

FINGERPRINT: Summarizing Cluster Evolution in Dynamic Environments

Eirini Ntoutsis, Institute for Informatics, Ludwig-Maximilians University of Munich, Germany

Myra Spiliopoulou, University of Magdeburg, Germany

Yannis Theodoridis, University of Piraeus, Greece

ABSTRACT

Monitoring and interpretation of changing patterns is a task of paramount importance for data mining applications in dynamic environments. While there is much research in adapting patterns in the presence of drift or shift, there is less research on how to maintain an overview of pattern changes over time. A major challenge is summarizing changes in an effective way, so that the nature of change can be understood by the user, while the demand on resources remains low. To this end, the authors propose FINGERPRINT, an environment for the summarization of cluster evolution. Cluster changes are captured into an “evolution graph,” which is then summarized based on cluster similarity into a fingerprint of evolution by merging similar clusters. The authors propose a batch summarization method that traverses and summarizes the Evolution Graph as a whole and an incremental method that is applied during the process of cluster transition discovery. They present experiments on different data streams and discuss the space reduction and information preservation achieved by the two methods.

Keywords: *Change Detection, Change Monitoring, Change Summarization, Cluster Evolution, Cluster Summarization, Data Streams, Dynamic Environments*

INTRODUCTION

Data streams are used in many modern applications and impose new challenges for the data management systems because of their size and high degree of variability. One of the challenges is the efficient detection and monitoring of changes in the underlying population. For example, changes in the patterns known to a network intrusion detection system may indicate that intruders test new attacks and abandon old, already known (and blocked) intrusion patterns.

In general, monitoring of change is essential for applications demanding long-term prediction and pro-action.

Cluster models are commonly used as a tool for studying the dynamics of a population. In recent years actually, due to the dynamic nature of data, it has been recognized that clusters upon the data of many real applications are affected by changes in the underlying population of customer transactions, user activities, network accesses or documents. A lot of research has been devoted in adapting the clusters to the changed population. Recently, research has expanded to encompass tracing and understanding of

DOI: 10.4018/jdwm.2012070102

the changes themselves, as means of gaining insights on the population; see for example the survey of Spiliopoulou (2011) on the evolution of social networks. Understanding change is also important when taking strategic decisions: Consider, for example, a business analyst who studies customer profiles; understanding how such profiles change over time would allow for a long-term proactive portfolio design instead of reactive portfolio adaptation. While much research has been recently devoted to pattern change detection, little work has been done on the efficient maintenance of the pattern changes.

The maintenance and summarization of pattern changes upon a stream is a new problem. Summarization of data (rather than patterns); however, has been studied extensively: Popular summarization methods include histograms and wavelets, and there is much work on the efficient maintenance of these structures and on the adaptation of their contents when data change; however, these methods do not show how the data change nor do they maintain the changes themselves. There is also research on storing, modifying and querying patterns in inductive or conventional databases (e.g., Bartolini, Ciaccia, Ntoutsis, Patella, & Theodoridis, 2004); however, those approaches have not been designed for patterns over streams and, although there is provision for modifying patterns when new data arrive, there are no solutions on the efficient maintenance of changes over time. Finally, there are methods for pattern change detection (e.g., Aggarwal, 2005; Mei & Zhai, 2005; Spiliopoulou, Ntoutsis, Theodoridis, & Schult, 2006), in which different types of change can be identified and highlighted; however, the efficient long-term maintenance of the changes over an “infinite” stream is not considered.

Evolution is a permanent characteristic of streamed data, thus long-term perusal requires a space-efficient accommodation of the evolving patterns and a representation that highlights remarkable changes while suppressing trivial pattern perturbations. In this study, we propose a graph representation of pattern changes/transitions and two algorithms that condense this graph into a “fingerprint” - a structure in which

similar patterns are efficiently summarized, subject to an information loss function.

The rest of the paper is organized as follows: Related work is discussed in the upcoming section. We then present our graph model for the representation of cluster transitions. The criteria for the summarization of cluster changes and the actual summarization methods are presented afterwards. Experiments are presented in the next section. Finally, the last section concludes our work.

RELATED WORK

Relevant to our work is the work on data summarization, stream clustering and change detection. We review these areas hereafter and point out how we differentiate.

Summarization Methods

Summarization for a set of transactions with categorical attributes is studied by Chandola and Kumar (2005). In one of their methods, they derive summaries by clustering the transactions, extracting the feature/value pairs supported by all transactions in a cluster and using them as cluster summaries. They do not address the issue of cluster change upon a stream, but propose two metrics that characterize the output of the summarization algorithm, “compaction gain” and “information loss”. Quite naturally, our metrics are similarly motivated and have similar names. However, they summarize static data using clusters, while we summarize evolving clusters upon an “infinite” data stream.

Summarization and change are considered by Ipeirotis, Ntoulas, and Gravano (2005), who study changes of database content summaries. They define as “content summary” for a database a set of keywords, weighted on their importance within the database. Meta-search services use such summaries to select appropriate databases, towards which they issue keyword-based queries. The reliability of such a summary deteriorates as the contents of the database change over time. So, the authors propose methods to quantify and detect summary changes. This

study addresses both the issue of summarization over the evolving database and the discovery of changes. However, the maintenance of the summaries themselves in a condensed form is beyond the scope of their work. On the other hand, the proposed FINGERPRINT method emphasizes on the summarization of the discovered population transitions.

The discovery and representation of cluster changes for text streams are studied in Mei and Zhai (2005). They apply soft clustering with mixture models at each time period, extract the representative keyword-list (“theme”) for each cluster and then monitor the evolution of these lists by tracing divergences between a current keyword list and past ones. Theme transitions are maintained on a “theme evolution graph,” which is then used to extract the life cycle of themes (through Hidden Markov Models). The graph structure is used to reflect pattern changes, but the maintenance of this ever-growing graph is not studied and the need for summarizing it without losing information is not anticipated.

Stream Clustering Methods

Relevant to our work is the work on stream clustering. Usually, storing an entire data stream or scanning a stream multiple times is impossible due to its tremendous volume (Gama, 2010; Farnstrom, Lewis, & Elkan, 2000). To this end, several clustering algorithms have been proposed which aggregate the stream online through some appropriate summary structure and cluster these summaries offline.

This rationale was first introduced in CluStream (Aggarwal, Han, Wang, & Yu, 2003); the summary structure, called micro-cluster, is a temporal extension of the cluster feature vector of BIRCH (Zhang, Ramakrishnan, & Livny, 1996). CluStream starts with k initial micro-cluster summaries and as new points arrive, the summaries are updated such that a total of k micro-clusters is maintained at each time point. The clusters are detected offline using a modified version of k -Means over summaries instead of raw data; the user chooses the summaries to be considered by specifying the time interval. Micro-clusters can be observed

as cluster summaries and are indeed designed to reduce space demand. Nonetheless, CluStream focuses on combining them into clusters rather than in summarizing them over time. Also, the information loss affected through summarization is not discussed. The same holds for DENstream (Cao, Ester, Qian, & Zhou, 2006) and DStream (Chen & Tu, 2007), which also follow the online-offline rationale.

Recently, a method has been proposed (Al-Mula & Al Aghbari, 2011) that clusters subsequences of a data stream in order to find frequent subsequences; this method, though, refers to multiple data streams.

Change Detection Methods

Change detection methods are also relevant to our work. Aggarwal (2005) models clusters through kernel functions and changes as kernel density changes at each spatial location of the trajectory. The emphasis is on computing change velocity and finding the locations with the highest velocity - the epicenters. This model of change is very powerful, but is restricted to data over a fixed feature space. Kalnis, Mamoulis, and Bakiras (2005) propose a special type of cluster change, the moving cluster, whose contents may change while its density function remains the same during its lifetime. They find moving clusters by tracing common data records (based on their IDs) between clusters of consecutive timepoints. Yang, Parthasarathy, and Mehta (2005) detect formation and dissipation events upon clusters of spatial scientific data. Their framework supports four types of spatial object association patterns (SOAP), namely Star, Clique, Sequence, and minLink, which are used to model different interactions among spatial objects. Such methods however, assume that the feature space does not change. Thus, they cannot be used for dynamic feature spaces, e.g., text stream mining, where features are usually frequent words. Furthermore, hierarchical clustering algorithms cannot be coupled with such a method.

Cluster transition modeling and detection methods are presented in the MONIC framework of Spiliopoulou, Ntoutsi, Theodoridis,

and Schult (2006), where both changes in the data and in the feature space are anticipated. Differently from the model of Aggarwal (2005), MONIC covers changes that involve more than one cluster (external transitions), such as split and absorption, allowing insights in the whole clustering. Internal transitions, i.e., changes within a single cluster (shrink, shift, etc.), are also supported. The transition tracking mechanism of MONIC is based on the contents of the underlying data stream, thus it is independent of the clustering algorithm and of the cluster summarization/labeling method; differently from Mei and Zhai (2005). For these reasons, we use the cluster transition model of MONIC as input to our methods for the summarization of cluster changes.

BUILDING THE EVOLUTION GRAPH

We model cluster evolution across a sequence of timepoints t_1, \dots, t_n and denote as ξ_1, \dots, ξ_n the clusterings discovered at those timepoints. A clustering $\xi_i, i > 1$ may be the result of a complete re-clustering at t_i or of the adaptation of clustering ξ_{i-1} . We further denote as d_i the substream of data records seen in the interval $(t_{i-1}, t_i]$ and as D_i the substream of records, on which ξ_i is based. Depending on whether data ageing is considered or not, D_i may be equal to the set of all records seen thus far ($D_i = \cup_{j=1}^i d_j$) or to the substream seen within a time window.

The *Evolution Graph* $EG = G(V, E)$ spans the whole period of observation (n timepoints). The set of nodes V corresponds to the set of clusters seen during this time period, i.e., $V = \{\xi_1, \dots, \xi_n\}$. The set of edges E contains the cluster transitions; for each $e = (v, v') \in E$, there is a timepoint $t_i, 1 \leq i < n$ such that $v \in \xi_i, v' \in \xi_{i+1}$. By this specification of the Evolution Graph, the edges connect nodes/clusters found at adjacent timepoints. An example is depicted in Figure 1: A dotted/green edge denotes a “split” of the source cluster to multiple

target clusters. A dashed/orange edge describes an “absorption”; the source cluster is contained in the target cluster. A solid/blue edge indicates a “survival”; the source cluster has survived into the target cluster with minor changes, such as changes in size or homogeneity (Spiliopoulou, Ntoutsis, Theodoridis, & Schult, 2006). This example graph depicts three types of cluster transitions: split, absorption and survival (Spiliopoulou, Ntoutsis, Theodoridis, & Schult, 2006). In the next subsection, we describe how we assign the semantics of those transitions to the edges of the graph.

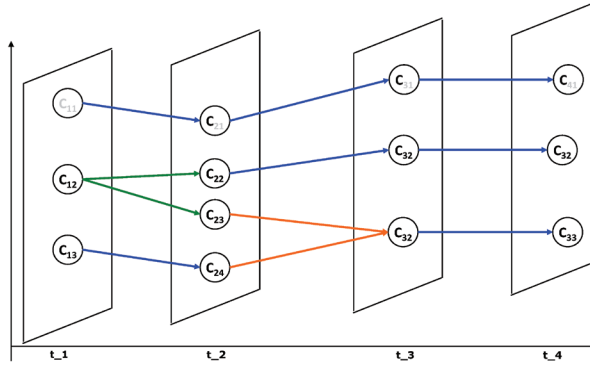
Semantics of the Graph Nodes

A node $c \in V$ represents a cluster found at timepoint t_i , i.e., belonging to clustering ξ_i . A node in the evolution graph is adorned with a “label” $c.label$ or \hat{c} , i.e., an intentional/summarized representation of its members. There are many elaborate summarized representations proposed in the literature, including micro-clusters (Aggarwal, Han, Wang, & Yu, 2003) and “droplets” over text data (Aggarwal & Yu, 2006). We opt for two simple representations, the cluster’s centroid for clusters over arbitrary numerical data and the cluster’s topic for clusters over text data.

Definition 1 (Centroid as Label). Let c be a cluster in an m -dimensional space of numerical properties. Its centroid is the vector of the mean values, $\hat{c} = \langle \mu_1, \dots, \mu_m \rangle$, where $\mu_l : 1 \leq l \leq m$ is the average of the data records’ values across the l^{th} -dimension.

Definition 2 (Keyword-based label). Let c be a cluster of text documents, where each document $d_i \in c$ is a vector in the feature space of the keywords $\{k_1, \dots, k_m\}$. The cluster label is defined as $\hat{c} = \langle w_{k_1}, \dots, w_{k_m} \rangle$, where w_{k_l} is the frequency of the l^{th} -keyword within c , if this frequency exceeds a boundary b and zero otherwise.

Figure 1. Example of an evolution graph (EG)



Semantics of the Graph Edges

An edge $e = (c, c') \in E$ denotes that a cluster $c \in \xi_i$ found at t_i has been “succeeded” by a cluster $c' \in \xi_{i+1}$ of the next timepoint. Succession means that among the clusters of ξ_{i+1} , the cluster c' is the one most similar to the cluster c . The semantics of cluster succession can be designed according to any of the approaches proposed for cluster evolution monitoring (e.g., Aggarwal, 2005; Mei & Zhai, 2005; Spiliopoulou, Ntoutsis, Theodoridis, & Schult, 2006). We have opted for the MONIC approach (Spiliopoulou, Ntoutsis, Theodoridis, & Schult, 2006) because it is independent of the clustering algorithm and can thus be used for any type of clusters, in contrary to, e.g., Aggarwal (2005). Also, it considers ageing of data which is important for streams.

In MONIC, cluster succession is based on the notions of cluster overlap and cluster matching: Let c be a cluster in clustering ξ_i at t_i and c' be a cluster in clustering ξ_{i+1} at t_{i+1} . The overlap of c and c' is defined as:

$$\text{overlap}(c, c') = \frac{|c \cap c'|}{|c|}$$

This means that the overlap depends on the members of c that are still remembered at t_{i+1} . Then, the “best match” or simply “match”

of c in ξ_{i+1} is the cluster $c' \in \xi_{i+1}$ that has the maximum overlap to c , subject to a threshold τ_{match} . If the threshold is not reached, then there is no match for c in ξ_{i+1} , i.e., c has disappeared/died. Depending on the match(es) of an old cluster among the new clusters, the old cluster may experience the following transitions:

1. *Survival*, denoted as $c \rightarrow c' : c \in \xi_i$ survives into $c' \in \xi_{i+1}$ iff c' is the match for c and there is no other cluster $z \in \xi_i$, for which c' is the match.
2. *Absorption*, denoted as $c \xrightarrow{c'} c' : c \in \xi_i$ is absorbed by $c' \in \xi_{i+1}$ iff c' is the match for c and there is at least one more cluster $z \in \xi_i$, for which c' is the match.
3. *Split*, denoted as $c \xrightarrow{c'} \{c_1, c_2, \dots, c_p\} : c \in \xi_i$ is split into $c_1, \dots, c_p \in \xi_{i+1}$, with $p > 1$, iff the overlap of c to each of these clusters exceeds a threshold τ_{split} and the overlap of all these clusters to together exceeds the match threshold τ .
4. *Disappearance*, denoted as $c \rightarrow \otimes : c \in \xi_i$ disappears if none of the above mentioned cases holds for ξ_{i+1} .

In our Evolution Graph, an edge is drawn from c to c' for each of the first three cases; if a cluster has no outgoing edges, then the fourth case has occurred. Further, we adorn the edges

with information on the transition type. In particular, let $e = (c, c') \in E$ be an edge from cluster $c \in \xi_i$ to $c' \in \xi_{i+1}$. Then, the edge e is adorned with a label $e.extTrans$ that describes the type of external transition as one of $\{survival, split, absorption\}$. If a cluster in ξ_i has no outgoing edge, it has disappeared. If a cluster in ξ_{i+1} has no ingoing edge, it has just emerged. For an emerged cluster, we form its cluster trace as follows:

Definition 3 (Cluster Trace). Let EG be an Evolution Graph captured for the timepoints t_p, \dots, t_n . For each emerged cluster c that appeared for the first time at t_i (i.e., a cluster without ingoing edge), we define its “cluster trace”, $trace(c) \equiv trace(c, t_i)$, as the sequence $\pi c_1 c_2 \dots c_m \phi$, where $c_i \equiv c$, $m \leq n - i$ and for each c_p , $i \geq 2$ there is an edge $e_i = (c_{i-1}, c_i)$ such that $e_i.extTrans = survival$. We denote the traceset of EG as T_{EG} .

For example, the traceset T_{EG} for the Evolution Graph of Figure 1 consists of the following sequences: (a) trace $\pi c_{11} c_{21} c_{31} c_{41} c_{51} \phi$, indicating that the emerged cluster c_{11} has survived across all five timepoints, (b) trace $\pi c_{22} c_{32} c_{42} c_{52} \phi$ of the emerged cluster c_{22} , one of the clusters to which c_{12} has been split and (c) the two-node traces $\pi c_{13} c_{24} \phi$ and $\pi c_{33} c_{43} \phi$. The other clusters c_{12}, c_{23}, c_{24} only existed for one timepoint and therefore built no traces.

Construction of the Evolution Graph

The Evolution Graph is built incrementally as new clusterings arrive at timepoints t_p, \dots, t_n . The pseudocode of the algorithm is depicted in Figure 2. When a new clustering ξ_i arrives at t_p , $i > 1$, our earlier algorithm MONIC (Spiliopoulou, Ntoutsis, Theodoridis, & Schult, 2006) is applied on the previous clustering ξ_{i-1} and the current one ξ_i (line 4): transitions between clusters of ξ_{i-1}, ξ_i are detected and an edge is added to the Evolution Graph for each detected transition adorned with information on the transition type (line 5). Clusters at ξ_i are also added as nodes

to the graph and labels are assigned to them (line 2). Note that MONIC uses the cluster contents (data members) for transition detection. Hence, we retain this information until the next timepoint only (line 6). So, the data members of the clusters at t_{i-1} are retained only until t_i , so as the transitions between clusters at t_p, t_{i-1} to be detected.

Algorithm Complexity

We turn now our attention to the complexity of the evolution graph construction algorithm. Adding a new clustering ξ_i into the Evolution Graph imposes the following costs: i) the computation of centroid labels for its cluster-members and ii) the computation of transitions between this clustering (ξ_i) and the clustering of the previous timepoint (ξ_{i-1}). Let $cost(c)$ be the cost of computing the centroid label of a cluster c . Then, the centroid labels computation for clustering ξ_i costs: $|\xi_i| \cdot cost(c)$. The cost of detecting transitions between the two consecutive clusterings ξ_p, ξ_i is the MONIC cost: $O(|\xi_i| \cdot |\xi_p| + |D_i| \cdot |D_p|)$ (Spiliopoulou, Ntoutsis, Theodoridis & Schult, 2006). Thus, the computational cost for adding clustering ξ_i into the evolution graph becomes:

$$addCost = O(|\xi_i| \cdot cost(c) + |\xi_i| \cdot |\xi_p| + |D_i| \cdot |D_p|)$$

If we consider an observation period of n timepoints, the building cost becomes:

$$EGbuildCost = (n-1) \cdot addCost$$

There is also the evolution graph storage cost, which refers to the space requirements for storing cluster labels, i.e., centroids, and their transitions. For a timepoint i , we should store the centroids of the clusters of the corresponding clustering ξ_i as well as the cluster transitions between this clustering and the previous timepoint clustering, t_j . If $|centroid|$ is the storage cost of a cluster centroid summary, then the storage cost for the ξ_i cluster centroid summaries is: $|\xi_i| \cdot |centroid|$. If $|edge|$ is the storage cost for an edge, then the storage cost for

Figure 2. The evolution graph (EG) construction algorithm

```

BuildEG()
Output:  $EG=G(V,E)$  //the Evolution Graph
begin
1. while a new clustering  $\xi_i$  arrives at  $t_i$  begin
2.    $EG.addNodes(\xi_i.nodes)$ ; //add  $\xi_i$  clusters in EG
3.   if ( $i==1$ ) then return EG; else  $j=i-1$ ;
4.    $E_{ji} = \text{MONICtransitions}(\xi_j, \xi_i)$ ; //detect transitions
5.    $EG.addEdges(E_{ji})$ ; //add transitions in EG
6.    $EG.updateNodes(\xi_j.nodes)$ ; //remove redundant information from  $\xi_j$ 
7. end;
8. return EG;
end

```

the transitions between ζ_i and ζ_j is:
 $|\zeta_i| * |\zeta_j| * |edge|$.

Considering the n -timepoints observation period, the storage cost becomes:

$$EG_{storageCost} = \sum_{i=1 \dots n} |\zeta_i| \cdot |centroid| + \sum_{i=1 \dots n-1, j=i+1} |\zeta_i| \cdot |\zeta_j| \cdot |edge|$$

In the above formula, the first term corresponds to the cost of storing the cluster centroid labels, whereas the second term corresponds to the cost of storing the cluster transitions.

Note that during Evolution Graph construction, we should also store the contents of the most recent clustering, ζ_j . This is because, in order to find transitions of the next incoming clustering ζ_i with respect to clustering ζ_j , the contents of ζ_j should be available. This cost is $O(|D_j|)$, where D_j are the data members of ζ_j .

QUERYING THE EVOLUTION GRAPH

The Evolution Graph contains a wealth of information regarding the evolution of the underlying population. Different queries might be imposed over the Evolution Graph so as to facilitate the end user to gain insights in the population and its evolution. We present some representative examples:

- **Forward History Queries:** *How does X evolve?*

Answer sketch: Start from X and follow its outgoing edge(s) until its descendant(s) disappear(s).

- **Backward History Queries:** *How did X emerge?*

Answer sketch: Start from X and follow its incoming edge(s) until its ancestor(s) appear(s) for the first time.

- **Comparison Queries:** *What other clusters have a similar forward history of transitions, or backward history, or both as X ?*

Answer sketch: Use the history of X as the query object and check whether the sequence of transitions that X has encompassed agrees with the sequence of transitions encompassed by some other cluster.

- **Impact Queries:** *Which clusters and at which timepoints have most influenced X into its current shape and content?*

Answer sketch: Assign an importance factor to each cluster Y participating in the history of X . Such a factor may be computed on the basis of, e.g., the overlap of Y with respect to its outgoing cluster/node on the history of X or the distance between Y and X as the number of in-between edges. Then, clusters can be ranked with respect to this importance factor.

SUMMARIZING THE EVOLUTION GRAPH TO ITS FINGERPRINT

The Evolution Graph captures the whole history of the population under observation and allows the study of cluster transitions and the inspection of cluster interrelationships. However, this graph is space consuming and highly redundant. Concretely, it contains information about each change but also contains information for clusters that did not change at all or have slightly changed. Hence, we summarize the Evolution Graph in such a way that cluster transitions are reflected but redundancies are omitted. For this, we summarize traces, i.e., sequences of cluster survivals, into “fingerprints”. These trace summaries constitute the “fingerprint” of the Evolution Graph. We first define the notion of “trace summary” and then present our algorithm FINGERPRINT in two variants, a batch and an incremental one.

Summarizing a Trace

The summarization process is applied over cluster traces (cf. Definition 3). Each trace T is traversed and the “removable” nodes are identified: These are the nodes that can be replaced by a smaller number of derived nodes, which are called “virtual centers” and are defined below.

Definition 4 (Virtual Center). Let $\langle c_1 c_2 \dots c_m \rangle$ be the trace of an emerged cluster c , $trace(c)$ and let $X = \langle c_j \dots c_{j+k} \rangle$ be a subtrace of this trace, i.e., a subsequence of adjacent nodes in the trace ($k \leq m - l$, $j \geq l$). We define the “virtual center” of X , $vcenter(X) = \hat{X}$ as a derived node composed of the averages of the labels of the nodes in X :

$$\hat{X}[i] = \frac{1}{|X|} \sum_{c_i \in X} \hat{c}[i]$$

where $[i]$ is the i^{th} dimension and \hat{c} denotes the label of cluster c . We use the notation

$c \mapsto \hat{X}$ to indicate that cluster $c \in X$ has been “mapped to” the virtual center \hat{X} .

If labels are centroid-based (Definition 1), \hat{X} is the center of the centroids of the clusters in X . If labels are keyword-based (Definition 2), \hat{X} contains the average frequencies of all frequent keywords in the clusters of X .

After introducing the virtual center as the summary of a subtrace, we define the summary of a trace: It consists of a sequence of nodes, each node being either an original cluster or a virtual center that summarizes a subtrace.

Definition 5 (Trace Summary). Let $T = \langle c_1 c_2 \dots c_m \rangle$ be a trace. A sequence $S = \langle a_j a_2 \dots a_k \rangle$ is a summary of T if and only if (a) $k \leq m$ and (b) for each $c_i \in T$ there is an $a_j \in S$ such that either $c_i = a_j$ or $c_i \mapsto a_j$, i.e., c_i belongs to a subtrace that was summarized to the virtual center a_j .

There are several possible summarizations of a trace, each one corresponding to a different partitioning of the trace into subtraces and consequently producing different virtual centers. We are interested in summarizations that achieve high space reduction while keeping information loss minimal. We generalize these objectives into functions measuring “space reduction” and “information loss,” as explained.

The replacement of a subtrace X by its virtual center \hat{X} results in storage space reduction, since less nodes are stored, but also in loss of information, since the original clusters are replaced by a “virtual center.” We model the information loss of each original cluster $c \in X$ as its distance from the virtual center \hat{X} to which it has been assigned after summarization:

$$ILoss_cluster(c, \hat{X}) = dist(\hat{c}, \hat{X}) \quad (1)$$

where $dist(\hat{c}, \hat{X})$ is the distance between the label of the original cluster \hat{c} and that of the virtual center \hat{X} .

The information loss for a cluster/node is now aggregated at the level of the trace, to

which the node belonged. The space reduction is also defined for traces.

Definition 6 (Information Loss). Let T be a trace and S be a summary of this trace. The “information loss” of T towards S is:

$$ILoss_trace(T, S) = \sum_{c \in T} ILoss_cluster(c, a_c) \quad (2)$$

where $a_c \in S$ corresponds to either the virtual center to which c is mapped after the summarization or to the cluster c itself. In the latter case, $ILoss_cluster(c, a_c) = 0$.

Definition 7 (Space Reduction). Let T be a trace and S be a summary of this trace. The “space reduction” of T towards S is the decrease in the number of nodes and edges that need to be stored:

$$\begin{aligned} SRReduction_trace(T, S) &= \frac{(|T| - |S|) + (|T| - 1 - (|S| - 1))}{|T| + |T| - 1} \\ &= \frac{2 \times (|T| - |S|)}{2 \times (|T| - 1)} \approx \frac{|T| - |S|}{|T|} \end{aligned} \quad (3)$$

where $|T|$ is the number of nodes in T and $|T| - 1$ the number of edges among its nodes (similarly for S).

This definition is similar to the definition of compaction gain in Chandola and Kumar (2005).

Next, we define the “fingerprint” of a trace as a summary, the virtual centers of which are proximal to the original cluster labels, subject to a distance upper boundary τ , so that the information loss effected through the replacement of a cluster by a virtual center is kept low.

Definition 8 (Fingerprint for a Trace). Let T be a trace and S be a summary of T . S is a “fingerprint” of T with respect to a distance threshold τ if and only if:

- (C1) For each node $c \in X$ replaced by a virtual center $a \in S$ it holds that $dist(\hat{c}, a) \leq \tau$ and,
- (C2) for each (sub)trace $\langle c_1 c_2 \dots c_k \rangle$ of T that has been summarized into a single virtual center a it holds that $\forall i = 1, \dots, k - 1 : dist(\hat{c}_i, \hat{c}_{i+1}) \leq \tau$.

By this definition, S is a fingerprint of T if it has partitioned T into subtraces of clusters that are similar to each other (condition C2) and each such subtrace has a virtual center that is close to all its original nodes (condition C1).

Once traces are summarized into fingerprints, the Evolution Graph can also be summarized, resulting in space reduction and information loss at the graph level.

Definition 9 (Fingerprint for an Evolution Graph). Let EG be an Evolution Graph and T_{EG} be its traceset. For each trace $T \in T_{EG}$, let S_T be its fingerprint (Definition 8), subject to a distance threshold τ on the distance among centroids. The set $S_{EG} := \{S_T : T \in T_{EG}\}$ is the “fingerprint of the Evolution Graph.” It effects a space reduction $SR(EG, S_{EG})$ and an information loss $IL(EG, S_{EG})$:

$$SR(EG, S_{EG}) = \sum_T SRReduction_trace(T, S_T) \quad (4)$$

$$IL(EG, S_{EG}) = \sum_T ILoss_trace(T, S_T) \quad (5)$$

We next present the two variants of our algorithm FINGERPRINT. The variant batchFINGERPRINT creates the fingerprint of an Evolution Graph by partitioning traces in such a way that their fingerprints can be built. This variant requires that the Evolution Graph is first constructed and stored as a whole. Then, we present the online variant incFINGERPRINT, which builds the fingerprints of the traces incrementally as new cluster transitions are detected, i.e., without requiring the construction of the Evolution Graph first.

Batch Summarization of the Graph

The batchFINGERPRINT summarizes an Evolution Graph EG by identifying its traces, building a fingerprint for each trace and substituting the traces in EG with their fingerprints. The batchFINGERPRINT satisfies the two conditions of Definition 8 by applying two heuristics on each (sub)trace T :

- *Heuristic A*: If T contains adjacent nodes that are in larger distance from each other than τ , then the pair of adjacent nodes c , c' with the maximum distance is detected and T is then partitioned into T_1 , T_2 so that c is the last node of T_1 and c' is the first node of T_2 .
- *Heuristic B*: If T satisfies condition $C2$ but contains nodes that are in larger distance from the virtual center than τ , then T is split as follows: The node c that has the maximum distance from $vcenter(T)$ is detected and T is partitioned into T_1 , T_2 so that c is the last node of T_1 and its successor c' is the first node of T_2 .

Heuristic A deals with violations of condition $C2$ and *Heuristic B* deals with violations of condition $C1$ for (sub)traces that already satisfy $C2$. We show the algorithm in Figure 3.

The batchFINGERPRINT creates a fingerprint of the Evolution Graph by traversing the graph, extracting its traces (line 1, condition $C2$) and summarizing each of them (line 4). The “produced” fingerprints of the traces are added to the fingerprint graph FEG (line 5). This operation encapsulates the attachment of a summarized trace to the graph by redirecting the ingoing/outgoing edges of the original trace towards the ends of the summarized trace.

The batchFINGERPRINT invokes *summarize_HeuristicA()* which recursively splits the trace into subtraces according to *Heuristic A* until $C2$ is satisfied. If the trace consists of only one node, then this node is returned (line 1). Otherwise, we test whether the trace contains

nodes whose labels are further off each other than the threshold τ (line 2). If $C2$ is satisfied, then *summarize_HeuristicB()* is invoked (line 8): It checks for condition $C1$ and returns the fingerprint of the (sub)trace input to it. If $C2$ is violated, the trace is partitioned according to *Heuristic A* (lines 3,4) and *summarize_HeuristicA()* is invoked for each partition (lines 5, 6). Finally, the summarized (sub)traces are concatenated (line 7) and returned. This concatenation operation restores or redirects the edges across which the split (line 4) was performed

The recursive function *summarize_HeuristicB()* operates similarly. It takes as input a (sub)trace T that has more than one nodes and satisfies condition $C2$. It builds the virtual center for T according to Definition 4. It then checks condition $C1$ by comparing the distance of the virtual center from each node to τ (line 2). If τ is not exceeded, the virtual center is returned (line 3). Otherwise, T is split at the node that has the highest distance from the virtual center, according to *Heuristic B* (lines 5, 6). The *summarize_HeuristicB()* is invoked for each partition (lines 7, 8). The returned fingerprints are concatenated into the fingerprint of T .

Incremental Summarization of the Graph

The batch summarization algorithm of Figure 3 requires as input the complete Evolution Graph, before building its fingerprint. This is resource-intensive, since the graph is growing continuously. We have therefore designed incFINGERPRINT, a fingerprint construction algorithm that summarizes the traces incrementally and does not require the a priori construction of the Evolution Graph. We show incFINGERPRINT in Figure 4.

Our incFINGERPRINT invokes (in line 1) our earlier algorithm MONIC (Spiliopoulou, Ntoutsis, Theodoridis, & Schult, 2006): MONIC compares the current clustering ζ (timepoint t_i) to the most recent one ζ (timepoint t_{i-1}), identifies the cluster transitions and returns them as a set E_i of labeled edges, according to

Figure 3. Offline summarization of the evolution graph with batchFINGERPRINT

```

batchFINGERPRINT(EG)
Input: the Evolution Graph EG
Output: FEG, a fingerprint of EG
1. traverse the EG and extract its traces into T;
2. FEG =  $\emptyset$ ;
3. for each trace  $T \in \mathcal{T}$  do
4.   FT = summarize_HeuristicA(T);
5.   FEG.addTrace(FT);
6. end-for

summarize_HeuristicA(T)
Input: a trace T
Output: a fingerprint of the trace
1. if  $|T| == 1$  then return T;
2. if C.b is not satisfied then
3.   find  $c \in T$  such that
      $\forall (y, z) \in T_{.1} : \text{dist}(y, z) < \text{dist}(c, c_{next})$  and
      $\forall (y, z) \in T_{.2} : \text{dist}(y, z) < \text{dist}(c, c_{next})$ ;
4.   split T into  $T_{.1} = \langle c_1, \dots, c \rangle$  and  $T_{.2} = \langle c_{next}, \dots, c_k \rangle$ ;
5.   FT_{.1} = summarize_HeuristicA(T_{.1});
6.   FT_{.2} = summarize_HeuristicA(T_{.2});
7.   return  $\langle \text{FT}_{.1} \cdot \text{FT}_{.2} \rangle$ ;
8. else return summarize_HeuristicB(T);
9. endif

summarize_HeuristicB(T)
Input: a trace T
Output: a fingerprint of the trace
1.  $v = \text{vcenter}(T)$ ;
2. if  $\forall y \in T : \text{dist}(y, v) < \tau$  then return v; // Condition C1
3. else
4.   find  $c \in T$  such that  $\text{dist}(c, v) = \max\{\text{dist}(y, v) | y \in T\}$ ;
5.   split T into  $T_{.1} = \langle c_1, \dots, c \rangle$  and  $T_{.2} = \langle c_{next}, \dots, c_k \rangle$ ;
6.   FT_{.1} = summarize_HeuristicB(T_{.1});
7.   FT_{.2} = summarize_HeuristicB(T_{.2});
8.   return  $\langle \text{FT}_{.1} \cdot \text{FT}_{.2} \rangle$ ;

```

previous subsection. The source of each edge corresponds to a node that is already in the fingerprint graph *FEG*. It is stressed that MONIC operates on the clusterings rather than the cluster labels retained in the nodes of the fingerprint graph. So, from line 2 on, incFINGERPRINT transfers information about the detected transitions in the *FEG*, summarizing survivals wherever possible. The result is an already summarized version of the Evolution Graph.

For each edge $e = (x, y)$, incFINGERPRINT examines whether e is a survival transition (line 3), i.e., whether e is part of a trace. If not, *FEG* is expanded by adding the cluster y and the edge e (lines 4, 5). We do not add the whole cluster;

we only retain its label (cf. “Change detection methods”).

If $e = (x, y)$ does belong to a trace, incFINGERPRINT checks whether the labels of x and y are similar to each other, according to condition C2 of Definition 8 (line 6). Since cluster x has already been added to *FEG*, we access its label $x.label$ directly, while the label of cluster y must be computed as \hat{y} . If condition C2 is not satisfied, the *FEG* is expanded by y and e as before. If finally, C2 is satisfied, then y and e do not need to be added to *FEG*. Instead, x and y are summarized into their virtual center v (line 10) and the node x is replaced by v (line 11). This means that all edges pointing to x are redirected to v .

Figure 4. Online summarization of the evolution graph with incFINGERPRINT

```

incFINGERPRINT(FEG)
Input: FEG // the fingerprint built so far
       $\zeta$  // the most recent clustering, build at timepoint  $t_{i-1}$ 
       $\xi$  // the current clustering, build at the current timepoint  $t_i$ 
Output: FEG // the updated fingerprint
1.  $E_i = \text{MONICtransitions}(\zeta, \xi)$ ;
2. for each edge  $e = (x, y) \in E_i$  do
3.   if  $e.\text{extTrans} \neq \text{"survival"}$  then
4.     FEG.addNode(y);
5.     FEG.addEdge(e);
6.   else if  $\text{dist}(x.\text{label}, \hat{y}) \geq \tau$  then // C2 is violated
7.     FEG.addNode(y);
8.     FEG.addEdge(e);
9.   else
10.     $v = \text{vcenter}(x, y)$ ;
11.    FEG.replaceNode(x, v);
12.   endif
13. end-for
14. return FEG;

```

The incFINGERPRINT not need to check for condition *C1*, since the distance of the virtual center of two nodes is less than the distance between the two nodes as a whole; the latter is less than τ by virtue of condition *C2*. This algorithm operates locally, treating pairs of adjacent nodes only, instead of whole traces. Thus, it has the advantage of not requiring the a priori construction of the Evolution Graph.

EXPERIMENTAL RESULTS

The goal of our experiments is to measure the space reduction and information loss for the two different summarization techniques and for different values of the centroid similarity threshold τ (cf. Definition 8) that governs the summarization process.

Datasets

We experimented with two numerical datasets, the *Network Intrusion dataset* and the *Charitable Donation dataset*, and with one text dataset, the *ACM H2.8 dataset*. The two numerical datasets have been used often in stream experiments (e.g., Aggarwal, Han, Wang,

& Yu, 2003; Cao, Ester, Qian, & Zhou, 2006). The text dataset has been used in the experiments of MONIC (Spiliopoulou, Ntoutsis, Theodoridis, & Schult, 2006). With respect to their evolving nature, the first dataset is rapidly evolving, the second one is relatively stable, while the third one evolves in an unbalanced way - one of the classes grows faster than the others.

The *Network Intrusion dataset (KDD Cup '99)* contains TCP connection logs from two weeks of LAN network traffic (424,021 records). Each record corresponds to a normal connection or an attack. The attacks fall into four main categories: DOS, R2L, U2R, and PROBING. So, we set the number of clusters to 5, including the class of normal connections. We used all 34 continuous attributes for clustering and removed one outlier point, as in Aggarwal, Han, Wang, and Yu (2003). We turned the dataset into a stream by sorting on the data input order. We assumed a uniform flow in speed of 2,000 instances per time period. For data ageing, we assumed a sliding window of 2 time periods/ timepoints.

The *Charitable Donation dataset (KDD Cup '98)* contains information (95,412 records) on people who have made charitable donations

in response to direct mailings. Clustering identifies groups of donors with similar donation behavior. Similar to Farnstrom, Lewis, and Elkan (2000), we used 56 out of the 481 fields and set the number of clusters to 10. As with the previous dataset, we used the data input order for streaming and assumed a uniform flow with 200 instances per time period. For data ageing, we used a sliding window of size 2.

The *ACM H2.8 document set* is the set of documents inserted between 1997 and 2004 in the ACM Digital Library, category H2.8 on "Database Applications." This dataset contains publications on (1) data mining, (2) spatial databases, (3) image databases, (4) statistical databases and (5) scientific databases. It further contains (6) uncategorized documents, i.e., those assigned in the parent class "database applications [only]." The subarchive consists of the documents whose primary (or secondary) class is one of these 6 classes (4,920 records). It evolves in an unbalanced way (Spiliopoulou, Ntoutsi, Theodoridis, & Schult, 2006): The category (1) is larger than all the others together and grows faster than the others. For the experiments in Spiliopoulou, Ntoutsi, Theodoridis, and Schult (2006), only the title and a list of keywords were considered for each document. We have used the same vectors and the same clustering algorithm, bisecting K-means, for $K=10$. We turned the data into a stream by using the publication date for the ordering of the records. We considered $n=7$ timepoints corresponding to the 7 publication years from 1998¹ to 2004; the corresponding data batches have different sizes varying from 837 in 1998 to 617 in 2004. The size of the sliding window was set to 2. A cluster label consists of the terms appearing in more than 60% of the cluster's vectors.

Example Traces and Fingerprints

To highlight the behavior of the summarization algorithms, we depict here some traces from the *ACM H2.8 dataset* and their corresponding fingerprints in Table 1, as produced by our summarization algorithms.

In 1998, we observe a new cluster with the label "*information system*" which is further developed in the next two years, 1999 and 2000. Its trace is:

$$trace(c_{1998_2}) = \prec c_{1998_2}, c_{1999_6}, c_{2000_3} \succ$$

The notation c_{y_i} refers to the i^{th} cluster of year y , $i = 1 \dots 9$ (cluster 0 is the garbage cluster)². At all time points, the cluster centroids contain the terms "*information*" and "*system*" but with different frequencies, as follows:

$$\hat{c}_{1998_2} = \langle information(0.96); system(0.61) \rangle$$

$$\hat{c}_{1999_6} = \langle information(0.88); system(0.74) \rangle$$

$$\hat{c}_{2000_3} = \langle information(0.76); system(0.78) \rangle$$

Both summarization algorithms condense this trace into a single virtual center. The batch algorithm creates this new node v in one step: $\hat{v} = \langle information(0.87); system(0.71) \rangle$, whereas the incremental algorithm first summarizes \hat{c}_{1998_2} and \hat{c}_{1999_6} into a virtual center: $\hat{v}_0 = \langle information(0.92); system(0.68) \rangle$, and then summarizes \hat{v}_0 and \hat{c}_{2000_3} into a new virtual center: $\hat{v}' = \langle information(0.84); system(0.73) \rangle$.

A further cluster that emerged in 1998 had the one-term label $\langle analysis(1.0) \rangle$. In 1999, it was split into two clusters, one labeled $\langle mining(1.0); datum(0.74) \rangle$ and one cluster with no label (garbage cluster). The former survived for two periods, thus resulting in the trace:

$$\prec c_{1999_8}, c_{2000_4}, c_{2001_6} \succ.$$

The information delivered without summarization is:

$$c_{1998_9} \rightarrow \{c_{1999_4}, \prec c_{1999_8}, c_{2000_4}, c_{2001_6} \succ\}.$$

Table 1. Example clusters and traces in the ACM H2.8 dataset

1998	1999	2000	2001	2002
information(0.96); system(0.61)	information(0.88) system(0.74)	information(0.76) system(0.78)		
		image(0.9) retrieval(0.72)	image(0.9) retrieval(0.71)	image(0.9) retrieval(0.69)
	---	knowledge(0.89) discovery(0.85) datum(0.62)	knowledge(0.91) discovery(1.0) datum(0.63)	
analysis(1.0)	mining(1.0) datum(0.74)	mining(1.0) datum(0.81)	mining(1.0) datum(0.84)	

Instead, the summarization delivers the fingerprint $c_{1998_0} \xrightarrow{\subset} \{c_{1999_4}, \hat{v}\}$, where \hat{v} is the summary of the trace $\prec c_{1999_8}, c_{2000_4}, c_{2001_6} \succ$.

Although the dataset is small (it spans only 7 time points), the above examples clearly display the advantages of cluster monitoring and summarization. Due to the re-occurrence of clusters (corresponding to real world concepts like information systems, image retrieval, etc.), the fingerprint produces an efficient and effective dataset overview over time. Such an overview is much more “readable” by the end user.

Space Reduction and Information Loss

We evaluated the space reduction achieved by the batch and the incremental summarization methods for different values of the centroid

distance threshold τ . The results for the two numerical datasets are depicted in Figure 5.

As expected, for both incFINGERPRINT and batchFINGERPRINT, the space reduction increases for larger values of τ , because less proximal centroids can be merged. The two algorithms achieve similar space savings although batchFINGERPRINT shows slightly higher values for most values of τ , especially in the *Network Intrusion dataset*.

The total space reduction for each dataset depends of course on the number of survivals per se: Among the total of 1,195 clusters/nodes generated for the Network Intrusion dataset, only 400 nodes participate in traces (~33%); the space reduction values achieved by both algorithms are in the range [21%, 33%] of the total size of the Evolution Graph. The Evolution Graph of the Charitable Donation dataset contained 4,770 clusters, of which 614 were involved in traces (~13%); the space reduction

Figure 5. Impact of threshold τ on space reduction for the Network Intrusion dataset (left) and the Charitable Donation dataset (right)

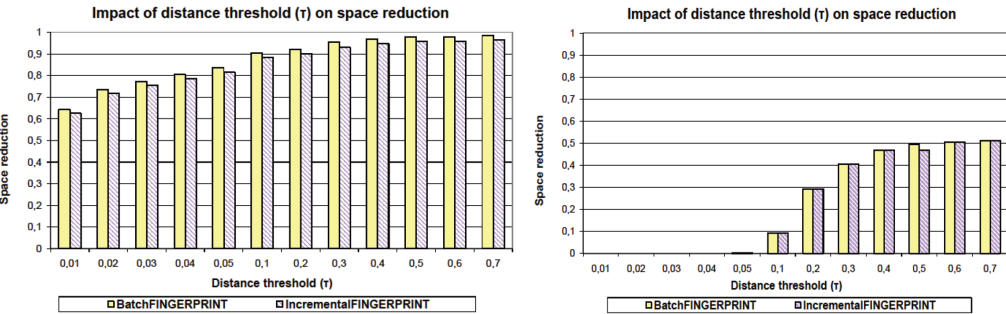
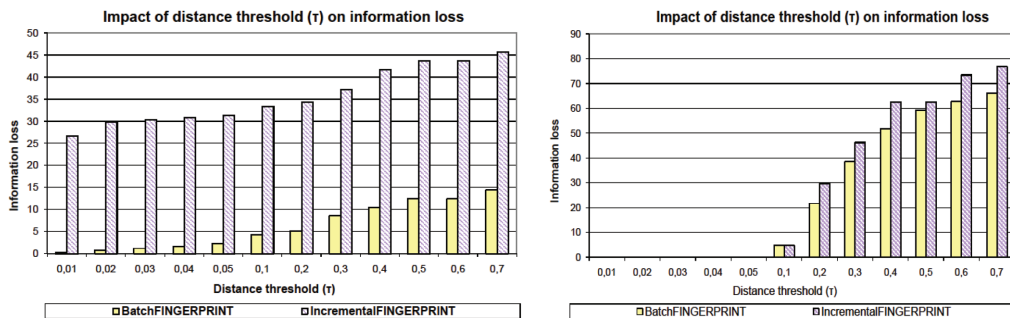


Figure 6. Impact of threshold τ on information loss for the Network Intrusion dataset (left) and the Charitable Donation dataset (right)



over the whole graph was thus no more than 7%. For the ACM H2.8 sub-archive, 24 out of 70 nodes were involved in traces (~34%), so that the space reduction over the whole graph ranged between 9% and 33%.

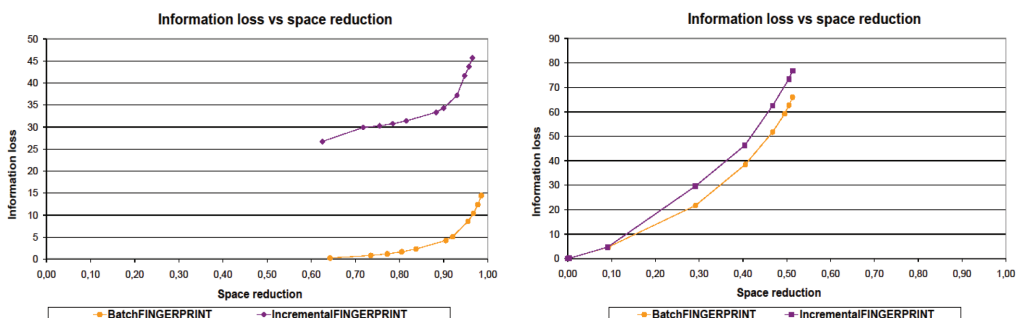
In Figure 6, we depict the information loss affected upon the datasets when summarizing incrementally versus in batch. The information loss increases with τ , since a larger τ implies that less similar centroids can be merged.

For the Charitable Donation dataset, the information loss incurred by the incremental algorithm is slightly higher than for the batch algorithm but follows the same curve for different values of τ . For the Network Intrusion dataset, the performance difference is dramatic: While the batch algorithm achieves a

very low information loss, the incremental algorithm performs very poorly. A possible explanation for the poor performance of incFINGERPRINT in the Network Intrusion dataset is the volatility of the dataset: It is likely that the survived clusters were unstable and not very similar to each other. Hence, incFINGERPRINT produced virtual centers that were not very close to the original pairs of centroids, while batchFINGERPRINT managed to build better virtual centers among multiple adjacent centroids.

In Figure 7, we show the joint curves of space reduction and information loss for the two datasets and for different values of the centroid distance threshold τ .

Figure 7. Correlation between information loss and space reduction for different values of τ for the Netw Intrusion dataset (left) and the Charitable Donation dataset (right)



CONCLUSION AND OUTLOOK

We have studied the effective summarization of cluster changes over an evolving stream of data. We modeled cluster transitions in a graph structure, the Evolution Graph, and proposed two algorithms that summarize it into a “fingerprint.” A fingerprint is a condensed representation, in which less informative cluster transitions are suppressed. We derived functions that measure the effected information loss and space reduction and we presented heuristics that drive the summarization process. One of our algorithms summarizes the Evolution Graph as a whole, while the other creates the graph’s fingerprint incrementally, during the process of cluster transition discovery. We have run experiments on three real datasets and have seen that incFINGERPRINT achieves similar space reduction to batchFINGERPRINT, but the information loss may be much higher depending on the volatility of the dataset.

The batch algorithm batchFINGERPRINT shows better performance comparing to the incFINGERPRINT algorithm, but it requires the whole dataset of transitions as an input. A hybrid summarization algorithm using both an online and an offline component is worth pursuing. Another interesting direction is modeling and investigation of the impact of the quality and stability of the original clustering on the summarization process. In this work, we have concentrated on the summarization of cluster survivals. A survival is the transition of a cluster to a similar successor cluster. Instead of placing constraints on the similarity among clusters, we want to study models of information loss for the summarization of arbitrary cluster transitions, so that only the most informative changes are delivered to the end-user.

ACKNOWLEDGMENTS

Eirini Ntoutsis is supported by an Alexander von Humboldt Foundation fellowship for postdocs (<http://www.humboldt-foundation.de/>). Part of the work of this author was done while with the University of Piraeus, Greece supported by the

Heracleitos program co-funded by the European Social Fund and national resources (Operational Program for Educational and Vocational Training II – EPEAEK II). The source of our implementations can be downloaded from <http://infolab.cs.unipi.gr/people/ntoutsis/fingerprint/>.

REFERENCES

- Aggarwal, C. C. (2005). On change diagnosis in evolving data streams. *IEEE Transactions on Knowledge and Data Engineering*, 17(5), 587–600. doi:10.1109/TKDE.2005.78
- Aggarwal, C. C., Han, J., Wang, J., & Yu, P. (2003). A framework for clustering evolving data streams. In *Proceedings of the 29th International Conference on Very Large Data Bases* (Vol. 29, pp. 81-92). New York, NY: ACM.
- Aggarwal, C. C., & Yu, P. (2006). A framework for clustering massive text and categorical data streams. In *Proceedings of the Sixth SIAM International Conference on Data Mining* (Vol. 124, pp. 479-483). Philadelphia, PA: Society for Industrial Mathematics.
- Al-Mulla, R., & Al Aghbari, Z. (2011). Incremental algorithm for discovering frequent subsequences in multiple data streams. *International Journal of Data Warehousing and Mining*, 7(4), 1–20. doi:10.4018/jdwm.2011100101
- Bartolini, I., Ciaccia, P., Ntoutsis, I., Patella, M., & Theodoridis, Y. (2004). A unified and flexible framework for comparing simple and complex patterns. In *Proceedings of the 8th European Conference on Principles and Practice of Knowledge Discovery in Databases* (pp. 496-499). New York, NY: Springer.
- Cao, F., Ester, M., Qian, W., & Zhou, A. (2006). Density-based clustering over an evolving data stream with noise. In *Proceedings of the Sixth SIAM International Conference on Data Mining* (Vol. 124, pp. 479-483). Philadelphia, PA: Society for Industrial Mathematics.
- Chandola, V., & Kumar, V. (2005). Summarization – compressing data into an informative representation. In *Proceedings of the Fifth IEEE International Conference on Data Mining* (pp. 98-105). Washington, DC: IEEE Computer Society.
- Chen, Y., & Tu, L. (2007). Density-based clustering for real-time stream data. In *Proceedings of the 13th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 133-142). New York, NY: ACM.

- Farnstrom, F., Lewis, J., & Elkan, C. (2000). Scalability for clustering algorithms revisited. *SIGKDD Explorations*, 2(1), 51–57. doi:10.1145/360402.360419
- Gama, J. (2010). *Knowledge discovery from data streams*. Boca Raton, FL: CRC Press.
- Ipeirotis, P., Ntoulas, A., & Gravano, L. (2005). Modeling and managing content changes in text databases. In *Proceedings of the 21st International Conference on Data Engineering* (pp. 606-617). Washington, DC: IEEE Computer Society.
- Kalnis, P., Mamoulis, N., & Bakiras, S. (2005). On discovering moving clusters in spatio-temporal data. In *Proceedings of the 9th International Conference on Advances in Spatial and Temporal Databases* (pp. 364-381). Berlin, Germany: Springer-Verlag.
- Mei, Q., & Zhai, C. (2005). Discovering evolutionary theme patterns from text: An exploration of temporal text mining. In *Proceedings of the 11th ACM SIGKDD International Conference on Knowledge Discovery in Data Mining* (pp. 198-207). New York, NY: ACM.
- Spiliopoulou, M. (2011). Evolution in social networks: A survey. In Aggarwal, C. (Ed.), *Social network data analytics* (pp. 147–173). Boston, MA: Kluwer Academic. doi:10.1007/978-1-4419-8462-3_6
- Spiliopoulou, M., Ntoutsis, I., Theodoridis, Y., & Schult, R. (2006). MONIC: Modeling and monitoring cluster transitions. In *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery in Data Mining* (pp. 706-711). New York, NY: ACM.
- Yang, H., Parthasarathy, S., & Mehta, S. (2005). A generalized framework for mining spatio-temporal patterns in scientific data. In *Proceedings of the 11th ACM SIGKDD International Conference on Knowledge Discovery in Data Mining* (pp. 716-721). New York, NY: ACM.
- Zhang, T., Ramakrishnan, R., & Livny, M. (1996). BIRCH: An efficient data clustering method for very large databases. In *Proceedings of the 11th ACM SIGKDD International Conference on Knowledge Discovery in Data Mining* (pp. 103-114). New York, NY: ACM.

ENDNOTES

- ¹ The timepoint 1998 includes publications of both 1997 and 1998, since the former contains only a small number of publications.
- ² Cluster identifiers are generated by the clustering algorithm at each timepoint.

Eirini Ntoutsis is a Post-doctoral researcher at the Ludwig Maximilian University of Munich, Germany, in the Database Systems Group of Prof. Hans-Peter Kriegel. She received her PhD (2008) in Informatics from the University of Piraeus, Greece. She also holds an MSc in Computer Science (2003) and a diploma in Computer Engineering & Informatics (2001) both from the Computer Engineering & Informatics Department (CEID), University of Patras, Greece. Her research focuses on different aspects of knowledge management and can be summarized as pattern extraction, change detection and monitoring over complex dynamic data. She has worked with different types of data arising from different application fields, including text, web, retail industry, robotics, spatio-temporal, user preferences, high-dimensional data. She has published 40 refereed articles in scientific journals and conferences in the above areas. She has presented tutorials on "Mining the Volatile Web" at the ECML PKDD 2005 and on "Mining Complex Dynamic Data" at the ECML PKDD 2011. She has served as a program committee member for ECML-PKDD 2011, SAC 2011-12, KDML 2011 and BASNA 2010-11. She was awarded a post-doctoral fellowship by the Alexander von Humboldt Foundation (2010-2011) and a doctoral scholarship by the Heracleitos program (2003-2007). She worked as a researcher in several European projects, namely PANDA (FP6/IST, 2001-04) and GeoPKDD (FP6/IST, 2005-09). She has worked as an R&D engineer at the Computer Technology Institute (RA-CTI) and as a Data Mining expert at the Hellenic Telecommunications Organization (OTE).

Myra Spiliopoulou is Professor of Business Information Systems in the Faculty of Computer Science of the Otto-von-Guericke University Magdeburg, Germany. Her group "Knowledge Management & Discovery" (KMD) works on data mining, stream mining and web mining for dynamic environments, and develops methods for model adaptation and model monitoring under drift. Her research on topic monitoring, social network monitoring and analysis of complex dynamic data has been published in renowned international conferences and journals. In 2011 she has presented a tutorial on "Mining Complex Dynamic Data" at the ECML PKDD 2011, and has been serving as Workshops Chair for the IEEE Data Mining Conference (ICDM'11, Vancouver, Canada), and as PC Area Chair for the ECML/PKDD 2011 (Athens, Greece). In the coming year she is serving as Senior Reviewer for the SIAM Data Mining Conference (SDM'12), the ECML PKDD 2012 and the ASONAM 2012, and she is PC Co-Chair of the German Classification Society Conference (GfKI'12, Hildesheim, Germany). Her work on "Evolution in Social Networks: A Survey" in the book "Social Network Data Analytics" (ed. Charu Aggarwal) has appeared recently at Springer Verlag.

Yannis Theodoridis is Assoc. Professor at the Department of Informatics, University of Piraeus, where he currently leads the Information Management Lab. Born in 1967, he received his Diploma (1990) and PhD (1996) in Electrical and Computer Engineering, both from the National Technical University of Athens, Greece. Before joining the University of Piraeus, he was member of the research staff at the Hellenic Research Foundation (1997-98) and the Computer Technology Inst. (1999-2002). His research interests include Data Science (management, analysis, mining) for mobility data, whereas he teaches databases, data mining and GIS at under- and post-graduate level. Apart from several national-level projects, he is or was scientist in charge and coordinator of two European projects, namely PANDA (FP6/IST, 2001-04) and CODMINE (FP6/IST, 2002-03), and principal investigator in GeoPKDD (FP6/IST, 2005-09), MODAP (FP7/ICT, 2009-12; member of the mgmt board), MOVE (COST, 2009-13; vice-chair of the mgmt committee), DATASIM (FP7/ICT, 2011-14) and SEEK (FP7/PEOPLE, 2012-15). He has served as general co-chair for SSTD'03, ECML/PKDD'11 and PCI'12, vice PC chair for IEEE ICDM'08, organizing chair for the 2010 summer school on "Mobility, Data Mining, and Privacy", member of the editorial board of the International Journal on Data Warehousing and Mining – IJDWM (2005-), and member of the SSTD endowment (2010-). He has offered several tutorials in top conferences (with the most recent at EDBT'09) and invited lectures in Greece and abroad (including PhD/MSc courses at Venice, Milano, KAUST, Aalborg and Trento) on Mobility Data Management and Data Mining topics. He has co-authored three monographs and more than 100 refereed articles in scientific journals and conferences, receiving more than 800 citations.