \dots, \alpha_k=a_k}(x_1, \dots, x_k) = \begin{bmatrix} \cos(\theta_1) \cdots \cos(\theta_k) \\ \sin(\theta_1) \cdots \sin(\theta_k) \end{bmatrix} \begin{bmatrix} a_1 & & \\ & \ddots & \\ & & a_k \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_k \end{bmatrix}$$

It is known that linear mapping does not break clusters but may cause overlapped clusters [33, 18]. Since  $\alpha$ -mapping is linear, there are no “broken clusters” in the visualization, i.e., the visual gaps between the point clouds reflect the real gaps between the clusters in the original high-dimensional space. All we need to do is to separate the possibly overlapped clusters, which can be achieved with the help of dynamic visualization through interactive operations.

The mapping is adjustable by  $\alpha_i$ . By tuning  $\alpha_i$  continuously, we can see the effect of the  $i$ -th dimension to cluster distribution in a series of smoothly changing visualizations, which can provide important clustering clues. The dimensions that are important to clustering will cause *significant changes* to the visualization as the corresponding  $\alpha$  values are continuously changed.

$\alpha$ -mapping based visualization is implemented in the VISTA subsystem as shown in Figure 5. The coordinates are arranged around the display center and the  $\alpha$ -widgets are designed for interactively adjusting each  $\alpha$  value. However, the above visual design also limits the number of dimensions that can be visualized and manually manipulated. In the current prototype of VISTA, the users can comfortably manually render up to 50 dimensions. Although the system can visualize more than 50 dimensions, we suggest using the semi-automated rendering method that will be introduced in section 4.3.

## 4.2 The Rules for Interactive Visual Rendering

To understand the basic visual rendering rules, we should investigate the dynamic properties of the visualization model, especially, the most important interactive operation –  $\alpha$ -parameter adjustment (or simply,  $\alpha$ -

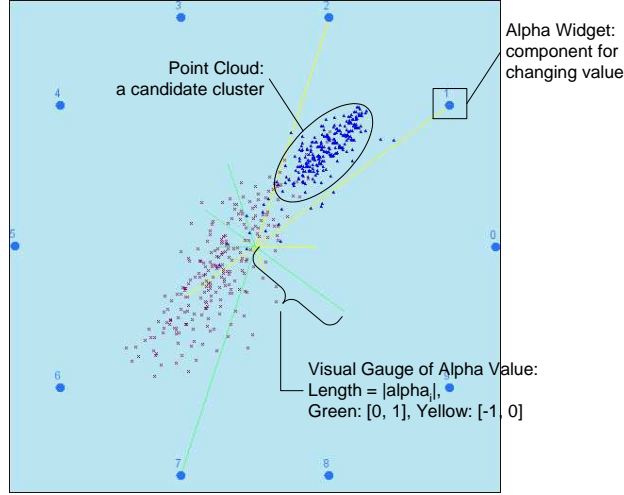


Figure 5: An implementation of  $\alpha$ -mapping

adjustment).  $\alpha$ -adjustment changes the parameters defined in Eq. (2). Each change refreshes the visualization in real time (about a couple of hundred milliseconds, depending on different hardware configuration and the size of the dataset), generating dynamically changing visualizations.  $\alpha$ -adjustment enables the user to find the *dominating dimensions*, to observe the dataset from different perspectives, and to distinguish the real clusters from overlapping in continuously changing visualizations. Therefore, we should investigate the dynamic properties of  $\alpha$ -adjustment operation.

Continuous  $\alpha$ -parameter adjustment of one dimension reveals the effect of this dimension on the entire visualization. Let  $X(x_1, \dots, x_k)$  and  $Y(y_1, \dots, y_k)$ ,  $x_i, y_i \in [-1, 1]$  represent any two normalized points in  $k$ -D space. Let  $\|\vec{v}\|$  represent the length of vector  $\vec{v}$ . We define the *visual distance* between  $X$  and  $Y$  is:

$$\begin{aligned} vdist(X, Y) &= \|A(x_1, \dots, x_k, \alpha_1, \dots, \alpha_i, \dots, \alpha_k) - A(y_1, \dots, y_k, \alpha_1, \dots, \alpha_i, \dots, \alpha_k)\| \\ &= \|(c/k) \sum_{i=1}^k \alpha_i (x_i - y_i) \vec{s}_i\| \end{aligned} \quad (4)$$

which means if  $x_i$  and  $y_i$  are close, changing  $\alpha_i$  does not change the visual distance between  $X$  and  $Y$  a lot – the dynamic visual result is that  $X$  and  $Y$  are moving together when  $\alpha_i$  changes. Meanwhile, neighboring points in  $k$ -D space also have similar values in each dimension as Euclidean distance is employed. Thus, we can conclude that the neighboring points in  $k$ -D space, which should belong to one cluster, not only are close to each other in 2D space, but also tend to move together in any  $\alpha$ -adjustment; while those points that are far away from each other in  $k$ -D space may move together in some  $\alpha$ -adjustment but definitely not in all  $\alpha$ -adjustments. This property makes  $\alpha$ -adjustment very effective in revealing the visual cluster overlaps introduced by  $\alpha$ -mapping. In addition, point movement can also reveal the value distribution of individual dimension. If we adjust the  $\alpha$  value of the dimension  $i$  only, the point movement can be represented by:

$$\begin{aligned} \Delta(i) &= A(x_1, \dots, x_k, \alpha_1, \dots, \alpha_i, \dots, \alpha_k) - A(x_1, \dots, x_k, \alpha_1, \dots, \alpha'_i, \dots, \alpha_k) \\ &= (c/k)(\alpha_i - \alpha'_i)x_i \vec{s}_i \end{aligned} \quad (5)$$

which means that the points having larger  $x_i$  will be moving “faster” along the  $i$ -th coordinate, and those having the similar  $x_i$  moving in a similar way. The initial setting of  $\alpha$  values may not reveal the distribution of an individual dimension as 6 shows. However, by looking at the density centers (the moving point cloud) along the  $i$ -th axis as  $\alpha_i$  changes, we can easily estimate the value distribution in  $i$ -th dimension. In Figure 6, we sketch that point movement and point distribution can be interpreted intuitively with each other.

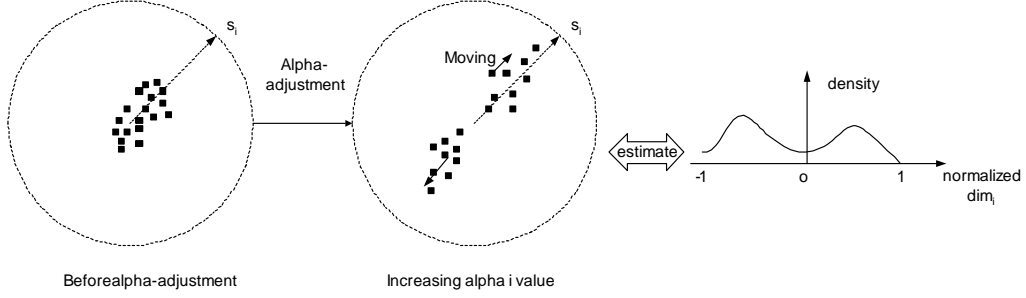


Figure 6:  $\alpha$ -adjustment, dimensional data distribution and point movement

In interactive visual rendering, some dimensions show “significant change” on visualization in continuous  $\alpha$ -adjustment, i.e., changing its  $\alpha$  value results in distinct point clouds moving in different directions, or causes the emergence of visible “gaps” between the point clouds. These dimensions play important roles in visual cluster rendering and thus we name them as “visually dominating dimensions”, and the others as “the fine-tuning dimensions”. The dominating dimensions usually have skewed distributions, where more than one distinctive mode exist on the distribution curve. For example, dimensions with near uniform distribution are definitely not dominating dimensions and dimensions with normal distribution are also less likely to be dominating but possible useful in the refinement.

Since the main goal of VISTA interactions is to distinguish the possible visual cluster overlapping, we can apply the following rules in visual rendering:

**Visual Rendering Rule 1.** *Sequentially try rendering each dimension. If the dimension is a visually dominating dimension, increase its  $\alpha$  value to certain degree so that the main point clouds are satisfactorily distinguished.*

**Visual Rendering Rule 2.** *Use the fine-tuning dimensions to polish the visualization. Adjust their  $\alpha$  values finely so that the visualization clearly shows the cluster boundaries.*

Guided by the above simple visual rendering rules, a trained user can easily find the satisfactory visualizations. While combined with the cluster labels generated by coarse automatic clustering algorithms (for example, K-Means algorithm), the rendering becomes even easier. During the rendering process, we can intuitively validate the algorithmic clustering results and conveniently incorporate the domain knowledge into the clustering process [6], which are difficult for most automated clustering algorithms.



### 4.3 Semi-automated Rendering

When the number of dimensions grows to a considerably large number ( $> 50$  and  $< 100$  dimensions to the VISTA system), manually rendering the dimensions becomes a difficult job. In the VISTA subsystem, we provide a semi-automated rendering method to automate the rendering of this type of datasets. Together with the visual rendering rules we have presented, the semi-automated rendering method can be quite efficient.

Concretely, our semi-automated rendering is performed in two stages: automatic random rendering (ARR) followed by automatic dimension-by-dimension rendering (ADDR). A simple version of random rendering is defined as follows. Let a dataset  $S$  have  $d$  dimensions and  $N$  records. Each dimension  $i$  ( $1 \leq i \leq k$ ) is associated with an initial  $\alpha$  value, say  $\alpha_i$ . Random rendering can be done in any number of rounds until some rough pattern of cluster distribution is observed. In each round, the  $\alpha_i$  value is changed by a small constant amount  $\epsilon$ , ( $0 < \epsilon < 1$ ), but the direction (increase or decrease) is randomly chosen for each  $\alpha_i$ . Since the  $\alpha$  values are bounded by 1 and -1, the change is “bounced” back at the bounds. By changing the  $\alpha$  values in this way, rather than randomly assigning them in each round, we can observe that the visualization is more smoothly changed. This type of continuity between the nearby visualizations is important to the user, since the user’s reaction might be slower than the change of visualization. When a nice rough pattern is observed, a few successive visualizations will be similar to the observed one, allowing the user to stop ARR around the satisfactory pattern.

After a rough pattern is observed in random rendering, we switch the automated rendering from ARR to ADDR for further refinement. In ADDR, for each dimension  $i$ ,  $\alpha_i$  is continuously and uniformly changed between  $[-1, 1]$ , namely,  $\alpha_i$  increases by  $\epsilon$  at each step from -1 to 1, and decreases by  $-\epsilon$  from 1 to -1. When a more refined cluster visualization is accepted by the user, ADDR for dimension  $i$  is stopped and moved to the next dimension. ARR in the first stage helps to quickly find some sketch of the cluster distribution and ADDR in

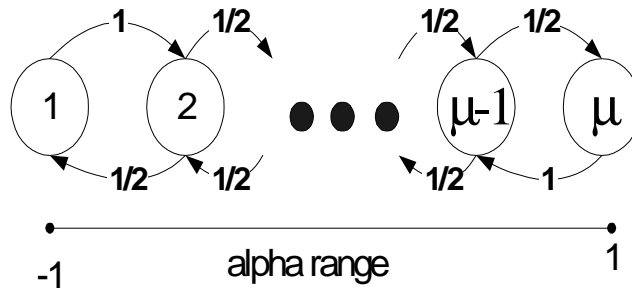


Figure 7: Markov model of random rendering of one dimension.

the second stage refines the sketch to get the final visualization. Essentially, the first stage provides the main saving with respect to the number of interactions required to find a satisfactory sketch of cluster distribution pattern, and is the dominating factor in determining how efficient the entire rendering will be. Therefore, we below focus on analyzing and improving the performance of the first stage.

Without loss of generality, we simplify the model as follows: each increase/decrease will move the  $\alpha$  to certain

fixed points, which is solely determined by the value  $\epsilon$ . For example, if  $\epsilon = 0.2$ , the serial of points would be  $\{-1, -0.8, -0.6, \dots, 0.8, 1\}$ . Suppose that there are  $\mu$  such points, including the two endpoints -1, and 1. For simplicity, we call this set of points the  $\mu$  set of points or  $\mu$  points for short.

Let  $P_j$ ,  $1 \leq j \leq \mu$  be the probability of setting  $\alpha$  to be one of the  $\mu$  points between  $[-1, 1]$ . We can model the ARR process with a Markov chain (Figure 7). It follows that  $2P_1 = P_2 = \dots = P_{\mu-1} = 2P_\mu$  [43], which implies that ARR almost uniformly sets  $\alpha_i$  to all values in  $[-1, 1]$  (except the two endpoints, which have lower probability).

Now we define the  $\alpha$  setting of the sketch visualization. Suppose that a sketch of cluster distribution can be observed with the set of  $\alpha_i$  value ranges, i.e., as long as the  $\alpha_i$  value within the corresponding range, a satisfactory cluster visualization will be observed. We model such a subrange for  $\alpha_i$  with  $\lambda_i = [\lambda_{i1}, \lambda_{i2}]$ , and  $|\lambda_i|$  as the number of the  $\mu$  points that fall into the range  $\lambda_i$ . Therefore, the probability that one ARR operation finds the satisfactory sketch can be estimated by the equation 6.

$$P = \frac{|\lambda_1|}{\mu} \cdot \frac{|\lambda_2|}{\mu} \dots \frac{|\lambda_k|}{\mu} = \frac{\prod_{i=1}^k |\lambda_i|}{\mu^k} \quad (6)$$

The above equation implies two important factors in terms of the efficiency of ARR. First, the number of effective dimensions,  $d$ , is in fact less than  $k$  varying from datasets to datasets. As the rendering rule 1 suggests, only the “dominating dimensions” are significant to rendering. In other words, the  $|\lambda_i|/\mu$  for the minor dimensions can be approximately treated as 1. Second, the individual coverage rate  $|\lambda_i|/\mu$  can be increased by reducing the effective  $\alpha$  ranges. Based on the analysis of  $\alpha$ -adjustment (recall Section 4), smaller  $\alpha_i$  values tend to hide the distribution detail over the dimension  $i$ , good for polishing, but the larger  $\alpha$  values help to distinguish visual cluster overlapping. Thus, in the ARR stage, we can choose to let ARR focus on the reduced ranges, say  $[-1, -\beta]$  and  $[\beta, 1]$  where  $0 < \beta < 1$ , for the dominating dimensions.

As observed in experiments, the rate of effective subrange  $\lambda_i$  to the reduced range is often quite large, and there are likely more than one  $\Lambda = (\lambda_1, \dots, \lambda_k)$  range combinations that can visualize the sketch of cluster distribution. Therefore, combined with the rendering rules, it is quite efficient to use ARR as the first step in rendering very high dimensional datasets.

However, ARR is not sufficient to find a detailed cluster visualization. A detailed cluster visualization might confine  $\lambda_i$ s to much smaller subranges, which requires the second stage, ADDR, to refine the sketch visualization obtained by ARR. Our experiments show that by using the combination of ARR and ADDR, the cluster visualization of census dataset (68 dimensions) can be captured in around 10 minutes.

## 5 ClusterMap Labeling

In the labeling phase of iVIBRATE, we use the adaptive ClusterMap labeling algorithm to effectively extend the intermediate clustering results to the entire large dataset. The concepts of the ClusterMap and the extended ClusterMap are discussed in the paper [5]. Thus, we only provide an overview of the ClusterMap design to

make this paper self-contained. We refer the readers to [5] for further details.

## 5.1 Encoding and Labeling Clusters with ClusterMap

ClusterMap is a convenient cluster representation derived from the VISTA cluster rendering subsystem. When visual cluster rendering produces satisfactory visualization, we can set the boundaries of a cluster by drawing a visual boundary to enclose it. Each cluster is assigned with a unique cluster identifier. After the cluster regions are marked, the entire display area can be saved (represented) as a 2D byte array (Figure 8). Each cell in the 2D array is labeled by an identifier – a cluster ID ( $>0$ ) if it is within cluster region, or the outlier ID ( $=0$ ), otherwise. Since the size of array is restricted by the screen size, we do not need a lot of space to save it. For example, the display area is only about  $688 \times 688$  pixels on  $1024 \times 768$  screen, slightly larger for higher resolution, but always bounded by a few mega pixels. As shown in Figure 8, the Cluster Map array is often a sparse matrix too, which can be stored more space-efficiently if necessary. Figure 9 is a visually defined ClusterMap of the 4D “iris” dataset. The boundaries of cluster C1, C2 and C3 were defined interactively.

In addition to the 2D array, we need also to save the mapping parameters for the labeling purpose. These parameters are listed in Table 1.

$Cmax_j, Cmin_j$	The max-min bounds of each column, $j = 1 \dots k$
$(x_0, y_0)$	The center of the visualization
$\alpha_j$	The $k$ $\alpha$ parameters, $j = 1 \dots k$
$\theta_j$	The $k$ $\theta$ parameters, $j = 1 \dots k$
$c$	The scaling factor

Table 1: Mapping parameters for ClusterMap representation

ClusterMap representation has several advantages. First, in most situations, ClusterMap provides more details than the centroid-based or representative-point-based cluster representation. Thus it is possible to preserve the precision of intermediate clustering results in the labeling phase. Secondly, the cluster boundaries can be adjusted conveniently to adapt to any special situations or to incorporate domain knowledge as we did with the VISTA system. Thirdly, with ClusterMap the outliers can be better distinguished. We shall see later that ClusterMap can also be conveniently used to adapt the extension of cluster boundary.

ClusterMap representation can be applied directly in the *basic ClusterMap labeling*. It works as follows. After the ClusterMap representation is loaded into the memory, each item in the entire large dataset is scanned and mapped onto one ClusterMap cell. The mapping follows the same mapping model used in the visual rendering system, which applies the max-min normalization followed by  $\alpha$ -mapping. Suppose that the raw large dataset is stored on disk in form of  $N$ -row by  $k$ -column table. We rewrite the formulas as follows:

$$\text{Normalization: } x'_{ij} = \omega_j * (x_{ij} - Cmin_j) - 1 \quad (7)$$

$$\omega_j = 2 / (Cmax_j - Cmin_j)$$

$$\alpha\text{-mapping: } x_i = \sum_{j=1}^k \psi_x(j) x'_{ij} - x_0, y_i = \sum_{j=1}^k \psi_y(j) x'_{ij} - y_0 \quad (8)$$

where  $\psi_x(j) = c\alpha_j \cos(\theta_j)/k$ ,  $\psi_y(j) = c\alpha_j \sin(\theta_j)/k$  and  $\omega_j$  can be pre-computed, and other parameters, such as  $c$  and  $\theta_j$  are the same as defined in the visualization model.

Sequentially, the algorithm reads the  $i$ -th item  $(x_{i1} \dots x_{ik})$  from the  $k$ -D raw dataset, normalizes and maps it with formulas (3) and (4) to a 2D cell coordinate  $(x_i, y_i)$ . From the cell  $(x_i, y_i)$ , we can find a cluster ID label or outlier label depending on the ClusterMap definition.

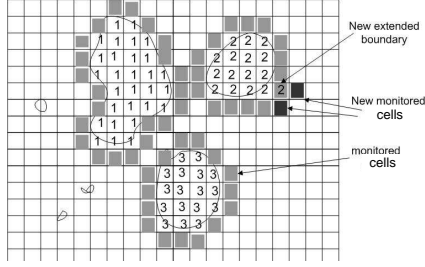


Figure 8: ClusterMap with the monitored area where  $\epsilon = 1$

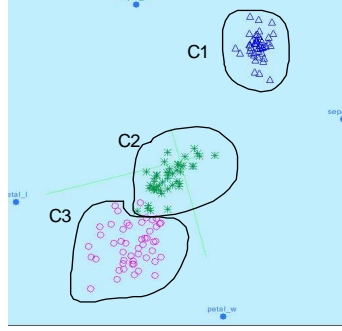


Figure 9: ClusterMap of the 4D “iris” data

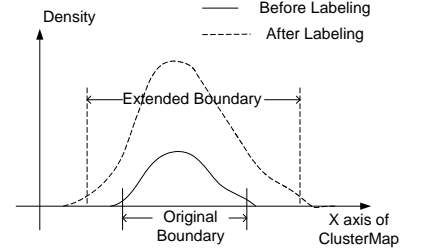


Figure 10: Boundary extension — the cross section of a typical evolving cluster in ClusterMap.

## 5.2 Adaptive ClusterMap for Boundary Extension

In the basic ClusterMap labeling, we assume that the cluster boundary defined on the sample set will not change significantly during the labeling process. However, with the increase of labeled items, the density in the original boundary areas will increase as well. Thus, the original boundary defined on sample set shall be extended to some extent. An example has been shown in Section 1.2 (Figure 2). Boundary extension encloses the nearby outliers into the clusters and may require the mergence of two nearby clusters if they become overlapped. Figure 10 sketches the possible extension in ClusterMap in terms of point density with the attending of labeled items.

Boundary extension is maintained by monitoring the point density around the boundary area. We have the initial boundary defined in ClusterMap representation. We name the cells within the cluster boundary as the “cluster cells” and the cells around the boundary area as the “boundary cells”. The initial boundary cells are precisely defined as within a short distance  $\epsilon$  away from the initial boundary. All non-cluster cells are “outlier cells” including the boundary cells. We define the density of a cell on the map as the number of items having been mapped to this cell. The density of boundary cells should be monitored in order to make decision on boundary extension. A threshold density of cell,  $\delta$ , is defined as two times of the average density of outlier cells. If the density of a boundary cell grows to  $\delta$  with the attending of labeled items, the boundary cell is turned into a cluster cell, which results in the extension of boundary. The non-cluster cells within the  $\epsilon$ -distance from the old boundary cell become the new boundary cells. Since the boundary is on the 2D cells, we can use cell as the basic distance unit and “city block” distance [50] as the distance function to define the  $\epsilon$ -distance.  $\epsilon$  is often a small number, for example, 1 or 2 blocks from the current boundary.  $\delta$  grows as the density in the non-cluster area grows, so that the measuring of boundary extension keeps consistent with the density of

non-cluster area.

To support the above adaptive algorithm, we need to extend the basic structure of ClusterMap. First, for each cell, we need one more field to indicate whether it is a monitored non-cluster cell or not. We also need to keep track of the number of points falling onto each cell. This information is saved at a “Density Map”. Since the average noise level will inevitably rise with the increase of labeled items,  $\delta$  should be periodically updated according to the average noise level. For detailed discussion of the ClusterMap algorithm, please refer to the paper [5]

The Adaptive ClusterMap labeling algorithm can be performed in two scans: The first scan generates an extended ClusterMap and the second scan can be performed to build up a R-tree index on the map for efficient access to the items on disk. After the first scan, the adjusted ClusterMap can be checked with the ClusterMap Observer to identify the anomalies (section 6.2) and to make the final refinement. The second scan is very helpful for many clustering applications that involve similarity search [36] and data indexing based on the clustering structure, which requires efficient access of the cluster members.

### 5.3 Complexity of the Labeling Algorithms

The two key factors that measure the effectiveness of the labeling algorithms are *accuracy* and *computational cost*. We have discussed the issues related to labeling outliers and irregular clusters, which can affect the accuracy of labeling algorithms greatly. The computational cost is important since the labeling algorithm has to deal with the entire large dataset. One way to estimate the cost of a labeling algorithm is to count the number of necessary multiplications. For example, one  $k$ -D Euclidean distance calculation costs  $O(k)$ . In this section, we analyze the cost of the four labeling algorithms: CBL, RPBL, Basic ClusterMap, and Adaptive ClusterMap.

Based on the formulas 7 and 8 given in section 5.1, we can roughly estimate the cost of the basic ClusterMap labeling. Map reading and parameter reading cost little constant time due to the limited small map size. For each item in the dataset, max-min normalization costs  $O(k)$  as shown by formula 7.  $\alpha$ -mapping function costs  $O(k)$  respectively to calculate the  $x, y$  coordinates with formula 8. Locating the cell in ClusterMap to get the corresponding cluster ID costs constant time. Hence, the total cost for the entire dataset is  $O(3kN)$ , where  $N$  is the number of rows in the dataset. While the Adaptive ClusterMap runs with the two scans, the cost is roughly two times of the Basic ClusterMap labeling, i.e.,  $O(6kN)$ .

When kd-tree [17] or other multi-dimensional tree is used to organize the representative points or centroids, we get the near-optimal complexity for the distance-comparison based labeling algorithms. Let  $r$  be the number of clusters and  $m$  be the number of representation points per cluster. The cost to find the nearest neighbor point in kd-tree is at least  $\log_2(rm)$  distance calculation for RPBL or  $\log_2(r)$  for CBL. For a typical RPBL as reported in the CURE paper, the number of representative points has to be greater than 10 ( $m \geq 10$ ) in order to roughly describe the regular non-spherical cluster shapes (mainly, the elongated shapes). The number should increase substantially if the irregular cluster shapes are detected. Conservatively, the cost of RPBL will be at least  $4kN$ , normally a little higher than that of the basic ClusterMap. The cost of CBL should be around  $\log_2(r)kN$  or

Algorithm	Complexity	LDS data ( $N=1M, k=5, r=5$ )	Census data ( $N=1M, k=68, r=3$ )
CBL	$[kN \log_2(r), rkN]$	4.67	56
RPBL	$[kN \log_2(rm), rmkN]$	6.5	65
Basic ClusterMap	3kN	4.42	54
Adaptive ClusterMap	$\sim 6kN$	8.69	108

Table 2: Cost estimation of the three algorithms

$rkN$  if  $r$  is small and tree structure will increase the cost.

Both CBL and RPBL need a small amount of memory,  $O(rk)$  and  $O(rmk)$ , respectively. Let  $w$  be the width and  $h$  be the height of the 2D ClusterMap, the basic ClusterMap will need  $O(wh)$  memory, which counts for several megabytes in reality. Correspondingly, the adaptive ClusterMap needs about  $O(2wh)$  memory.

Table 2 summarizes the complexities. The cost on two large datasets, LDS and Census data which will appear in the later sections, are also listed in Table 2 to give a feeling of the real cost. For both datasets,  $m$  is set to 20 and the time unit is second.

In summary, the Adaptive ClusterMap uses a little more time and space to label the datasets but this small extra cost can bring the huge benefits as we have discussed. Moreover, we are able to apply more visualization methods when integrating the ClusterMap method with the other two phases, which we will present in the following sections.

## 6 Integrating the Three Phases

Integrating the three phases of iterative clustering process (Sampling, Visual Cluster Analysis, and Adaptive ClusterMap Labeling) under the iVIBRATE framework presents an interesting and unique challenge. Since the phases are interconnected in sequence, without proper operations in the earlier phases, the later phases might be affected in terms of clustering quality and the errors could be propagated and aggravated in the later phases. In this section, we investigate two important issues in integrating the three phases. First, we study the effect of sampling phase on the later two phases, primarily the impact on determining the max-min bounds from samples (section 6.1) and exploring the small clusters hiding in outliers (section 6.3). Second, we analyze the effect of Visual Cluster Analysis on the quality of Adaptive ClusterMap algorithm, and develop the anomaly monitoring in the disk-labeling phase (section 6.2).

### 6.1 Determine the Max-min Bounds from Samples

The VISTA subsystem requires to determine the max-min bounds for normalization, denoted by “ $C_{max_j}$ ” and “ $C_{min_j}$ ”. These bounds are used not only in the rendering phase by the  $\alpha$  function but also in the labeling phase by the ClusterMap algorithms, thus these bounds should be kept unchanged throughout the three-phase clustering process. Max-min normalization is the first step in the VISTA visualization model (Section 4.1), which prepares the data for  $\alpha$ -mapping without loss of any information for visual cluster rendering. However,

since the max-min bounds are obtained from the sample set, and they may differ from the actual bounds for the entire datasets. When the bounds derived from the sample dataset does not represent a good approximation of the actual bounds of the entire dataset, such inappropriate estimation of bounds may cause additional errors in the labeling phase. Therefore, an important problem to be addressed is how to determine the proper bounds based on the sample set that are effective for both the cluster rendering phase and the Adaptive ClusterMap labeling phase.

The effect of inappropriate bounds is twofold. First, if the max and the min bounds are too tight (i.e., the two bounds are too close to one another), even though they enclose all samples, there might be high out-of-bounds rates for the entire large dataset, which increases the amount of errors generated at the labeling phase. On the other hand, if the max-min bounds are too loose, most values are scaled down to a narrow range and the difference between the values cannot be observed efficiently in cluster rendering. Recall that the  $\alpha$ -mapping in equation 2 of Section 4.1 shows two possible ways that we can adjust the visual parameters in order to observe the visual difference between different values: one is to adjust the  $\alpha_i$  values ( $1 \leq i \leq k$ ) and another is to alter the scaling factor  $c$ . Since the  $\alpha$  values ( $\alpha_1, \dots, \alpha_k$ ) are restricted in the range of  $[-1, 1]$  for the purpose of efficient interactive rendering, thus sometimes we might have to adjust the scaling factor  $c$  to a large value, which, however, could enlarge the entire visualization improperly so that some part of visualization falls out of the visualization display area.

The first problem is addressed by the relationship between the sample value bounds and the sample size — if we use the sample value bounds as the max-min bounds, analytically, how many sample points do we need in order to find the bounds that are also good for all data points? The problem of bounds estimation based on the sample data can be formalized as follows.

Let  $n$  denote the size of the sample dataset and  $p$  denote the probability of points in the entire dataset covered by the sample bounds. Since bounds estimation for each column is independent, without loss of generality, we can treat the values from one column as samples of a random variable  $X$  of size  $n$ . We now estimate the bounds for the random variable  $X$  with the sample set, so that the bounds cover  $100p$  percent of the distribution of  $X$  with certain confidence. This problem can be modeled as *Tolerance Interval* [9].

**Definition 1.** A tolerance interval  $(r \leq X \leq s)$  with tolerance coefficient  $\gamma$  is a random interval. Its range  $[r, s]$  includes at least  $100p$  percent of distribution with the probability  $\gamma$ .

In our case, we fix the two end points as the two order statistics,  $X_{(1)}$  and  $X_{(n)}$ , i.e. the max and min value of the sample, for easy processing. The above definition then can be rephrased as:

$$P[P(X_{(1)} < X < X_{(n)}) \geq p] = \gamma \quad (9)$$

Let  $F_X(x)$  be the distribution function of  $X$  and  $U_{(n)}$  and  $U_{(1)}$  be the max and min values of  $n$  uniform samples in  $[0, 1]$ .  $P(X_{(1)} < X < X_{(n)})$  is equal to  $F_X(X_{(n)}) - F_X(X_{(1)}) = U_{(n)} - U_{(1)}$ . Therefore, without knowing the distribution of  $X$ , we can find the distribution of  $U_{(n)} - U_{(1)}$  instead, which is solely related to the order statistics of uniform distribution. Let  $U = U_{(n)} - U_{(1)}$ , it is easy to find the joint distribution with order

statistics. We can get the density distribution of  $U = U_{(n)} - U_{(1)}$ ,  $f_U(u, v)$  as follows.

$$f_U(u, v) = n(n-1)u^{n-2}(1-u) \quad (10)$$

Now, since  $\gamma = P(U \geq p)$ , we can get the following relation between  $\gamma$ ,  $p$ , and  $n$ .

$$\gamma = \int_p^1 n(n-1)u^{n-2}(1-u)du \quad (11)$$

The right side of the equation is the *incomplete Beta function* (represented as `betainc(1-p, 2, n-1)` in Matlab). Fixing one of the three parameters  $\gamma$ ,  $p$ , and  $n$ , we can infer the relation between the other two parameters. We are more interested in the range of the sample size for a large  $p$  so that the sample bounds can cover almost all points in the entire dataset. By setting  $p$  to a very high probability, 0.999, we get the relationship between  $\gamma$  and sample size  $n$  as Figure 11 shows.

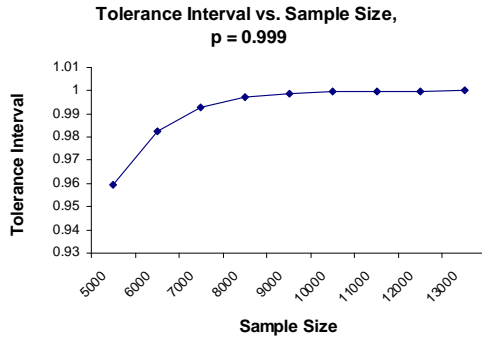


Figure 11: The relation between  $\gamma$  and  $n$ ,  $p = 0.999$

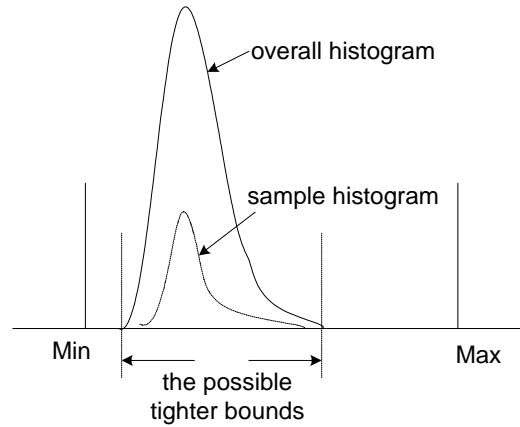


Figure 12: The possible tighter bounds for skewed distribution

At sample size  $n \approx 13,000$ , the tolerance level almost reaches 100%, which means that we can confidently say that the max-min value bounds of a sample set in size of  $n = 13,000$  or larger are the bounds which cover 99.9% records in the entire dataset. Note that the number 13,000 is induced without any assumption about the data distribution and the sample size for real datasets. For real datasets that have some special distributions, the sample size should be smaller as shown in section 8.

Our experiments with the first prototype of iVIBRATE shows that its VISTA cluster rendering subsystem can comfortably handle up to 50,000 items with 30-50 dimensions in near real time, in a common computer system environment (for example, CPU 1.5Ghz, memory 256M) [6]. Thus, the VISTA subsystem can also comfortably render a large sample set ( $n > 13,000$ ), which contains the max-min bounds for the entire dataset.

As the second problem addressed, the initial max-min bounds based on the sample value bounds might also be too wide, when the distribution is skewed as shown in Figure 12. In a skewed distribution, almost all points are located within a narrow range, with small amount of points far away from the center. This can frequently



happen in most real datasets. In this case, if we simply use the sample value bounds for normalization, the cluster rendering subsystem may not work efficiently. In order to avoid such situations, we need to check the histogram of data columns to identify the skews. If the skews are found, we may need to check the histogram for the entire dataset to carefully narrow down the bounds. With the histogram of the entire dataset, we can identify how much we can move the bounds to the center in order to make the out-of-bounds rate less than  $1 - p$ .

Based on the above analysis, we suggest the following steps to choose the normalization bounds for each column.

1. Sample the dataset to get a sample set in size of 13,000 or more;
2. Find the max-min bounds of sample set as the initial normalization bounds for each column and build the histograms for the columns with the sample set;
3. If some columns have very skewed distribution with some outliers, we build the histograms for these columns with the entire dataset. The loose bounds can be narrowed down according to the histograms for entire dataset, while maintaining the out-of-bounds rate as  $1 - p$ , e.g.,  $p=0.999$ .

The cost of choosing the proper bounds for the dimensions is quite acceptable. If no skew is found in the second step, the total cost is  $O(n)$ ,  $n$  is the sample size. Otherwise it is  $O(N)$ ,  $N$  is the size of the large dataset. Since this is a one-time process,  $O(N)$  is still not bad.

## 6.2 Monitoring the Anomalies in ClusterMap Labeling

Boundary extension can behave abnormally due to low sample rate, imprecise rendering result or inappropriate setting of initial cluster boundary. We first discuss two possible anomalies in the ClusterMap labeling process, and then introduce the methods to monitoring and handling these anomalies.

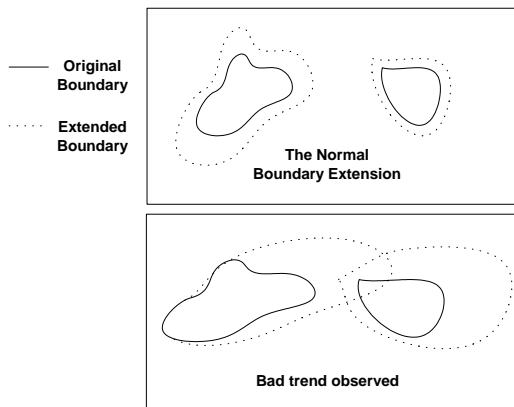


Figure 13: Anomalies that require fine adjustment of  $\alpha$  values

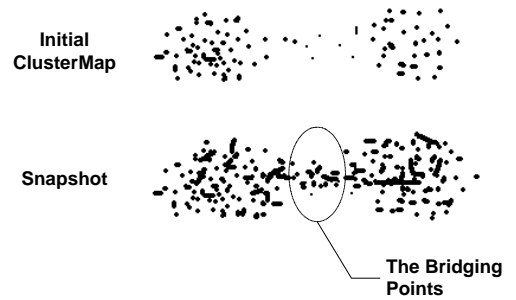


Figure 14: The bridging points

- The first anomaly is the vague cluster boundary. The cluster boundary becomes vague soon after labeling certain amount points, while normal boundary extension should be slow and happen uniformly around the boundary throughout the entire labeling process. There are two situations that can cause the anomaly, illustrated by Figures 13 and 14. First, the initial distances between the clusters are not defined appropriately due to the sample size or lack of visual refinement, which requires the user to tune the initial ClusterMap, e.g., adjusting the  $\alpha$  parameters slightly in VISTA. This is illustrated by Figure 13. Second, the other situation is the “bridging points” between the clusters, which are not dense enough in the sample set but they may form the “bridge” that connects the clusters later in the labeling process, as shown in Figure 14. The user has to make decision based on the domain knowledge to either split the clusters or merge them on the ClusterMap.
- The second anomaly is that the ordering of data records on disk may affect the boundary extension. For example, a sequence of data is mapped to a focused boundary area at early stage of labeling thus the boundary is extended in burst, but later on no more points are mapped to that area. As a result, this area is falsely extended as a part of cluster. We observe that this error only happens in the situation where such particular data records are stored together and the labeling is done according the original ordering of data records on disk. This anomaly can be avoided by accessing the data records in perturbed sequence. We use a method named “sequence perturbation” (Figure 15). To put it simply, if the large data file is regarded as a block file, we equally divide the datasets into  $s$  sequences of blocks. In each processing window, we read the data blocks at the head of each sequence and perturb the sequence of the records in the blocks. Labeling these sequence-perturbed records can reduce the risk of non-uniformity considerably in data ordering.

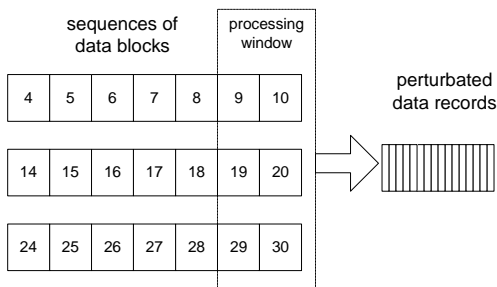


Figure 15: Record perturbation

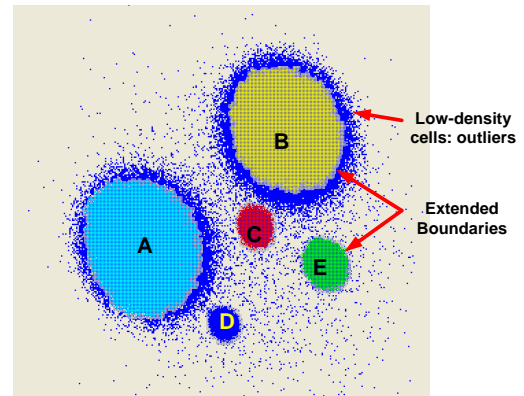


Figure 16: A snapshot of labeling a large dataset, visualized with ClusterMap Observer

In general, these anomalies can be monitored with the “snapshots” of labeling, which are visualized by the tool “ClusterMap Observer”. Snapshots are a series of evolving ClusterMaps and density maps, which incorporate the boundary extension, saved at some time interval during the first scan of Adaptive ClusterMap labeling. The user can observe the snapshots with ClusterMap Observer. If the anomalies are observed, the user can terminate

the labeling process early and returns to VISTA subsystem to adjust the original ClusterMap. Figure 16 shows a snapshot of labeling a large dataset after 10 million records are labeled. The noisy areas around cluster A and B are not labeled as cluster cells since the density of these cells does not reach the threshold. Whether these cells should be included into clusters or not, may depend on the user’s requirement. However, the extended boundary can always be edited with ClusterMap Observer, which makes the entire labeling process very flexible and manageable.

### 6.3 Detect and Explore the Small Clusters

Missing small clusters is most likely caused by low sample rate, e.g. less than 1% of the entire dataset. The small clusters may start to emerge as the labeling proceeds, which could be detected in the snapshots of labeling.

If there are small clusters emerging, we can use the following filtering method to confirm and explore the small clusters in detail. In the labeling phase, we run the Adaptive ClusterMap labeling to label all records, and then extract the outliers only from the large dataset for visual rendering. If the outlier dataset is still large, it is sampled and rendered in VISTA cluster rendering system again. Since the size of the outlier dataset is usually much smaller than that of the original dataset, one additional sampling for the outlier dataset is often sufficient to discover the small clusters in it. Similarly, the observed small clusters are marked in an *additional* ClusterMap. We can repeat this process until the size of the outliers becomes negligibly small. This process might result in a couple of additional ClusterMaps representing the small clusters at different detail levels. These ClusterMaps are used together to effectively label the interested small clusters.

In this process, the user can always control the “drill-down” level and the size of interested small clusters. More flexibly, the user can select any interested area of ClusterMap and zoom in to observe the possible small clusters in the corresponding portion of data only. This can be done iteratively, which results in a general extended iVIBRATE framework for hierarchically exploring the clusters in very large datasets (Figure 17). In short, under the iVIBRATE framework, users have more flexibility in monitoring and exploring the details in clustering structure. To our knowledge, no one of the existing approaches has provided such flexibility.

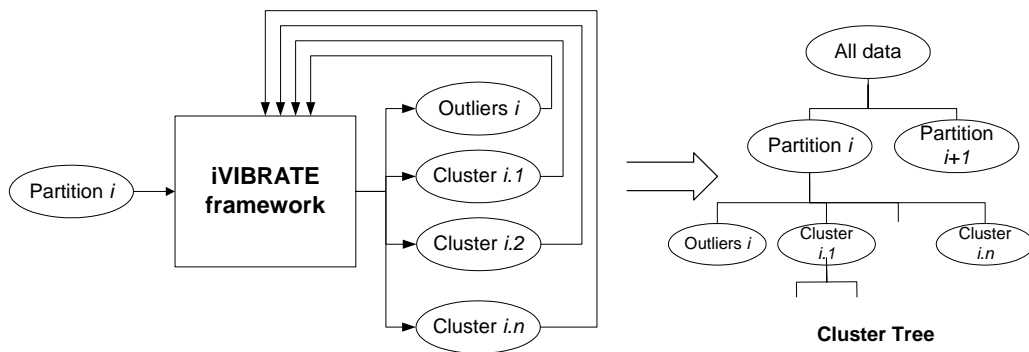


Figure 17: Iterative exploration in the extended iVIBRATE framework

## 7 EXPERIMENTS

This section presents three sets of experiments. The first set of experiments shows the effectiveness of VISTA visual clustering rendering in finding irregularly shaped clusters. The second set of experiments shows that ClusterMap can handle outliers and irregularly shaped clusters with low computational cost. The third set of experiments demonstrates the advantage of Adaptive ClusterMap in dealing with large datasets. The results show that this visualization powered framework *i*VIBRATE is more reliable and flexible in clustering very large datasets.

### 7.1 Datasets and Experiment Setup

The first set of experiments are conducted on a number of well-known datasets that can be found in UCI machine learning database <sup>1</sup>. These datasets, although small or median in size, have irregular cluster distribution, which is an important factor for testing the effectiveness of the VISTA system. We carefully choose these datasets with the following three factors in mind: 1) the current version of VISTA system only concerns the datasets having a manageable number of numerical attributes; 2) clusters in most of the datasets are not in regular spherical shape, the size of cluster may vary in a big range, and the distance between clusters can be so close that the algorithmic approaches can easily fail to distinguish; 3) the existing class labels can effectively indicate the irregular clusters. For easy comparison, we also ignore the tiny clusters in some datasets, for example, in “ecoli” data and “shuttle” data.

Dataset	$N$	# of dim.	# of clusters
Breast-w	699	10	2
Crx	690	15	2
Ecoli	336	7	8
Hepatitis	155	19	2
Ionosphere	351	34	2
Iris	151	4	3
Wine	178	12	3
Shuttle.test	14500	9	7

Table 3: The datasets used in visual rendering.

Two datasets are used for the second set of the experiments. One is the simulated dataset DS1 used in CURE [20]. DS1 is a 2D dataset having five regular clusters, including three spherical clusters, two connected elliptic clusters, and many outliers. In our experiments, DS1 is used to evaluate the effect of outliers on the labeling algorithms. The second dataset is the “shuttle” dataset (STATLOG version, test dataset) introduced in the first set of experiments. It is a 9-dimensional dataset with very irregular cluster distribution. There are seven clusters in this dataset, among which one is very large with approximately 80% of data items, and two are moderately large with approximately 15% and 5% of data items, respectively. Others are tiny clusters and thus ignored in comparison. “Shuttle” dataset is used to evaluate the effect of irregular clusters on the labeling process. These

<sup>1</sup><http://www.ics.uci.edu/~mllearn/Machine-Learning.html>

<sup>2</sup>There are totally 8 attributes in Ecoli data, but one is the name of E.Coli, which is discarded in clustering.

two datasets should show how ClusterMap avoids the common problems of the traditional algorithms.

In the third set of experiments, a simulated 5-dimension large dataset LDS with one million records is designed to test the performance of Adaptive ClusterMap on very large datasets. Figure 26 shows a 10K sample set visualized with VISTA system. LDS simulates 5 clusters – three are approximately spherical, and the other two are in irregular shape. There are also about 1% outliers. LDS is well-designed so that we can approximately predefine the control labels for entire dataset with small errors. This dataset is used to evaluate the effect of all of the three factors: outliers, irregular clusters, and boundary extension.

The three labeling algorithms, CBL, RPBL, and ClusterMap are implemented in C++. RPBL is based on the boundary points generated by the CURE clustering algorithm, which was known as a fine RPBL adapted for non-spherical cluster. We run CURE clustering to get the boundary points with the following parameters: the number of representative points is 20 and alpha (the shrink factor) is set to 0.5 as suggested. We also use ANN (Approximate Nearest Neighbor) C++ library from University of Maryland at College Park to construct *kd*-trees for RPBL and CBL in order to improve the performance of nearest neighbor search.

## 7.2 Visual Cluster Rendering

In this section we will introduce the experimental result concerning the power of visual cluster rendering system in finding clusters. The VISTA visual clustering system was implemented in Java <sup>3</sup>.

When we finish the interactive cluster rendering, we mark the cluster areas, in which the points are respectively labeled with the cluster ID. With the original labels in the datasets, we can define the items that are wrongly clustered as the errors, the number of which divided by the size of the dataset is the error rate of visual cluster rendering result.

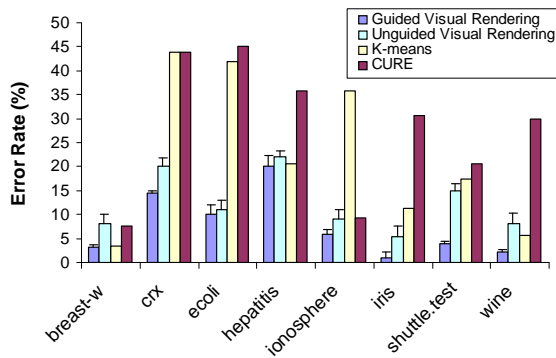


Figure 18: Comparison of error rates on the experimental data

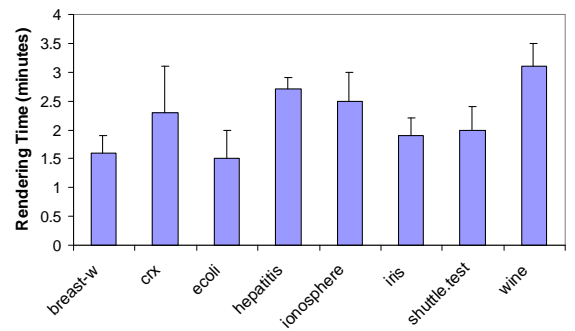


Figure 19: Visual rendering cost (GVR) on the experimental data

First, we use unguided visual rendering (UGV), where no external information is incorporated, to find the visual partition. Unguided visual rendering only depends on the visually observed dense-point areas and the gaps between the areas to discern the clusters and the cluster boundaries. Since there is visual bias on the

<sup>3</sup><http://disl.cc.gatech.edu/VISTA>

visualization, the visual rendering sometimes may trap in local minima, where some visual cluster overlaps are not distinguished. It has been shown that solely depending on visual information to define the clusters is error-prone [10, 24] due to the mapping for visualization.

We could possibly avoid the local minima by incorporating some external information, either from the result of the clustering algorithms or from the domain knowledge. This results in the second rendering method “guided visual rendering (GVR)”. In our experiment, ten of labeled items from each cluster are randomly selected serving as the “landmarks” in GVR.

We also compare the visual rendering with two algorithms, K-means and CURE algorithms to see the possible improvement. CURE clustering [20] is recognized as one that can deal with irregular cluster shapes in some degree, and K-means is the most popular algorithm. Both algorithms use the normalized data as input. The K-means results shown in Figure 18 are the best result in ten runs.

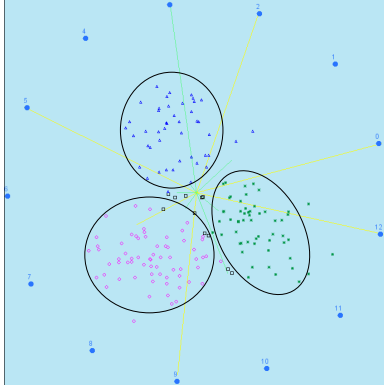


Figure 20: Visualization of “wine” data

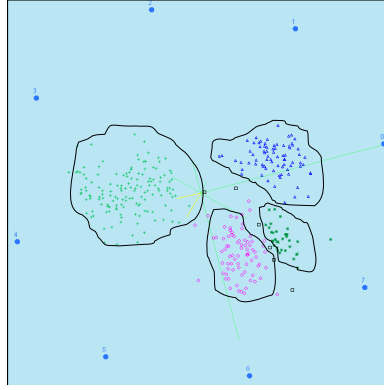


Figure 21: Visualization of “Ecoli” data

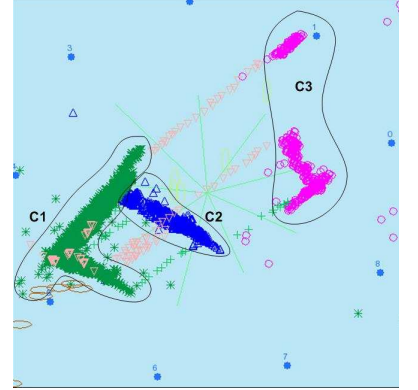


Figure 22: Visualization of “shuttle” data

The experimental result shows that neither CURE or K-means can deal irregularly shaped clusters very well, and UGV may also trap in some local minima for some datasets, for example, “breast-w (breast-cancer-wisconsin)” and “wine” (Figure 20). However, by combining the external information we can improve the UGV result to some extent. For example, “Iris” and “Ecoli” (Figure 21) datasets have very clear cluster structure thus the UGV and GVR yield almost the similar results. But in rendering “shuttle” (Figure 22), we need the help of the landmark points to distinguish the clustering structure, as the domain-specific clusters are formed by merging or splitting some irregular clusters.

In addition, the average interaction time (GVR) of five trained users (Figure 19) indicates that it is not difficult for a trained user to find a satisfactory visualization with the help of the visual rendering rules introduced in section 4.2.

### 7.3 Outliers and Irregular Cluster Shapes

We have discussed three different labeling algorithms: Representative-Point Based Labeling (RPBL), Centroid-Based Labeling (CBL), and ClusterMap (basic ClusterMap in this set of experiments). In this section we study

the performance and error rates in labeling the outliers and irregularly shaped clusters.

We run VISTA to get the ClusterMaps in map resolution of  $688 \times 688$ . The cost to rebuild a ClusterMap structure in memory is about 340~360ms, which can be ignored for processing very large datasets. Both DS1 and shuttle datasets show that the estimation of costs is appropriate – all three are linear and the basic ClusterMap is almost the fastest one (Figure 23). Linear complexity is good enough for a one-time building process.

The DS1 dataset is used to show the effect of outliers on the algorithms. The result shows that the error rate of ClusterMap is much lower than the other two algorithms. From the visualization of labeling results, we observed that CBL can suffer from both the variant cluster sizes, the distance between clusters and the outliers. Particularly, we take a look at the visualization of RPBL result (Figure 24), which clearly shows that the outliers are labeled as the members of the nearby cluster and some points in the cluster boundary are incorrectly labeled too.

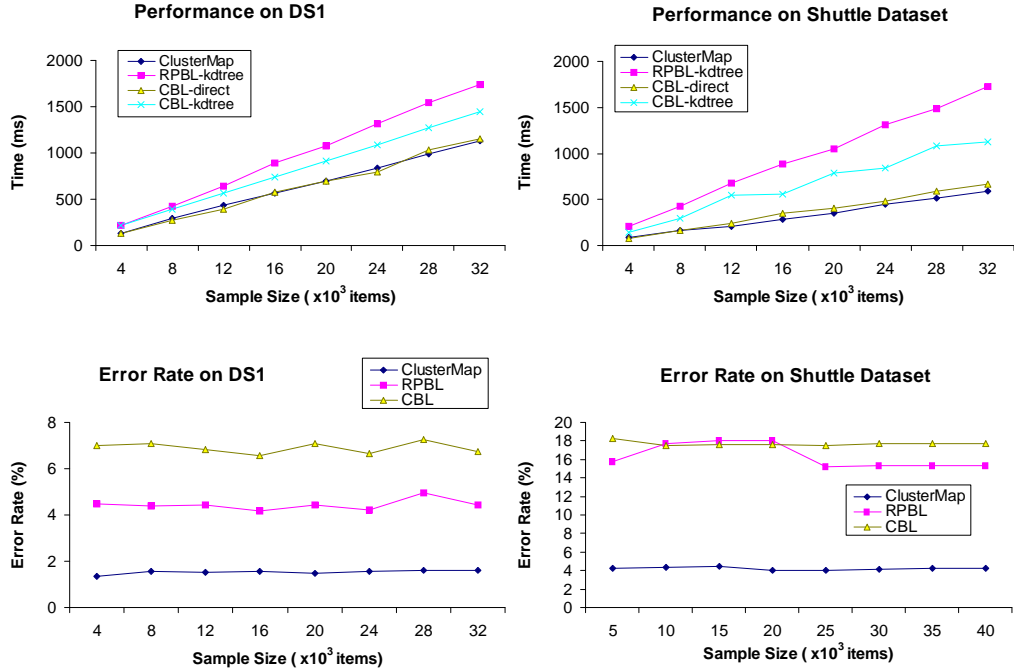


Figure 23: Labeling outliers and irregularly shaped clusters

“Shuttle” dataset has very irregular clusters. Without the incorporation of domain knowledge, the intermediate clustering will not be satisfactory. If not impossible, it is very difficult for the automated clustering algorithms to incorporate such important external information. We use the shuttle dataset to show that the error in intermediate clustering result will propagate to and might be amplified in labeling phase with the existing labeling algorithms. With the increase of labeled records, the error rate of RPBL increases to the level similar to that of CBL, due to its lack of ability dealing with the very irregular cluster shapes. The basic ClusterMap labeling keeps consistent with the VISTA cluster rendering result and thus has much lower error.

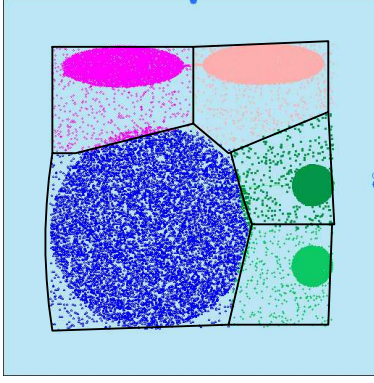


Figure 24: Outliers are labeled as the members of the nearby clusters

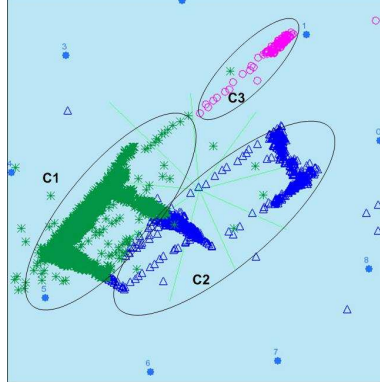


Figure 25: Visualization of the inaccurate RBPL result on Shuttle data

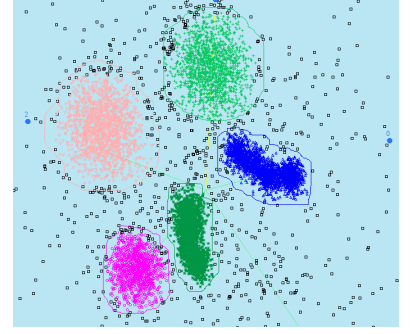


Figure 26: Visualization of 10K samples of LDS

## 7.4 Adaptive ClusterMap on Large Dataset

This experiment on the large dataset LDS mainly shows the scalability and effectiveness of Adaptive ClusterMap algorithm, especially in dealing with outliers and boundary extension.

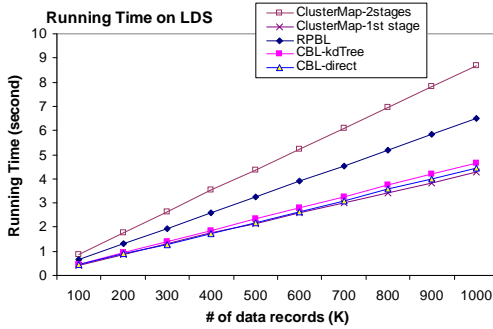


Figure 27: Performance on LDS

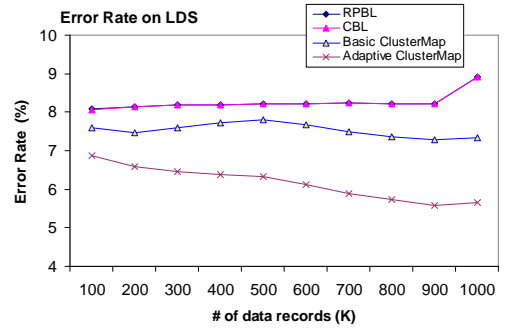


Figure 28: Error rate on LDS

The progressive plots of time cost and error rate are shown in Figure 27 and 28. The performance curve shows that the two-stage cost of Adaptive ClusterMap labeling will be a little greater than the other algorithms, however, it still keeps linear in terms of the number of labeled records. If the second stage is not necessary for some applications as we discussed, the first stage will only cost as little as CBL. Therefore, the Adaptive ClusterMap labeling is still scalable to the number of records.

From the error rate curves (Figure 28), we can observe the difference and the trend between the four algorithms. Basically, RPBL and CBL have higher error than the ClusterMap labeling algorithms due to the lack of ability dealing with outliers and imprecise boundary definition. Since the clusters are well-separated, as the number of labeled records increases, RPBL and CBL basically have the similar labeling results and thus the error rates are very close (they are overlapped in the figure). The basic ClusterMap does not consider the boundary extension, therefore, has higher error than Adaptive ClusterMap. The error of Adaptive ClusterMap tends to decrease with the increasing number of labeled records, because the more the labeled records are incorporated, the better the



Adaptive ClusterMap labeling can approximate the real cluster boundary.

## 8 Exploring Clusters in Very Large Census Data: a Comprehensive Example

In this section, we analyze the clusters in a real large dataset – the 1990 Census dataset using the iVIBRATE framework. The procedure revisits most of the concepts, methods, and models we have presented in the previous sections. Concretely, it includes max-min bounds checking and refining in the sampling phase, rendering the sample dataset with the VISTA subsystem in the visual cluster analysis phase, using adaptive ClusterMap labeling to label the entire dataset in the disk-labeling phase, and analyzing the small clusters missed by sampling. We will start with the description of the dataset,

### 8.1 Dataset

The very large “US Census 1990 Data<sup>4</sup>” is used in this empirical study to show the effectiveness and flexibility of the iVIBRATE framework. This dataset is a discretized version of the raw census data, originally used by the paper [39] in studying the relationship between the sampling approach and the effectiveness of Expectation-Maximization (EM) based clustering algorithms for very large datasets. This dataset is very large in terms of both the number of records and the number of attributes. After dropping many of the less useful attributes in the raw dataset, the total number of preserved attributes still reaches 68. It also contains more than 2 million (2,458,284) records, about 352 megabytes. Since the dataset is a discretized version, we also run an entropy-based categorical clustering/validation algorithm (ACE algorithm for finding a sub-optimal partition and BkPlot for determining the best Ks) [7] under the iVIBRATE framework to cross-validate the result.

### 8.2 Sampling and Bounds

In section 6.1, we have formally analyzed the out-of-bounds rate in terms of the sample size. Note that the analysis is good for any data distributions, which actually results in a quite conservative estimate to the appropriate sample size (around 10,000 sample records) guaranteeing the global bounds included in the sample dataset. In order to verify this observation, we try two sets of sample datasets: 5K-record and 10K-record, respectively. Each set consists of 10 sample sets, generated by uniformly sampling the entire census dataset. By checking the bounds, we find only two 5K-record sample sets having 1 and 2 attributes, respectively, missing the maximum discretized values, which cause less than 0.1% out-of-bounds rate for the entire dataset. However, all of the 10K-record datasets include the global bounds. This demonstrates that the estimated threshold in section 6.1 is indeed a conservative number. By checking the histograms, we confirm that it is unnecessary to adjust the initial max-min bounds for effective rendering. Therefore, both 5K-record and 10K-record sample datasets should be fine for effective visual rendering.

---

<sup>4</sup>In UCI KDD Archive <http://kdd.ics.uci.edu/databases/census1990/USCensus1990.html>

### 8.3 Visual Cluster Rendering

Due to the large dimensionality, manually rendering the data with the dimension-by-dimension method is not recommended for the census dataset. The new semi-automated rendering method: automatic random rendering (ARR) followed by automatic dimension-by-dimension rendering (ADDR) (section 4.3) is used in visually rendering a set of 10K-record sample census datasets in this experiment.

Recall that ARR is a randomized process, visual rendering with ARR can result in different visualizations. We conduct visual rendering 5 times to see the difference between the rendering results. By comparing the initial visualization results, the two visually well-separated clusters are confirmed as the same clusters in all of the five renderings. We list some numbers about the interactive visual rendering in Table 4. The rendering time is the wall-clock time with the unit of half minute. Due to the high dimensionality and dense cloud, the average rendering time is about 10 minutes much higher than those shown in Figure 19. Automatic random rendering (ARR) can quickly identify the sketch of cluster distribution (2 clusters) in about 2 minutes, while automatic dimension-by-dimension rendering (ADDR) should take longer time to refine the visualization. In Table 4, only the number of records in C1 and C2, are shown, the rest of the 10K records are regarded as outliers. "Shared" represents the percentage of the points in current rendering result that are also shared by the other rendering results.

No.	Random Rendering (minutes)	Automatic DDR (minutes)	C1 (pts)	C1 Shared (%)	C2 (pts)	C2 Shared (%)
1	2.0	8.0	2001	93.0%	7546	87.4%
2	1.5	9.0	2198	84.7%	7610	86.7%
3	2.0	7.0	2070	89.9%	7252	90.9%
4	2.5	8.5	2040	91.2%	7403	89.1%
5	1.5	7.0	2189	85.0%	7508	87.8%

Table 4: Rendering the 10K-record sample census dataset.

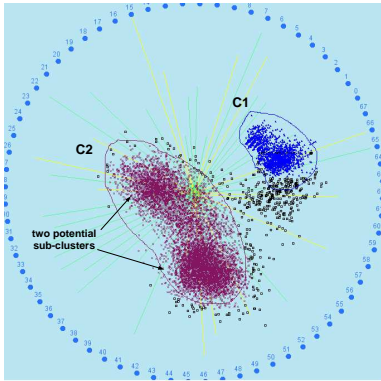


Figure 29: Result of visual rendering 1

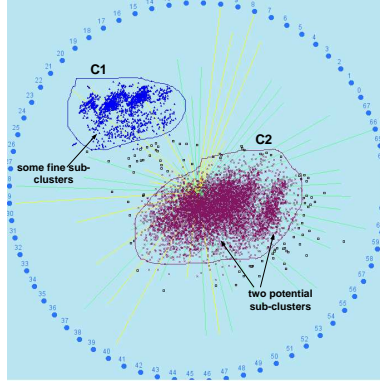


Figure 30: Result of visual rendering 2

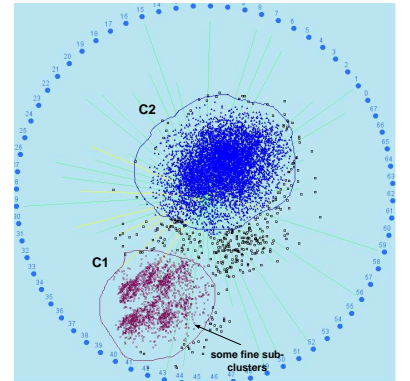


Figure 31: Result of visual rendering 3

We show three of the five visualizations in Figure 29, 30, and 31. All three clearly show the two major clusters. Figure 29 and 30 also show that C2 could potentially have two subclusters. Cluster validation with the entropy criterion [7] shows that the "best K" for clustering census dataset should be 3 or 2 (Figure 34), which implies that the initial visual clue about the possible existence of three clusters could be true. Similarly, Figure 30 and

31 show that C1 may have some fine structure, which we can also further refine if necessary.

By exploring the clusters C1 and C2 separately, we identify that C2 indeed contains two sub-clusters. We do not show the separated rendering result but the refined visualization of three clusters only (Figure 32). The result is confirmed by the entropy-based cluster validation method. Figure 33 shows that the visually separated clusters match well with the clusters labeled by the entropy-based ACE categorical clustering algorithm [7]. Note that from the algorithmic result, we are not able to identify the outliers and the initial cluster boundary, but with the VISTA system we are able to interactively define the initial cluster boundary. Figure 32 also shows the initial cluster boundary which is used to generate the ClusterMap for the labeling phase.

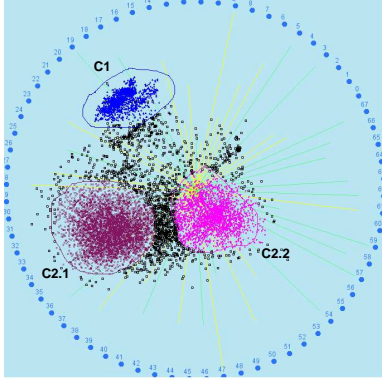


Figure 32: Final result of visual rendering with initial boundary – three clusters

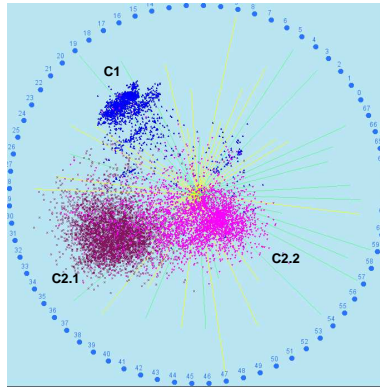


Figure 33: Visualization of ACE clustering result

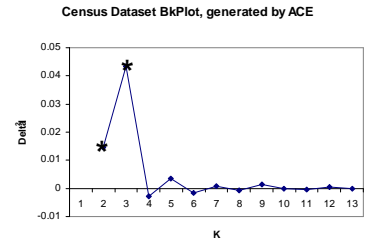


Figure 34: The “Best K” plot for census dataset

In summary, from the above procedure, we show that the combination of the automatic clustering/validation algorithms and the visual rendering method can really help to cross-validate and improve each other. Therefore, cluster analysis under the iVIBRATE framework can provide truly insightful results with higher confidence level.

## 8.4 ClusterMap Labeling

Visually clustering the census data in the iVIBRATE framework is a procedure of visual data exploration (Unguided Visual Rendering as defined earlier.), in the sense that no initial clustering clues, such as domain knowledge, are provided. For the LDS dataset in Section 7.4, we can predefine the exact cluster labels with little error and compare the labeling results with the exact cluster labels. However, because of the lack of control labels for the census dataset, we need to change the way of evaluating the labeling algorithms.

As we have initially observed, Adaptive ClusterMap labeling, together with monitoring with “ClusterMap Observer”, gives much less labeling error and the result becomes closer to the exact cluster labels with the increase of labeled records. Without the exact control labels for the census dataset, we want to see the difference between the results of other labeling algorithms and the Adaptive ClusterMap only, which should be consistent with our interpretation about the features of the labeling algorithms in Section 7.4.

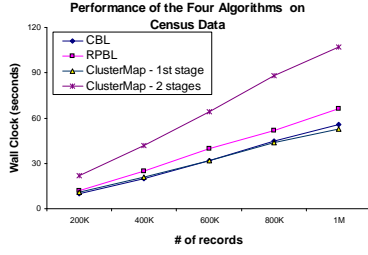


Figure 35: Cost of labeling the census dataset

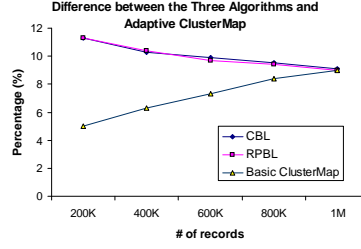


Figure 36: Difference of the labeling results

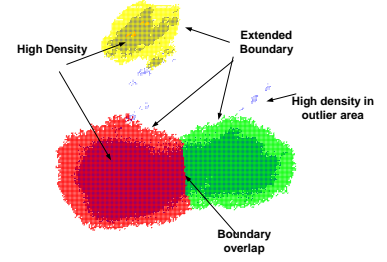


Figure 37: The ClusterMap after labeling 1M records

Figure 35 shows the performance curve of the three labeling algorithms similar to what we observed for the LDS dataset. Figure 36 shows the percentage of difference between the compared algorithms and the adaptive ClusterMap labeling. Since the three clusters are reasonably shaped and distributed, RPBL and CBL have no big difference as shown in Figure 36. The difference between RPBL/CBL and Adaptive ClusterMap tends to decrease with the increase number of labeled items, due to the automatic boundary extension including more and more points previously regarded as outliers. At the point a large number of records, for example, 1 Million records, are labeled, most of the difference should come from outliers. Similarly, the increase of difference between the basic ClusterMap and the adaptive version is solely caused by boundary extension.

Figure 37 demonstrates the cluster and boundary definition after labeling 1 million records. The boundary is evenly extended around the original boundary, therefore, the initial ClusterMap definition is very good. C2 and C3 tend to merge, but C1 is still clearly separated from C2 and C3. Whether to merge C2 and C3 or not may depend on the domain knowledge.

## 8.5 Exploring the Small Clusters

Some outliers emerge with high density as labeled in Figure 37. Therefore, it might be worthwhile to visualize the outliers only. Labeling the first million records results in 70632 records, about 7% of the total, labeled as outliers. We sample the outliers again to gain a 10K-record sample set, which is visualized as Figure 38. The initial visualization with the same parameters used in labeling clearly shows that clusters could emerge in the areas A1, A2, A3. By rendering this dataset, we find 5 clusters (the left panel in Figure 38). O1 is mapped to the boundary of the original cluster C2, therefore it should be outliers around C2. O2 and O3 mapped to the areas A2 and A3 respectively, showing the evidence of new clusters in these areas. O4 and O5 are mapped to the same area A1, which visualizes the delicate structure hiding in A1.

Observing the proximity of the small clusters to the three identified main clusters in Figure 38, we find that O4 and O5 may have close relation to the main cluster C1, and O2 and O3 seem some kind of extension of the main cluster C2.2. Merging them to the main clusters or not should depend on the domain knowledge.

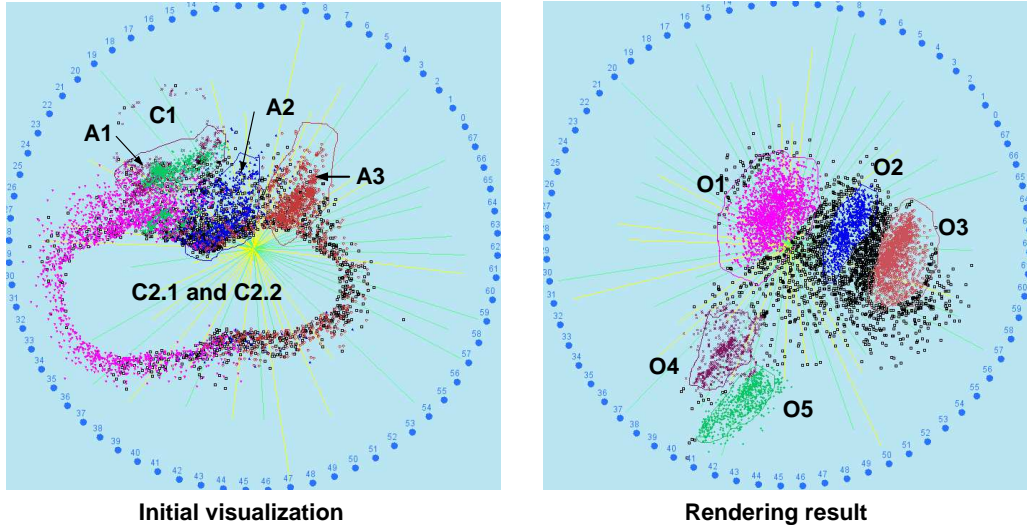


Figure 38: Exploring possible small clusters hiding in outliers

In summary, the adaptive ClusterMap labeling and the small cluster analysis with the census dataset demonstrate the effectiveness and flexibility of iVIBRATE framework in processing very large datasets. With the iVIBRATE framework, the cluster-boundary extension, outlier labeling, and small cluster detection, can be performed and monitored in a visual way, which greatly improves the precision of clustering and increases the confidence about the clustering result.

## 9 CONCLUSION

We have described iVIBRATE — an interactive-visualization based three-phase framework for clustering large datasets, focusing on two main components: VISTA visual cluster rendering system and Adaptive ClusterMap Labeling. The VISTA subsystem allows users to interactively view the potential clusters in continuously changing visualizations. More importantly, it can incorporate, validate, and refine the results of any automated clustering algorithms. It is especially effective for dealing with irregularly shaped clusters. To improve the efficiency of visual rendering, we give the heuristic rendering rules and propose a new semi-automated rendering method: automatic random rendering plus automatic dimension-by-dimension rendering for higher dimensional ( $>50D$  and  $<100D$ ) datasets. The adaptive ClusterMap labeling subsystem preserves the irregular cluster boundary defined in VISTA subsystem, clearly distinguishes the outliers and provides effective mechanisms for fine-tuning and flexible refinement of boundary extensions during the disk labeling phase. We also discussed the issues and solutions in integrating the sampling approach and the two main components into iVIBRATE framework, which help to maintain the reliability of the entire framework.

Experimental results show that, by incorporating visual cluster rendering and visualization-based disk labeling, the iVIBRATE approach not only considerably improves the quality of intermediate clustering results on the representative set with irregular shaped clusters, but also effectively extends the results to the entire large dataset with low computational cost.

The iVIBRATE framework presents a very flexible and intuitive approach for exploring clusters in large datasets, which is also an open framework for incorporating existing clustering methods with the power of visualization. We believe there will be many interesting issues in applying and extending the iVIBRATE framework in the future.

## 10 ACKNOWLEDGMENTS

This research is partially supported by NSF CNS CCR, NSF ITR, DoE SciDAC, DARPA, CERCS Research Grant, IBM Faculty Award, IBM SUR grant, and HP Equipment Grant.

## References

- [1] ANKERST, M., BREUNIG, M. M., KRIEGEL, H.-P., AND SANDER, J. OPTICS: Ordering points to identify the clustering structure. *Proc. of ACM SIGMOD Conference* (1999).
- [2] BAEZA-YATES, R., AND RIBEIRO-NETO, B. *Modern Information Retrieval*. Addison Wesley, 1999.
- [3] BARBARÁ, D., DUMOUCHEL, W., FALOUTSOS, C., HAAS, P. J., HELLERSTEIN, J. M., IOANNIDIS, Y. E., JAGADISH, H. V., JOHNSON, T., NG, R. T., POOSALA, V., ROSS, K. A., AND SEVCIK, K. C. The new jersey data reduction report. *IEEE Data Eng. Bull.* 20, 4 (1997), 3–45.
- [4] BRADLEY, P. S., FAYYAD, U. M., AND REINA, C. Scaling clustering algorithms to large databases. *Proc. of ACM SIGKDD Conference* (1998), 9–15.
- [5] CHEN, K., AND LIU, L. Clustermap: Labeling clusters in large datasets via visualization. *Proc. of ACM Conf. on Information and Knowledge Mgt. (CIKM)* (2004).
- [6] CHEN, K., AND LIU, L. VISTA: Validating and refining clusters via visualization. *Information Visualization* 3, 4 (2004).
- [7] CHEN, K., AND LIU, L. The “best k” for entropy-based categorical clustering. *Proc. of Intl. Conf. on Scientific and Statistical Database Management (SSDBM)* (2005).
- [8] CLEVELAND, W. S. Visualizing data. *AT&T Bell Laboratories* (1993).
- [9] CONOVER, W. *Practical Nonparametric Statistics*. John Wiley and Sons, 1998.
- [10] COOK, D., BUJA, A., CABRERA, J., AND HURLEY, C. Grand tour and projection pursuit. *Journal of Computational and Graphical Statistics* 23 (1995).
- [11] COX, T. F., AND COX, M. A. A. *Multidimensional Scaling*. Chapman&Hall/CRC, 2001.
- [12] CUTTING, D. R., KARGER, D. R., PEDERSEN, J. O., AND TUKEY, J. W. Scatter/gather: a cluster-based approach to browsing large document collections. *Proc. of ACM SIGIR Conference* (1992).

- [13] DHILLON, I. S., MODHA, D. S., AND SPANGLER, W. S. Visualizing class structure of multidimensional data. *the 30th Symposium on the Interface: Computing Science and Statistics* (1998).
- [14] DUBES, R. C., AND JAIN, A. K. Validity studies in clustering methodologies. *Pattern Recognition* (1979).
- [15] ESTER, M., KRIEGEL, H.-P., SANDER, J., AND XU, X. A density-based algorithm for discovering clusters in large spatial databases with noise. *Second International Conference on Knowledge Discovery and Data Mining* (1996).
- [16] FALOUTSOS, C., AND LIN, K.-I. D. FastMap: A fast algorithm for indexing, data-mining and visualization of traditional and multimedia datasets. *Proc. of ACM SIGMOD Conference* (1995).
- [17] FREIDMAN, J. H., BENTLEY, J. L., AND FINKEL, R. A. An algorithm for finding best matches in logarithmic expected time. *ACM Trans. on Mathematical Software* 3, 3 (1977).
- [18] GALLIER, J. *Geometric Methods and Applications for Computer Science and Engineering*. Springer-Verlag, 2000.
- [19] GRAY, J. What is next? a few remaining problems in information technology. *Proc. of ACM SIGMOD Conference* (1999).
- [20] GUHA, S., RASTOGI, R., AND SHIM, K. CURE: An efficient clustering algorithm for large databases. *Proc. of ACM SIGMOD Conference* (1998).
- [21] GUHA, S., RASTOGI, R., AND SHIM, K. ROCK: A robust clustering algorithm for categorical attributes. *Proc. of IEEE Intl. Conf. on Data Eng. (ICDE)* (1999).
- [22] HALKIDI, M., BATISTAKIS, Y., AND VAZIRGIANNIS, M. Cluster validity methods: Part I and II. *SIGMOD Record* 31, 2 (2002), 40–45.
- [23] HINNEBURG, A., AND KEIM, D. A. An efficient approach to clustering in large multimedia databases with noise. *Proc. of KDD98* (1998).
- [24] HINNEBURG, A., KEIM, D. A., AND WAWRYNIUK, M. Visual mining of high-dimensional data. *IEEE Computer Graphics and Applications* (1999).
- [25] HOFFMAN, P., GRINSTEIN, G., MARX, K., GROSSE, I., AND STANLEY, E. Dna visual and analytic data mining. *IEEE Visualization* (1997).
- [26] HUBER, P. J. Massive data sets workshop: The morning after. *Massive Data Set* (1996).
- [27] INSELBERG, A. Multidimensional detective. *IEEE Symposium on Information Visualization* (1997).
- [28] IWAYAMA, M., AND TOKUNGA, T. Cluster-based text categorization: a comparison of category search strategies. *Proc. of ACM SIGIR Conference* (1995).



- [29] JAIN, A. K., AND DUBES, R. C. *Algorithms for Clustering Data*. Prentice hall, 1988.
- [30] JAIN, A. K., AND DUBES, R. C. Data clustering: A review. *ACM Computing Surveys* 31 (1999).
- [31] KANDOGAN, E. Visualizing multi-dimensional clusters, trends, and outliers using star coordinates. *Proc. of ACM SIGKDD Conference* (2001).
- [32] KARYPIS, G., HAN, E.-H. S., AND KUMAR, V. Chameleon: hierarchical clustering using dynamic modeling. *IEEE Computer* 32, 8 (1999).
- [33] KEIM, D. Visual exploration of large data sets. *ACM Communication* 44, 8 (2001).
- [34] LARKIN, J., AND SIMON, H. Why a diagram is (sometimes) worth ten thousand words. *Cognitive Science* 11 (1987).
- [35] LEBLANC, J., WARD, M. O., AND WITTELS, N. Exploring n-dimensional databases. *IEEE Conference on Visualization* (1990).
- [36] LI, C., CHANG, E., GARCIA-MOLINA, H., AND WIEDERHOLD, G. Clustering for approximate similarity search in high-dimensional spaces. *IEEE Trans. on Knowledge and Data Eng.* 14, 4 (2002).
- [37] LITTLEFIELD, R. J. Using the GLYPH concept to create user-definable display formats. *National Computer Graphics Association* (1983).
- [38] LIU, H., AND MOTODA, H. *Feature Extraction, Construction and Selection: A Data Mining Perspective*. Kluwer Academic Publishers, 1998.
- [39] MEEK, C., THIESSON, B., AND HECKERMAN, D. The learning-curve sampling method applied to model-based clustering. *Journal of Machine Learning Research* 2 (2002), 397–418.
- [40] NEWMAN, M. E. The structure and function of complex networks. *SIAM Review* (2003).
- [41] PALMER, C., AND FALOUTSOS, C. Density biased sampling: An improved method for data mining and clustering. *Proc. of ACM SIGMOD Conference* (2000).
- [42] RAMASWAMY, L., GEDIK, B., AND LIU, L. A distributed approach to node clustering in decentralized peer-to-peer networks. *IEEE Transactions on Parallel and Distributed Systems (TPDS)* 16, 9 (Sept. 2005).
- [43] ROSS, S. M. *Introduction to Probability Models*. Academic Press, 2000.
- [44] ROUSSINOV, D., TOLLE, K., RAMSEY, M., AND CHEN, H. Interactive internet search through automatic clustering: an empirical study. *Proc. of ACM SIGIR Conference* (1999).
- [45] SCHUTZE, H., AND SILVERSTEIN, C. Projections for efficient document clustering. *Proc. of ACM SIGIR Conference* (1997).
- [46] SHARMA, S. *Applied Multivariate Techniques*. Wiley&Sons, 1995.



- [47] SHEIKHOESLAMI, G., CHATTERJEE, S., AND ZHANG, A. Wavecluster: A multi-resolution clustering approach for very large spatial databases. *Proc. of Very Large Databases Conference (VLDB)* (1998).
- [48] SHNEIDERMAN, B. Inventing discovery tools: Combining information visualization with data mining. *Information Visualization 1* (2002).
- [49] SILVERSTEIN, C., AND PEDERSEN, J. O. Almost-constant-time clustering of arbitrary corpus subsets. *Proc. of ACM SIGIR Conference* (1997).
- [50] SONKA, M., HLAVAC, V., AND BOYLE, R. *Image Processing, Analysis and Machine Vision*. Brooks/Cole Publishing, 1999.
- [51] VITTER, J. S. Random sampling with a reservoir. *ACM Trans. on Mathematical Software 11*, 1 (1985).
- [52] WARD, M. Xmdvtool: Integrating multiple methods for visualizing multivariate data. *IEEE Visualization* (1994).
- [53] WILLETT, P. Recent trends in hierarchic document clustering: A critical review. *Information Processing and Management 24*, 5 (1988).
- [54] XU, X., ESTER, M., KRIEGEL, H.-P., AND SANDER, J. A distribution-based clustering algorithm for mining in large spatial databases. *Proc. of IEEE Intl. Conf. on Data Eng. (ICDE)* (1998).
- [55] YANG, L. Interactive exploration of very large relational datasets through 3d dynamic projections. *Proc. of ACM SIGKDD Conference* (2000).
- [56] YANG, Y., AND PEDERSEN, J. O. A comparative study on feature selection in text categorization. *Proceedings of ICML-97, 14th International Conference on Machine Learning* (1997), 412–420.
- [57] ZAMIR, O., AND ETZIONI, O. Web document clustering: a feasibility demonstration. *Proc. of ACM SIGIR Conference* (1998).
- [58] ZHANG, T., RAMAKRISHNAN, R., AND LIVNY, M. BIRCH: An efficient data clustering method for very large databases. *Proc. of ACM SIGMOD Conference* (1996).