

Hierarchical Clustering of Time Series Data Streams

Pedro Pereira Rodrigues, João Gama and João Pedro Pedroso

Abstract—This paper presents and analyzes an incremental system for clustering streaming time series. The Online Divisive-Agglomerative Clustering (ODAC) system continuously maintains a tree-like hierarchy of clusters that evolves with data, using a top-down strategy. The splitting criterion is a correlation-based dissimilarity measure among time series, splitting each node by the farthest pair of streams. The system also uses a merge operator which reaggregates a previously split node, in order to react to changes in the correlation structure between time series. The split and merge operators are triggered in response to changes in the diameters of existing clusters, assuming that, in stationary environments, expanding the structure leads to a decrease in the diameters of the clusters. The system is designed to process thousands of data streams that flow at high-rate. The main features of the system include update time and memory consumption that do not depend on the number of examples in the stream. Moreover, the time and memory required to process an example decreases whenever the cluster structure expands. Experimental results on artificial and real data assess the processing qualities of the system, suggesting competitive performance on clustering streaming time series, exploring also its ability to deal with concept drift.

Index Terms—Data stream analysis, clustering streaming time series, incremental hierarchical clustering, change detection.

I. INTRODUCTION

IN recent days, information has grown importance as the web spread out, with constant increase in communication capabilities, creating a global network interaction of both data and processes. The traditional setting for data analysis had turned gathering data into one of the most difficult tasks in data mining applications. Nowadays, we face the opposite situation. In fact, not rarely the amount of data available from a given source (e.g. sensor networks) is so high that traditional batch systems, which are usually based on memory storage and multiple readings of the same data, simply cannot be used. In recent real-world applications, data flows continuously from a *data stream* at high speed, producing examples over time. Traditional models cannot adapt to the high speed arrival of new examples [1]. This way, new algorithms have been developed that aim to process data in *real-time*. These algorithms should be capable of processing each example in constant time and memory, while they consistently supply a compact data description at each given moment [2]. In this context, quicker responses are usually requested. We need to

continuously maintain a decision model at any time, which should reflect the behavior of most recent data.

Time series data is perhaps the most common kind of data explored by data miners [3]. Clustering is probably the most frequently used data mining algorithm [4], used in exploratory data analysis. Data clustering techniques that work in real-time must allow the update of the clustering definition based only on the current model and on new examples. Besides the resources restrictions, one could be interested in the analysis of the clustering structure, as well as clusters' evolution over time. Systems should be able to refine the cluster structure whenever more information is available and to take into account that the structure can change over time. Among different techniques known in literature, hierarchical models are more versatile as they do not require an a priori definition of the number of clusters. From these, divisive methods seem to be the most appropriate to apply to an online procedure, building the clustering structure using a top-down strategy [5].

Most of the work in incremental clustering of data streams has been concentrated on example clustering rather than variable clustering. Clustering variables (e.g. time series) is a very useful tool for some applications, such as sensor networks, social networks, electrical power demand, stock market, etc. However, incremental clustering of variables is still not a completely covered issue. The standard approach to cluster variables in a batch scenario uses the *transpose* of the working matrix. In a data stream scenario, this is not possible because the transpose operator is a blocking operator [2]. For the task of clustering variables in streams new algorithms are required. At our best knowledge, this is one of the first proposals to achieve this goal.

The main objective of this work is to present an adaptive system to perform hierarchical clustering of variables, where each variable is a time series and each new example that is fed to the system is the value of an observation of all time series in a particular moment. The main characteristics of the system are incremental update with new examples, anytime output of the clustering structure, and the ability to detect and react to changes that may occur in it. This paper is a substantial extension of a previous one [6] published as a short paper in a data mining conference, where only a small overview of the system was presented.

In the next section, a review over incremental clustering analysis for data streams is performed. Section III introduces the proposed hierarchical approach for the task of clustering streaming time series, focusing on the decisions of splitting and aggregation. In section IV, experimental evaluation and real data application results are gathered, supporting the quality of the system. Moreover, scalability and sensitivity tests are presented to assess the robustness of the proposed method. A small discussion over the limitations of the system is presented in section V, and we finalize the exposition with section VI, where concluding remarks and future work are presented.

Manuscript received July 21, 2006; revised May 20, 2007; revised September 27, 2007; accepted November 6, 2007.

Pedro Pereira Rodrigues is with LIAAD - INESC Porto L.A. and the Faculty of Sciences of the University of Porto, Portugal.

E-mail: pprodrigues@fc.up.pt

João Gama is with LIAAD - INESC Porto L.A. and the Faculty of Economics of the University of Porto, Portugal.

E-mail: jgama@fep.up.pt

João Pedro Pedroso is with UESP - INESC Porto L.A. and the Faculty of Sciences of the University of Porto, Portugal.

E-mail: jpp@fc.up.pt

0000-0000/00\$00.00 © 2007 IEEE

II. RELATED WORK

Knowledge discovery systems are usually constrained by three limited resources: time, memory and sample size. In traditional applications sample size limitation was proved to be dominant since it often leads to overfitting. Nowadays, time and memory seem to be the bottleneck for machine learning applications, mainly the last one. Current data mining problems approach a new type of datasets, called *data streams*. According to Guha et al. [7], a *data stream* is an ordered sequence of points that can be read only once or a small number of times. For what is relevant in our work, we will consider reading data only once. This data model emerged from several new applications, such as sensor data, web clicks, credit card usage, multimedia data, etc., which required a different approach for the usual data mining problems. Incremental methods have been developed in various fields of data mining research (e.g. [1], [8]) which can cope with data stream analysis. Likewise, incremental clustering techniques have also emerged from research in the past years.

Clustering is usually taken as a batch procedure, statically defining the structure of objects. Many clustering techniques emerged from research, but we can group most of them in two major paradigms: *partitional* [9] and *hierarchical* [10] approaches. Nevertheless, *density* [11] and *grid-based* [12] systems also present promising research lines. A comprehensive overview on clustering analysis can be found in [9]. Hierarchical algorithms have a major advantage over partitional methods as they do not require a user-predefined number of target clusters.

We should stress that our goal is to perform clustering on the variable domain (the time series) and not to cluster examples, as most of the previous work. While in batch clustering this is not a challenging issue, as *examples* and *variables* can be easily transposed, in data streams clustering the standard matrix transposition is not applicable. In fact, most of the works published in clustering of data streams refer to example clustering. This is one of the first works referring to clustering variables in the data streams framework. In the following we present a concise review of the literature.

One of the first clustering systems to be developed, being both incremental and hierarchical, was the *COBWEB*, a conceptual clustering system that executes a hill-climbing search on the space of hierarchical categorization [13]. This method incrementally incorporates objects in a probabilistic categorization tree, where each node is a probabilistic concept representing a class of objects. The gathering of this information is made by means of the categorization process of the object down the tree, updating counts of sufficient statistics while descending the nodes, and executing one of several operations: classify an object according to an existent cluster, create a new cluster, combine two clusters or divide one cluster into several ones.

A. Partitioning Methods

Bradley et al. [14] proposed the *Single Pass K-Means*, an algorithm that aims at increasing the capabilities of *k-means* for large datasets. The main idea is to use a buffer where points of the dataset are kept in a compressed way. The *STREAM* [15] system can be seen as an extension of [14] which aims to minimize the sum of the squared differences (as in *k-means*) keeping as restriction the use of available memory. *STREAM* processes data into batches of m points. These points are then stored in a

buffer in main memory. After filling the buffer, *STREAM* clusters the buffer into k clusters. It then summarizes the points in the buffer by retaining only the k centroids along with the number of examples in each cluster. *STREAM* discards all the points but the centroids weighted by the number of points assigned to it. The buffer is filled in with new points and the clustering process is repeated using all points in the buffer. This approach results in a one-pass algorithm constant-factor approximation algorithm. The main problem is that *STREAM* never considers data evolution. The resulting clustering can become dominated by the older, outdated data of the stream. However, an interesting aspect of this algorithm is the ability to compress old information, a relevant issue in data stream processing.

Recent research developments are directed towards distributed algorithms for continuous clustering of examples over distributed data streams. Cormode et al. [16] proposed different strategies to achieve this goal, with local and global computations, in order to balance the communication costs. They considered techniques based on the *furthest point* algorithm [17], which gives a approximation for the radius and diameter of clusters with guaranteed cost of two times the cost of the optimal clustering. They also present the *parallel guessing* strategy, which gives a slightly worse approximation but requires only a single pass over the data. They conclude that, in actual distributed settings, it is frequently preferable to track each site locally and combine the results at the coordinator site.

B. Hierarchical Methods

One major achievement in this area of research was the *Balanced Iterative Reducing and Clustering using Hierarchies* system [18]. The *BIRCH* system builds a hierarchical structure of data, the *CF-tree*, a balanced tree where each node is a tuple (*Clustering Feature*). A clustering feature contains the sufficient statistics for a given cluster: the number of points, the sum of each feature-values and the sum of the squares of each feature-value. Each cluster feature corresponds to a cluster, being hierarchically organized in a *CF-tree*. Each non-leaf node in the tree aggregates the information gathered in the descendant nodes. This algorithm tries to find the best groups with respect to the available memory, while minimizing the amount of input and output. The *CF-tree* grows by aggregation with only one pass over the data, thus having complexity $O(N)$. Another use of the *CF-tree* appears in [19]. More than an algorithm, the *CluStream* is a complete system composed by two components, one *online* and another *offline*. Structures called *micro-clusters* are locally kept, having statistical information of data. These structures are defined as a temporal extension of *clustering feature* vectors presented in [18], being kept as images through time, following a pyramidal form. This information is used by the offline component that depends on a variety of user-defined parameters to perform final clustering by an iterative procedure.

The *CURE* system [20], *Clustering Using REpresentatives*, performs a hierarchical procedure that assumes an intermediate approach between centroid based and all-point based techniques. In this method, each cluster is represented by a constant number of points well distributed within the cluster, which capture the extension and shape of the cluster. This process allows the identification of clusters with arbitrary shapes. The *CURE* system also differs from *BIRCH* in the sense that, instead of pre-aggregating all the points, this system gathers a random sample of the dataset,

using Chernoff bounds in order to obtain the minimum number of examples.

C. Concept Change Detection

Most of the work in machine learning assumes that examples are generated at random according to some stationary probability distribution [21]. There are already several methods in predictive machine learning to deal with changing concepts (e.g. [21], [22]). Nevertheless, the notion of concept drift applied to clustering analysis is not directly derived from concept drift on the variables domain, as clustering structure may not be affected by variable's dynamics. Detecting concept drift as usually conceived for one time/order varying variable is not the same as detecting concept drift on the clustering structure of several time/order varying variables. These are usually points in the stream of data where the clustering structure gathered with previous data is no longer valid, since it no longer represents the new relations of dissimilarity between the streams. In this work, we also try to address this feature.

III. ONLINE DIVISIVE-AGGLOMERATIVE CLUSTERING

The task of clustering time series over data streams is not widely studied. In fact, at the best of our knowledge, this is the first proposal of a hierarchical approach to the problem, so we should start by formally introduce it. Afterwards, a complete description of the proposed system is presented, trying to address its main characteristics and innovations.

A. Clustering Streaming Time Series

Data streams usually consist of variables producing examples continuously over time. The basic idea behind clustering streaming time series is to find groups of variables that behave similarly through time. Let $X = \langle x_1, x_2, \dots, x_n \rangle$ be the complete set of n data streams and $X^t = \langle x_1^t, x_2^t, \dots, x_n^t \rangle$ be the example containing the observations of all streams x_i at the specific time t . The goal of a clustering system for multiple time series is to find (and make available at any time t) a partition P of streams, where streams in the same cluster tend to be more alike than streams in different clusters. In partitional clustering, searching for k clusters, the result at time t should be a matrix P of $n \times k$ values, where each P_{ij} is one if stream x_i belongs to cluster c_j and zero otherwise. Specifically, we can inspect the partition of streams in a particular time window from a starting time s until current time t , using examples $X^{s..t}$, which would give a temporal characteristic to the partition. In a hierarchical approach to the problem, the same possibilities apply, with the benefit of not having to previously define the target number of clusters, thus creating a structured output of the hierarchy of clusters. An example partition could be defined as $P^t = \{\{\{x_1\}, \{x_3, x_5\}\}, \{x_2, x_4\}\}$, stating that data streams x_1, x_3, x_5 have some similarity between them (more pronounced between x_3 and x_5), being at the same time somehow dissimilar from x_2 and x_4 .

B. A Hierarchical Approach

In this paper, the ODAC (*Online Divisive-Agglomerative Clustering*) system is presented, which is an algorithm for incremental clustering of streaming time series that constructs a hierarchical tree-shaped structure of clusters using a top-down strategy. The

leaves are the resulting clusters, with each leaf grouping a set of variables. The union of all leaves is the complete set of variables. The intersection of any two leaves is the empty set. The system encloses an incremental distance measure and executes procedures for expansion and aggregation of the tree-based structure, based on the diameters of the clusters.

The main assumption of the system is that decisions taken over a sample of the most recent data are, in the limit and under certain conditions, equivalent to those taken over an infinite set of observations. Given this, the system continuously monitors existing clusters' diameters over time. The diameter of a cluster is the maximum distance between variables of that cluster. For each existing cluster, the system finds the two variables defining the diameter of that cluster. At time t , if a given condition is met on this diameter, the system splits the cluster and assigns each of the chosen variables to one of the new clusters, becoming the *pivot* variable for that cluster. Afterwards, all remaining variables on the old cluster are assigned to the new cluster which has the closest pivot. New leaves start new statistics, assuming that only forthcoming information will be useful to decide whether or not this cluster should be split. Each node c_k will then represent relations between streams using examples $X^{i_k..s_k}$, where i_k is the time at which the node was created and s_k is the time at which the node was split (or current time t for leaf nodes). This feature increases the system's ability to cope with changing concepts as, later on, a test is performed to check if the previously decided split still represents the structure of variables. On stationary data streams, the overall intra-cluster dissimilarity should decrease with each split. This way, if a cluster is split into two child-leaves, the diameter of the new clusters should be less or equal than the diameter of the parent node. If the diameter of a leaf is greater than its parent's diameter, then the previously taken decision no longer reflects the structure of data. The system reaggregates on the cluster's parent, restarting statistics. The forthcoming sections describe the inner core of the system.

C. Incremental Dissimilarity Measure

The system must analyze distances between incomplete vectors, possibly without having any of the previous values available. Thus, these distances must be incrementally computed. Since we want to make decisions with statistical support, we will use the Hoeffding bound to support our decisions, forcing the criterion - the distance measure - to be scaled [23]. We use Pearson's correlation coefficient [24] between time series as *similarity* measure, as done by [25]. Deriving from the correlation between two time series a and b calculated in [26], the factors used to compute the correlation can be updated incrementally, achieving an exact incremental expression for the correlation:

$$\text{corr}(a, b) = \frac{P - \frac{AB}{n}}{\sqrt{A_2 - \frac{A^2}{n}} \sqrt{B_2 - \frac{B^2}{n}}} \quad (1)$$

The *sufficient statistics* needed to compute the correlation are easily updated at each time step: $A = \sum a_i$, $B = \sum b_i$, $A_2 = \sum a_i^2$, $B_2 = \sum b_i^2$, $P = \sum a_i b_i$. In ODAC, the dissimilarity between variables a and b is measured by an appropriate metric, the *Rooted Normalized One-Minus-Correlation*, given by

$$\text{rnomc}(a, b) = \sqrt{\frac{1 - \text{corr}(a, b)}{2}} \quad (2)$$

with range $[0, 1]$. We consider the cluster's *diameter* to be the highest dissimilarity between two time series belonging to the same cluster, or the variable variance in the case of clusters with single variables.

D. Growing the Hierarchy

The main procedure of the ODAC system is to grow a tree-shaped structure that represents the hierarchy of the clusters present in the data. In this system, each example is processed only once. The system incrementally updates, at each new example arrival, the sufficient statistics needed to compute the dissimilarity matrix, enabling its application to clustering data streams. The dissimilarity matrix for each leaf is only computed when it is being tested for splitting or aggregation, after receiving a minimum number of examples. When processing a new example, only the leaves are updated, avoiding computation of unneeded dissimilarities; this speeds up the process every time the structure grows.

1) *Splitting Criteria*: One problem that usually arises with this sort of models is the definition of a minimum number of observations necessary to assure convergence. A common way of doing this includes a user-defined parameter; after a leaf has received at least n_{min} examples it is considered ready to be tested for splitting. Another approach is to apply techniques based on the Hoeffding bound [23] to solve this problem. The Hoeffding bound has the advantage of being independent of the probability distribution generating the observations [1], stating that after n independent observations of a real-valued random variable r with range R , with confidence $1 - \delta$ the true mean of r is at least $\bar{r} - \epsilon$, where \bar{r} is the observed mean of the samples and

$$\epsilon = \sqrt{\frac{R^2 \ln(1/\delta)}{2n}} \quad (3)$$

As each leaf is fed with a different number of examples, each cluster c_k will possess a different value for ϵ , designated ϵ_k . Let $d(a, b)$ be the distance measure between pairs of time series and $D_k = \{(x_i, x_j) | x_i, x_j \in c_k, i < j\}$ be the set of pairs of variables included in a specific leaf c_k . After seeing n samples at the leaf, let $(x_1, y_1) \in \{(x, y) \in D_k | d(x, y) \geq d(a, b), \forall (a, b) \in D_k\}$ be the pair of variables with maximum dissimilarity within the cluster c_k and, in the same way, considering $D'_k = D_k \setminus \{(x_1, y_1)\}$, let $(x_2, y_2) \in \{(x, y) \in D'_k | d(x, y) \geq d(a, b), \forall (a, b) \in D'_k\}$. Let $d_1 = d(x_1, y_1)$, $d_2 = d(x_2, y_2)$ and $\Delta d = d_1 - d_2$ be a new random variable, consisting on the difference between the observed values through time. Applying the Hoeffding bound to Δd , if $\Delta d > \epsilon_k$, we can confidently say that, with probability $1 - \delta$, the difference between d_1 and d_2 is larger than zero, hence selecting (x_1, y_1) as the pair of variables representing the diameter of the cluster. That is,

$$d_1 - d_2 > \epsilon_k \Rightarrow \text{diam}(c_k) = d_1 \quad (4)$$

With this rule, the ODAC system will only split the cluster when the true diameter of the cluster is known with statistical confidence given by the Hoeffding bound. This rule triggers the moment when the leaf has been fed with enough examples to support the decision. Although a time series is not a purely random variable, we have decided to model the time series first-order differences in order to reduce the negative effect of autocorrelation on the Hoeffding bound. Moreover, with this

approach, the missing values can be easily treated with a zero value, considering that, when unknown, the time series is constant.

2) *Resolving Ties*: The rule presented in equation 4 redirects the research to a different problem. There might be cases where the two top-most distances are nearly or completely equal. To distinguish the cases where the cluster has many variables nearly equidistant from the cases where there are two or more highly dissimilar variables, a tweak must be done. Having in mind the application of the system to a data stream with high dimension, possibly with hundreds or thousands of variables, we turn to a heuristic approach. Based on techniques presented in [1] and [21], we introduce a parameter to the system, τ , which determines how long we will let the system check for the real diameter until we force the splitting and aggregation tests. At any time, if $\tau > \epsilon_k$, the system overrules the criterion of equation 4, assuming the leaf has been fed with enough examples, hence it should consider the highest distance to be the real diameter.

3) *Controlling the Growth*: To prevent the hierarchy from growing unnecessarily, we define another criterion that has to be fulfilled to perform the splitting. The splitting criterion should reflect some relation among the distances between variables of the cluster. Given this fact, we can impose a cluster to be split if it includes a high difference between $(d_1 - \bar{d})$ and $(\bar{d} - d_0)$, where d_0 stands for the minimum distance between variables belonging to the cluster and \bar{d} is the average of all distances in the cluster. In our approach, we relate the expression with the global difference $d_1 - d_0$. Our heuristic is the following: for a given cluster c_k , we choose to split this leaf into a node with two child-leaves if the following condition is met:

$$(d_1 - d_0) | d_1 + d_0 - 2\bar{d} | > \epsilon_k \quad (5)$$

This expression gives the global positioning of the mean with respect to the range of the existing distances, representing two inherent concepts: the farther the highest distance is from the minimum distance, the higher is the possibility of a split occurrence; also, the farther the mean distance is from the average of the maximum and minimum distances $((d_1 + d_0)/2)$, the higher is the probability of splitting.

4) *Expanding the Tree*: When a split point is reported, the pivots are variables x_1 and y_1 where $d_1 = d(x_1, y_1)$. The system assigns each of the remaining variables of the old cluster to the cluster which has the closest pivot. The sufficient statistics of each new cluster are initialized. The total space required by the two new clusters is always less than the one required by the previous cluster. Algorithm 1 sketches the splitting procedure.

E. Aggregating at Concept Drift Detection

The main setting of our system is the monitoring of existing clusters' diameters. In the development of a hierarchical structure of clusters, the overall intra-cluster dissimilarity should decrease with each split. This should be observed on stationary data. On stationary data streams, the diameter of a cluster decreases every time a split occurs. However, usual real-world problems deal with non-stationary data streams, where time series that were correlated in the past are no longer correlated to each other in the current time period. They may also be approaching time series of other clusters. The main problem here is to distinguish between scenarios where a leaf grows into a more elaborated structure or a change in the relevance of previous splits is found.

Algorithm 1 TestSplit

Input: A cluster c_k
Output: Boolean value stating if cluster c_k was split or not
 {Tests the cluster for splitting}
 1: let d_1, d_2, d_0 and \bar{d} be the distances previously defined for equations 4 and 5;
 2: **if** $d_1 - d_2 > \epsilon_k$ **or** $\tau > \epsilon_k$ **then**
 3: **if** $(d_1 - d_0)|d_1 - \bar{d} - (\bar{d} - d_0)| > \epsilon_k$ **then**
 4: create c_x and c_y with x_1 and y_1 as pivots, that is, $x_1 \in c_x \wedge y_1 \in c_y$;
 5: **for** each remaining variable $x_i \in c_k$ not yet assigned **do**
 6: **if** $d(x_i, x_1) \leq d(x_i, y_1)$ **then** $x_i \in c_x$;
 7: **else** $x_i \in c_y$;
 8: **end for**
 9: return *True*;
 10: **end if**
 11: **end if**
 12: return *False*;

The strategy that is adopted in this work is based on the analysis of the diameters. In this way, no computation is needed between the variables of different nodes¹. While the correlation structure between time series is stationary, a cluster split will never increase the clusters diameter. For each given leaf c_k , we should search to see if the split decision that created it still represents the structure of data. Thus, we shall test the diameters of c_k and c_k 's parent (c_j), assuming that the child's diameter should not be larger than the diameter of the parent node. We define a new random variable $\Delta a = \text{diam}(c_k) - \text{diam}(c_j)$. Applying the Hoeffding bound to this random variable, if $\Delta a > \epsilon$ then the diameter of the child node is confidently larger than the parent's diameter. Given this, we choose to aggregate on c_j if

$$\text{diam}(c_k) - \text{diam}(c_j) > \epsilon_{jk} \quad (6)$$

where $\epsilon_{jk} = \max(\epsilon_j, \epsilon_k)$, supporting the decision with the bound that was defined with less data. The system decreases the number of clusters as previous division is no longer supported, which means that it may not reflect the best divisive structure for recent data. The resulting leaf starts new computations and a concept drift is detected. Algorithm 2 introduces the algorithm. Figure 4 illustrates the evolution of a cluster structure in time-changing

¹Note that only the leaves are updated with new examples. Each node is associated with a specific time window representing the state of the time series in that period.

Algorithm 2 TestAggregate

Input: A cluster c_k
Output: Boolean value stating if cluster c_k was aggregated or not
 {Tests the cluster for aggregation}
 1: update dissimilarities of c_k , if needed, and $\epsilon_{jk} = \max(\epsilon_j, \epsilon_k)$
 2: **if** $\text{diam}(c_k) - \text{diam}(c_j) > \epsilon_{jk}$ **then**
 3: cut c_j 's sub-tree, turning c_j into a leaf with reset sufficient statistics;
 4: return *True*;
 5: **end if**
 6: return *False*;

data. The complete experience is presented on section IV-C.

F. Memory Usage and Time Complexity

The ODAC system presents the required features of an adaptive learning system. For each leaf in which the diameter is known (with confidence level given by the Hoeffding bound) the system tests for aggregation first, so in case of concept drift it will not start to grow unnecessarily. Algorithm 3 presents our method, merging the splitting with the aggregative procedure.

Algorithm 3 ODAC

Input: A set of streaming time series $X = \langle x_1, x_2, \dots, x_n \rangle$
Output: A hierarchical clustering structure S with leaves (clusters) $L = \langle l_1, l_2, \dots, l_m \rangle$
 1: **repeat**
 2: read new example X^t and update sufficient statistics on the leaves L ;
 3: **for** each leaf l_k not yet tested **do**
 4: update dissimilarities and the Hoeffding bound ϵ_k for this leaf;
 5: **if** TestAggregate(l_k) **or** TestSplit(l_k) **then**
 6: announce new structure S ;
 7: **end if**
 8: **end for**
 9: **until** EOF

Complexity analysis can be done with respect to memory usage and time consumption; in both, our system has some interesting features. A system which aims at efficiently clustering data streams must comply with constant memory usage [2]. In ODAC, the size needed to keep the sufficient statistics at each node with n variables is $O(n^2)$. Let us consider splitting this node into two new leaves, with n_1 and n_2 streams being assigned to each of them, respectively, where $n = n_1 + n_2$. Considering that $(n_1 + n_2)^2 > n_1^2 + n_2^2$, $\forall n_1, n_2 > 0$, although the space used by children nodes is still $O(n_1^2 + n_2^2) \in O(n^2)$ the reduction in the number of sufficient statistics is $O(n_1 n_2)$ with $n_1 n_2 \in [n-1, (n/2)^2]$. System's space complexity is quadratic on the number of variables but constant on the number of examples; this allows the handling of an infinite amount of examples, usually present in data streams. Hence, it is usable in the data streams paradigm, especially considering that the number of variables is constant.

A system which aims at efficiently clustering data streams must also comply with constant execution time, with respect to the number of examples [2]. When updating the system with a new example, the update of sufficient statistics results in $O(n^2)$ operations. This is, as expected, quadratic in the number of variables. The update does not depend on the number of examples seen. Therefore, the system's update time is constant in respect to the number of examples, satisfying the data stream requirements. The splitting procedure (algorithm 1) needs to compute the two maximum, the minimum and the mean of the dissimilarity matrix. This procedure is linear in the number of distances, hence $O(n^2)$, quadratic with the number of variables. When computing dissimilarities, assuming the worst case scenario where only one leaf exists, the number of dissimilarities computed is also $O(n^2)$, thus quadratic with respect to the number of variables. An important feature of this algorithm is that every time a split is performed on a

leaf with n variables, the global number of dissimilarities needed to be computed at the next iteration diminishes at least $n - 1$ (worst case scenario) and at most $(n/2)^2$ (best case scenario). As a result, the time complexity of each iteration of the system is constant with respect to the number of examples, and decreases with every split; it is therefore capable of addressing continuous data streams.

IV. EXPERIMENTAL EVALUATION

ODAC is an online clustering system for time series data streams. Since the scope of the system is very well defined, the experimental evaluation is performed aiming at the verification of specific hypotheses. First, if the system is applied to a set of time series with stationary clustering structure, the system should converge to the real structure of the streams. However, if the streams present dynamic behavior, then the system should detect changes in the clustering structure and adapt it accordingly. Finally, we should evaluate how the system performs on real data produced by applications which generate time series data streams, and how the system scales to the high number of time series usually produced in these contexts. Moreover, sensitivity tests should be performed in order to assess the sensitivity of the system to slight changes in parameters.

A. Evaluation Criteria for Clustering Validity

Usually, the criteria used to evaluate clustering methods concentrate on the quality of the resulting clusters, or their fitness to a given known structure. Given the hierarchical characteristic of our system, we are also interested in assessing the quality of the hierarchy constructed by the algorithm.

We can generally find three approaches to investigate clusters' quality, known as *external*, *internal* and *relative* criteria [4]. *External* criteria make an evaluation based on a pre-specified structure which reflects our prior knowledge or intuition over the data set. Usually used *external criteria* include the *Folkes and Mallows Index* [27] and the *Hubert's Γ Statistic* [28]. An *internal* criterion usually stands for an evaluation based on measures involving the data members themselves (for example, the *Divisive Coefficient* [10] or the *Cophenetic Correlation Coefficient* [9]). The third approach compares the resulting cluster structure with other clustering results for the same data set. Examples of *relative* indices are the *Dunn's Index* [29] and the *Modified Hubert's Γ Index* [4]. The results used to make this comparison may be gathered using the same algorithm with different specifications or using different algorithms.

For the set of experiments, we will focus on two *relative* criteria to find the best number of clusters and assess the quality of the resulting clusters, and an *internal* criterion to assess the hierarchical quality of the resulting clustering structure.

1) *Cluster Quality*: To validate our system, we consider the *MHT* - *Modified Hubert's Γ Statistic* and the *DVI* - *Dunn's Validity Index* in order to assess the quality of the resulting clusters. We have decided to use both as they represent a good complement of each other: the first one is especially motivated for detecting compact and well-separated clusters, but has the drawback of being dependent from the number of clusters; the second criterion does not depend on the number of clusters but is more unstable, as it is based on the single linkage distances and diameters. The

modified Hubert's Γ index is given by

$$MHT = \frac{1}{M} \sum_{i=1}^{N-1} \sum_{j=i+1}^N P_{ij} Q_{ij} \quad (7)$$

where $M = \frac{n(n-1)}{2}$, P is the proximity matrix and Q is a $N \times N$ matrix where each Q_{ij} is the distance between the representative points (centroids, medoids, etc.) of the clusters to which i and j belong. High values of this index represent compact and well-separated clusters. The Dunn's Validity Index is given by

$$DVI = \min_{i,j} \left\{ \frac{d(c_i, c_j)}{\max_k \{diam(c_k)\}} \right\} \quad (8)$$

where $d(c_i, c_j)$ is the single linkage dissimilarity function between two clusters and $diam(c_k)$ is the diameter of cluster c_k . High values of this index also represent compact and well-separated clusters. The *MHT* measure increases with the number of clusters, so the quest for the best value is made by looking at the plot of the measure along a different number of clusters and finding the *knee* in the plot (the number of clusters that produces the highest increase ΔMHT in the quality measure). The second index is independent of the number of clusters. This way, the highest value is considered the best one.

2) *Hierarchy Quality*: Hierarchical clustering methods, such as the ODAC system, allow the analysis of another quality measure, the *CPCC* - *Cophenetic Correlation Coefficient*, which measures quality in hierarchical structures. The *CPCC* is defined as

$$CPCC = \frac{\sum_{i=1}^{N-1} \sum_{j=1+1}^N C_{ij} P_{ij} - \mu_P \mu_C}{\sqrt{\sum_{i=1}^{N-1} \sum_{j=1+1}^N P_{ij}^2 - \mu_P^2 \sum_{i=1}^{N-1} \sum_{j=1+1}^N C_{ij}^2 - \mu_C^2}} \quad (9)$$

where C is the *Cophenetic Proximity Matrix* with each c_{ij} being the proximity level at which the two objects i and j appeared together in the same cluster for the first time and

$$\mu_P = \frac{1}{M} \sum_{i=1}^{N-1} \sum_{j=1+1}^N P_{ij} \quad \text{and} \quad \mu_C = \frac{1}{M} \sum_{i=1}^{N-1} \sum_{j=1+1}^N C_{ij} \quad (10)$$

The closer the value of this index is to 1, the better the match and the better the hierarchy fits the data. High values of any of these indexes indicate the presence of a good clustering structure.

B. System Evaluation on Stationary Environments

As few data streams are really available for free use, the first evaluation of the proposed system is made using synthetic data sets, created with specific characteristics, described in section IV-B.1. The system is then tested in real data, the PDMC Sensor Data Set.

For each data set (with n variables), 10 runs of *k-means* (R [30] implementation) are executed, each one with all the possible numbers of clusters, k . Quality measures are calculated in these runs and an average is considered to find the best number of clusters for the dataset. Afterwards, ODAC is applied in the same data set enabling a comparison between the final structure and the one provided by *k-means*. This procedure is performed for artificial and real data sets. Comparison with a batch DIANA [10]

system, using the same correlation-based distance measure, is also presented. On all experiments, the δ parameter of the Hoeffding bound is set to 0.05, determining a confidence level of 95%; the τ parameter of the ODAC system is set to 0.02. Output is presented in graphics where squared objects are used for leaves and round objects for nodes. The information inside each node represents: node number, diameter for that cluster and the number of examples seen by that node. For leaves, the variables included in the cluster are also displayed.

1) *Evaluation on Artificial Data*: The data sets used in this section were created using a time series generator that produces n time series belonging to a predefined number k of clusters with a noise constant β . Each cluster c_k has a pivot time series p and the remaining time series are created as $p + \lambda$ where

$$\lambda \sim U(-\beta p, \beta p) \quad (11)$$

We have created three data sets with ten variables each, especially prepared to test different hypothesis: a *closed* data set, where all ten variables are created by the same concept; the *two clusters* data set, where two well-defined clusters are created with five variables each and the *4C* data set, representing a more complex yet stationary structure of clusters. External criteria results are not shown here as all runs resulted in the expected clustering structure.

a) *Closed Data Set*: A hierarchical clustering procedure should always create a hierarchical structure of clusters. Nevertheless, if a data set represents a closed cluster, with all variables highly correlated, then we expect a single cluster as the system output. We have created a set of ten time series with 100K examples and tested different values for the parameter β of the artificial data generator. The final result is a single cluster. Although for $\beta \geq 0.3$ some splitting occurred, this was quickly reversed by aggregation. Only for values of $\beta \geq 0.9$ we observed spoiled results, as the cluster spread out, growing an unstable hierarchy of clusters. No meaning can be extracted from the quality measures. Based on this results, we have decided to set $\beta = 0.3$ for the following experiences.

b) *Two Clusters Data Set*: In order to study the splitting capabilities of ODAC when in presence of a simple cut point, we have built a data set with ten time series divided in two clusters:

$$\{\{a1, a2, a3, a4, a5\}, \{a6, a7, a8, a9, a10\}\}.$$

The results are presented on Table I, where the nc column presents the number of clusters with which the results are gathered. In the ODAC system, the correct outcome is achieved immediately. The *DVI* index on *k-means* marks the best (and real) number of clusters (and the real cluster structure). The *CPCC* value for the ODAC result is almost perfect.

c) *Four Clusters (4C) Data Set*: For an evaluation on a more complex structure of clusters, we have created a data set with 100K observations of 10 variables, with the inner cluster noise $\beta = 30\%$ and configuration

$$\{\{a1, a2\}, \{a3, a4, a5\}, \{a6\}, \{a7, a8, a9, a10\}\}.$$

The aim is to check for a good partitioning in order to find the best hierarchy. The *DVI* values on the *k-means* (Table I) mark the true number of clusters, also found by ODAC. Interestingly, the *CPCC* gave a low value, which might be explained by the known feature of the measure lightly favoring the balanced hierarchies.

TABLE I
K-MEANS AND ODAC QUALITY RESULTS FOR THE *Two Clusters* (TOP)
AND *Four Clusters* (BOTTOM) DATA SETS.

Two Clusters					
System	nc	MHI^Γ	ΔMHI^Γ	DVI	CPCC
<i>K-Means</i>	2	0.278	—	1.934	—
	3	0.286	0.008	0.728	—
	4	0.294	0.008	0.732	—
	5	0.298	0.004	0.730	—
<i>ODAC</i>	2	0.166	—	1.936	0.998
Four Clusters (4C)					
System	nc	MHI^Γ	ΔMHI^Γ	DVI	CPCC
<i>K-Means</i>	2	0.178	—	0.995	—
	3	0.366	0.188	1.000	—
	4	0.388	0.022	1.937	—
	5	0.392	0.004	0.734	—
<i>ODAC</i>	4	0.332	—	1.937	0.651

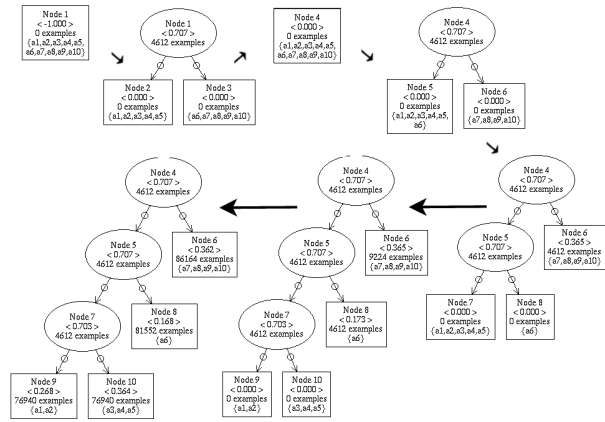


Fig. 1. ODAC system growth and final structure (4C Data Set).

The final structure here is not balanced, as can be seen in the final result presented in Figure 1.

2) *Evaluation on Real-World Data*: The Physiological Data Modeling Contest Workshop (PDMC) was held at the ICML 2004 and aimed at information extraction from streaming sensors data. The training data set for the competition² consists of approximately 10,000 hours of this data, containing several variables: userID, sessionID, sessionTime, characteristic[1..2], annotation, gender and sensor[1..9]. We have concentrated on sensors 2 to 9, since we were interested in finding the relations between different sensor data. Data was extracted by userID, resulting in several data sets of eight continuous variables.

a) *Comparing Users*: In this comparison, we have tried to find the right number of clusters for each user with more than 50K examples, possibly the same for all. We should note that the best result is defined differently for each of the quality measures. For the *MHI* statistic we must seek for the *knee* in the plot, as this measure increases with the number of clusters, whilst for the *DVI* measure the maximum value is used. Table II presents the *k-means* evaluation for PDMC datasets with userID = 1 and userID = 25 (80182 and 141251 observations). Apparently, from a conservative point of view, the best number of clusters is 3. Figure 2 sketches the ODAC system's final structure, for the same data sets. The system achieved the quality level presented in the

²Available at <http://www.bodymedia.com/support/TrainingSet.zip>

TABLE II
K-MEANS AND ODAC QUALITY RESULTS FOR THE PDMC DATA SETS,
FOR USERS 6 (TOP) AND 25 (BOTTOM).

UserID = 6					
System	nc	MHT [†]	Δ MHT [†]	DVI	CPCC
K-Means	2	0.141	—	0.300	—
	3	0.308	0.167	0.300	—
	4	0.311	0.003	0.198	—
	5	0.393	0.082	0.300	—
	6	0.395	0.002	0.156	—
ODAC	3	0.377	—	0.891	0.251

UserID = 25					
System	nc	MHT [†]	Δ MHT [†]	DVI	CPCC
K-Means	2	0.099	—	0.272	—
	3	0.360	0.261	0.325	—
	4	0.362	0.002	0.242	—
	5	0.380	0.018	0.242	—
	6	0.384	0.004	0.228	—
ODAC	3	0.191	—	1.026	0.479

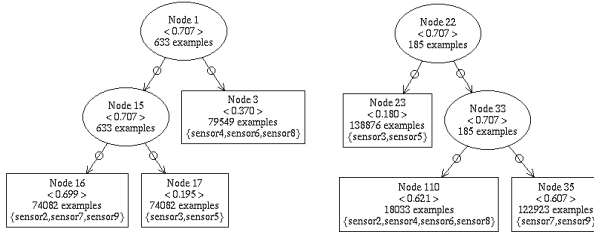


Fig. 2. ODAC final structure gathered with the physiological data with userID = 6 (left) and userID = 25 (right).

bottom line of Table II. The same structure is found, agreeing with the expected number of clusters given by the *k-means* experience. For the data set with userID = 25, again the value 3 for the number of clusters appears. This time, the structure found by ODAC is a bit different but also represents three clusters.

b) *Comparing ODAC and DIANA*: We have tried to find some cluster structure on the data variables, to compare it with a batch divisive analysis system, using the same dissimilarity measure for each user. For userID = 1, 93344 examples were processed. Our system evolved splitting and aggregating, until a stable structure was gathered after 55K examples, never changing from then on. Figure 3 shows the comparison with a batch DIANA system, using the same correlation-based distance measure. The final result is approximately equivalent to the structure achieved by the batch system using the same dissimilarity measure. For

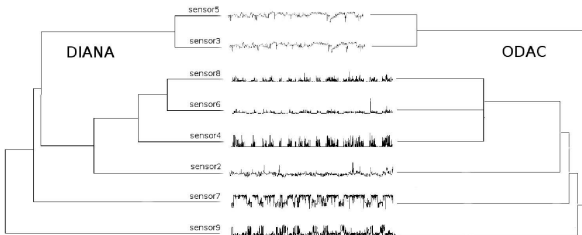


Fig. 3. Comparison between ODAC and DIANA (data for userID = 1). Although higher-level structure is different, at cluster level both structures are quite close. This plot presents the proposed quality of the incremental system when compared with a batch system.

other userIDs, results were very similar.

C. System Evaluation on Dynamic Environments

For a more complex evaluation of system's dynamics (splitting and aggregation), we have defined experiments to determine if the system is able to detect and adapt to changes in the clustering structure. We have generated a dataset by concatenating examples from two known structures, hence simulating a concept change. The data set is equivalent to 4C in the first 50K examples and the last 50K examples were generated from a different concept. The concepts used in the drifting data set shifted from

$$\{\{a1, a2\}, \{a3, a4, a5\}, \{a6\}, \{a7, a8, a9, a10\}\}$$

to

$$\{\{a1, a2, a3, a4, a10\}, \{a5, a6, a7\}, \{a8, a9\}\}$$

The evaluation criteria are: the definition of the clustering structure, the ability to detect change, the delay in detection, and the ability to react and to converge to the new concept.

The evolution of the clustering structure is sketched in Figure 4. The figure presents only the structures defined after an event (split or merge) has occurred. The figure contains also the number of examples read at the time of each snapshot and the time elapsed since concept drift. The system converged at the first concept's structure in 18448 examples. The concept drift occurs at example 50K and was detected 3220 examples later. The cluster structure collapsed to a single leaf at example 62448. Only after 9224 examples from the moment the complete structure collapsed, at example 71672 (21672 examples after the drift), the system reached the correct structure for the second concept, which remains stable until the end of the data set (100K). Simultaneously, we observed a dynamic behavior in the system, splitting and aggregating along the examples, until it stabilized on the final tree structure. The values of Table III are separated for each concept. The global results suggest good performance on splitting, concept drift detection and aggregation.

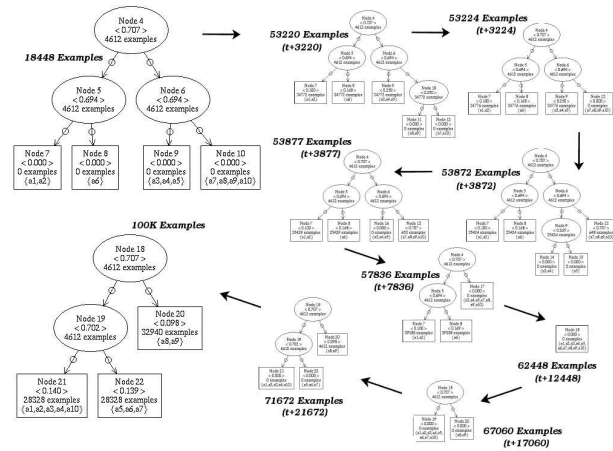


Fig. 4. ODAC structure evolution (concept drift data set). Start: First concept is defined for the data set; 18448 exs: ODAC reaches first concept's structure and stabilizes; 50000 exs (t): Concept drift occurs in the data set; 53220 exs (t + 3220): ODAC detects changes in the structure; 62448 exs (t + 12448, s): ODAC collapses all structure; 71672 exs (t + 21672, s + 9224): ODAC gathers second concept and stabilizes; End: Second concept remains in the data set and the correct final structure of the second concept was discovered.

TABLE III
K-MEANS AND ODAC QUALITY RESULTS FOR BOTH CONCEPTS OF THE
CONCEPT DRIFT DATA SET.

First Concept					
System	nc	MHI [†]	Δ MHI [†]	DVI	CPCC
K-Means	2	0.176	—	0.993	—
	3	0.321	0.145	0.996	—
	4	0.387	0.066	3.000	—
ODAC	5	0.387	0.000	0.712	—
	4	0.331	—	3.000	0.839
Second Concept					
System	nc	MHI [†]	Δ MHI [†]	DVI	CPCC
K-Means	2	0.278	—	0.996	—
	3	0.345	0.067	3.000	—
	4	0.345	0.000	0.710	—
ODAC	5	0.346	0.001	0.710	—
	3	0.146	—	3.000	0.801

To attest the ability to react to concept drift, we have created 30 data sets with the same structure, and monitored how long the system takes to react to drift and stabilize in the new concept, even with high levels of noise. The average number of examples needed to detect drift was around 4900 and the average number of examples needed to stabilize on the next concept was 8500 after the detection.

D. System Performance on Real-World Electrical Data

Time series of electrical data are one of the most widely studied sets of data, mainly for forecast purposes usually by means of neural networks. One big problem with this approach is the time consumption of the neural network's training procedure. This fact, in conjunction with high dimensionality (common electrical networks combine thousands of different time series) suggests the need to cluster time series. Usually, an expert is required for this job. Our system may present some benefits in this particular task.

From the raw data received at each sub-station, gathered during four months, we have extracted only the variables related to active power (according to the, possibly erroneous, variable ID). Each example represents the average value for the past five minutes of each variable. This resulted in a data set with 34734 examples (120 days, 14 hours and 30 minutes) with 887 different streams (variables). The system finds a complete structure of clusters over time, aggregating when necessary. The final structure can be analyzed in Figure 5. The image is rather large and is presented only as an example of the structure obtained as a whole, revealing a tree-like structure away from the list-like worst-case scenario. We present a sketch of the variables, for one of the clusters, in Figure 6. This cluster expresses good intra cluster correlation ($\mu = 0.629$, $\sigma = 0.227$), considering that this is real data from five variables along 7385 observations. Nevertheless, clusters with worse intra cluster correlation may be found with some noisy variables included. This is probably due to the lack of examples fed to these clusters, reducing the confidence on splitting. It is expected that, with more examples, these clusters would be split in the future.

1) *Scalability Evaluation*: A different set of sensors was used to conclude scalability characteristics of our system. Taking into account only the current intensity sensors, we have aggregated observations on an hourly basis over more than two and a half

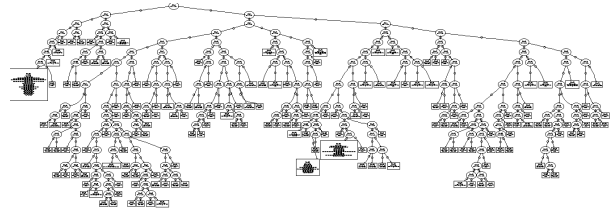


Fig. 5. Electrical demand data - active power. The resulting tree is quite large but is presented here to observe that the final structure is far from the worst case scenario (list-like).

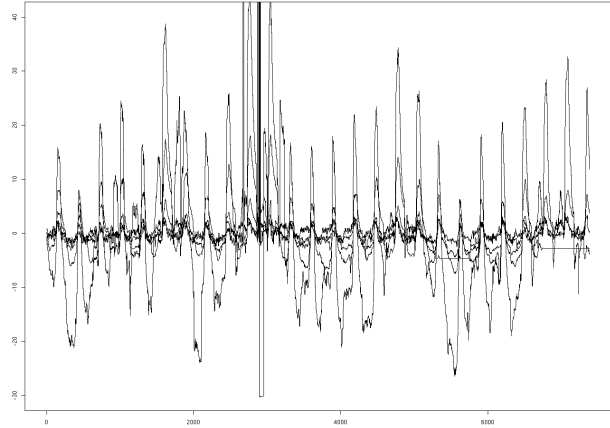


Fig. 6. Cluster from electrical data set: active power. In this cluster, five variables are kept together, along 7385 examples, corresponding to 25 days, 15 hours and 25 minutes. From this plot it is clear the similarity between time series, as the mean intra cluster correlation is 0.629 with standard deviation 0.227.

years. This data set represents 2700 sensors along 22364 observations. For this data set we monitored both update speed and memory usage along time, as these are two of the main positive arguments of the proposed system. For all measures, we present a tendency graph based on an average on a week window. We can note that the overall speed of the system (Figure 7) consistently grows along time, as the structure is growing, being even more pronounced in the update speed (Figure 8). It becomes clear the advantages of local computation in the speed-up achieved by splitting. Moreover, even with the increasing number of clusters to test, the total processing speed increases with splits.

This is, in fact, motivated by the reduction of stored values

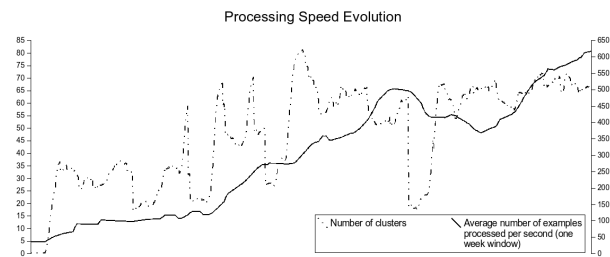


Fig. 7. Total processing speed evolution: This graph presents the evolution of the speed at which the system processes examples, both updating the sufficient statistics and testing clusters for splitting and aggregation, in terms of examples per second, averaged over a week window. The number of clusters is presented along time with respect to the secondary y-axis. Even with the increasing number of clusters to test, the total processing speed increases with splits.

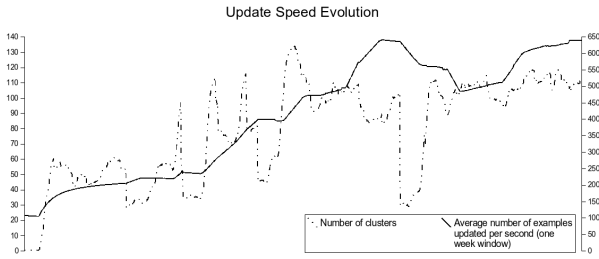


Fig. 8. Update speed evolution: This graph presents the evolution of the speed at which the system updates the sufficient statistics, in terms of examples per second, averaged over a week window. The number of clusters is presented along time with respect to the secondary y-axis. It becomes clear the advantages of local computation in the speed-up achieved by splitting.

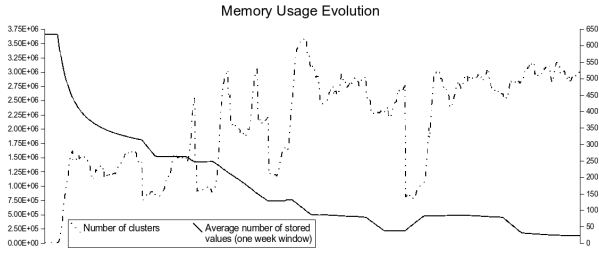


Fig. 9. Memory usage evolution: This graph presents the evolution of the global memory usage of system, in terms of values needed for sufficient statistics, averaged over a week window. The number of clusters is presented along time with respect to the secondary y-axis. The reduction on the stored values is significant: when the number of clusters is 1/5 of the number of sensors, the memory usage diminishes to 5% of the original size.

needed for sufficient statistics computation, as seen in Figure 9. The focus of the system on local computations reduces the problem of needing quadratic memory with respect to the number of sensors.

2) *Sensitivity to the δ Parameter:* We have used a subset of the current intensity sensors data to check the system's sensitivity to changes in the parameter δ used by the Hoeffding bound. An increase of this parameter represents a decrease on the confidence level of the Hoeffding bound. This data set represents 10 sensors along 6000 observations. Figure 10 plots the changes in the three measured validity indices for the selected sensors. The results on the Dunn's Validity Index and the Modified Hubert's Γ statistic indicate that δ has a small impact in the quality of the induced clusters. This parameter has more impact in the quality of the hierarchy found by the system, as the fluctuations in the CPCC statistic seem to indicate.

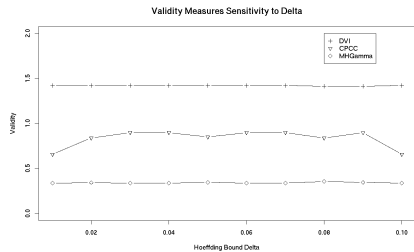


Fig. 10. Sensitivity of the system to changes in the δ parameter. This parameter seems to have some influence on the hierarchical quality of the structure (CPCC) without much influence on the clustering validity (DVI and MHI).

V. STRENGTHS AND LIMITATIONS

Hierarchical models present some characteristics that strengthen their usability, e.g. they do not need a predefined number of target clusters. In the case of data streams, the most relevant property of hierarchical clustering is that the computational resources (time and memory) needed to process each example reduce as the hierarchy grows. Moreover, this decrease in the number of computed dissimilarities is lower-bounded by the number of variables in the original cluster. This feature becomes even more relevant as the number of streaming series grows unbounded. In such cases, computation of dissimilarities at root level becomes a bottleneck, so the ability to grow a hierarchical structure and only compute dissimilarities locally to each cluster presents a much better approach to problems with thousands of time series.

However, in the data streams framework, the correlation structure between variables may change over time. This fact may require a re-structuring of the clustering hierarchy, that is, to move variables from one cluster to another. The main problem is that this *moving* operation must follow the path defined by the actual structure. For example, considering the following situation where the starting structure $\{\{a1, a2\}, \{a3, a4\}\}, \{\{a5, a6\}, \{a7, a8\}\}$ changes to $\{\{a1, a2, a5\}, \{a3, a4\}\}, \{\{a6\}, \{a7, a8\}\}$, although only one attribute ($a5$) has changed, the system has to completely destroy the current hierarchy and re-construct it from scratch. This behavior is expected to occur in any hierarchical model. However, we should note that this example is in fact the worst case scenario. Maintaining dissimilarity matrices in all nodes would enable quick responses to abrupt changes in the lowest levels. However, this approach would also require exponential time and space to process examples. Our approach, based on updating the correlation matrix only at leaves, is limited to changes that smoothly evolve over time. Overall, the gain in not computing all dissimilarities at every bunch of examples compensates the possible burden of re-constructing the hierarchy from scratch from time to time, if and when a concept drift might occur. Fuzzy clustering, where variables are shared in different clusters, seems to be a promising alternative. The main problem is how to achieve the trade-off between fast answer to changes and the time and space required to update the model.

VI. CONCLUSIONS AND FUTURE WORK

In this paper we have presented ODAC, a clustering system for streaming time series. ODAC uses a top-down strategy to construct a binary tree hierarchy of clusters with the goal of finding highly correlated sets of variables. A common measure of cluster quality is the cluster's diameter, which is defined as the highest dissimilarity between objects of the same cluster. The system evolves by continuously monitoring the clusters' diameters.

The examples are processed as they arrive, using a single scan over the entire data. The system incrementally computes the dissimilarities between time series, maintaining and updating the sufficient statistics at each new example arrival, updating only the leaves. The splitting criterion is supported by a confidence level given by the Hoeffding bound, which detects when the system has gathered enough information to confidently define the diameter of each individual cluster.

ODAC is designed towards thousands of data streams that flow at high-rate. Two main characteristics are update time and

memory consumption. Both reduce whenever the tree structure grows. This is a major achievement accomplished by ODAC since only dissimilarities at the leaves must be computed. This way, every time the system grows it becomes faster, overcoming the bottleneck of having to compute all dissimilarities at root level, which is known to have quadratic complexity on the number of streams.

The system includes an agglomerative phase, based on the diameters of existing clusters, also supported by the Hoeffding bound. The aggregation phase enables the adaptation of the cluster structure to smooth changes in the correlation structure of time series. At our best knowledge, this is the first incremental and hierarchical whole clustering system designed for high-speed time series, being able to dynamically adjust the clustering structure in reaction to environment changes.

Experimental results indicate that the performance is competitive for clustering time series, showing that the system is nearly equivalent to a batch divisive clustering on stationary time series. Experimental evaluation using artificial data sets shows the robustness of the system with respect to different time series clustering hypotheses, and the ability to adapt in the presence of evolving concepts. Results using real-world data sets show competitive performance when compared with batch clustering analysis. They also reveal good performance on finding the correct number of clusters, obtained by a bunch of runs of *k-means*.

When considering the expansion of the structure, the strict splitting of variables appears as a possible drawback, in the sense that a previous decision of moving a variable to a leaf, when there is no statistical confidence on the decision of assignment, may split variables that should be together. An approach based on fuzzy sets [31] would let forthcoming examples decide what to do with those variables. This approach is already scheduled as a future work task, as it seems to be a very important option in time series incremental clustering. Moreover, the computation effort forced by the high dimensionality of the data being processed may represent a drawback to the clustering procedure. Thus, techniques based on clipping [32] could also be considered.

ACKNOWLEDGMENT

Pedro P. Rodrigues is supported by a PhD grant by the Portuguese Foundation for Science and Technology (SFRH/BD/29219/2006). The authors also wish to thank the Pluri-annual financial support attributed to LIAAD, and the participation of projects ALES II (POSC/EIA/55340/2004) and RETINAE (PRIME/IDEIA/70/00078).

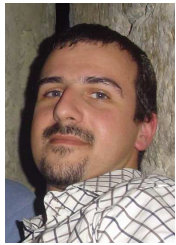
REFERENCES

- [1] P. Domingos and G. Hulten, "Mining high-speed data streams," in *Proceedings of the Sixth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM Press, 2000, pp. 71–80.
- [2] D. Barabási, "Requirements for clustering data streams," *SIGKDD Explorations*, vol. 3, no. 2, pp. 23–27, January 2002.
- [3] E. Keogh and S. Kasetty, "On the need for time series data mining benchmarks: A survey and empirical demonstration," in *Proceedings of the Eighth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM Press, July 2002, pp. 102–111.
- [4] M. Halkidi, Y. Batistakis, and M. Varziargiannis, "On clustering validation techniques," *Journal of Intelligent Information Systems*, vol. 17, no. 2-3, pp. 107–145, 2001.
- [5] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference and Prediction*. Springer Verlag, 2000.
- [6] P. P. Rodrigues, J. Gama, and J. P. Pedrosa, "ODAC: Hierarchical clustering of time series data streams," in *Proceedings of the Sixth SIAM International Conference on Data Mining*, J. Ghosh, D. Lambert, D. Skillicorn, and J. Srivastava, Eds. Bethesda, Maryland, USA: SIAM, April 2006, pp. 499–503.
- [7] S. Guha, A. Meyerson, N. Mishra, R. Motwani, and L. O'Callaghan, "Clustering data streams: Theory and practice," *IEEE Transactions on Knowledge and Data Engineering*, vol. 15, no. 3, pp. 515–528, 2003.
- [8] F. Ferrer, J. Aguilar, and J. Riquelme, "Incremental rule learning and border examples selection from numerical data streams," *Journal of Universal Computer Science*, vol. 11, no. 8, pp. 1426–1439, 2005.
- [9] A. K. Jain and R. C. Dubes, *Algorithms for Clustering Data*. Prentice-Hall, 1988.
- [10] L. Kaufman and P. J. Rousseeuw, *Finding Groups in Data: An Introduction to Cluster Analysis*. John Wiley and Sons, 1990.
- [11] M. Ester, H.-P. Kriegel, J. Sander, and X. Xu, "A density-based algorithm for discovering clusters in large spatial databases with noise," in *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining*, E. Simoudis, J. Han, and U. Fayyad, Eds. AAAI Press, 1996, pp. 226–231.
- [12] W. Wang, J. Yang, and R. R. Muntz, "STING: A statistical information grid approach to spatial data mining," in *Proceedings of the Twenty-Third International Conference on Very Large Data Bases*, M. Jarke, M. J. Carey, K. R. Dittrich, F. H. Lochovsky, P. Loucopoulos, and M. A. Jeusfeld, Eds. Morgan Kaufmann, 1997, pp. 186–195.
- [13] D. H. Fisher, "Knowledge acquisition via incremental conceptual clustering," *Machine Learning*, vol. 2, no. 2, pp. 139–172, 1987.
- [14] P. Bradley, U. Fayyad, and C. Reina, "Scaling clustering algorithms to large databases," in *Proceedings of the Fourth International Conference on Knowledge Discovery and Data Mining*. AAAI Press, 1998, pp. 9–15.
- [15] L. O'Callaghan, A. Meyerson, R. Motwani, N. Mishra, and S. Guha, "Streaming-data algorithms for high-quality clustering," in *Proceedings of the Eighteenth Annual IEEE International Conference on Data Engineering*. IEEE Computer Society, 2002, pp. 685–696.
- [16] G. Cormode, S. Muthukrishnan, and W. Zhuang, "Conquering the divide: Continuous clustering of distributed data streams," in *Proceedings of the 23rd International Conference on Data Engineering (ICDE 2007)*, 2007, pp. 1036–1045.
- [17] T. F. Gonzalez, "Clustering to minimize the maximum inter-cluster distance," *Theoretical Computer Science*, vol. 38, no. 2-3, pp. 293–306, 1985.
- [18] T. Zhang, R. Ramakrishnan, and M. Livny, "BIRCH: An efficient data clustering method for very large databases," in *Proceedings of the 1996 ACM SIGMOD International Conference on Management of Data*. ACM Press, 1996, pp. 103–114.
- [19] C. Aggarwal, J. Han, J. Wang, and P. Yu, "A framework for clustering evolving data streams," in *VLDB 2003, Proceedings of Twenty-Ninth International Conference on Very Large Data Bases*. Morgan Kaufmann, September 2003, pp. 81–92.
- [20] S. Guha, R. Rastogi, and K. Shim, "CURE: An efficient clustering algorithm for large databases," in *Proceedings of the 1998 ACM SIGMOD International Conference on Management of Data*, L. M. Haas and A. Tiwary, Eds. ACM Press, 1998, pp. 73–84.
- [21] J. Gama, P. Medas, and P. Rodrigues, "Learning decision trees from dynamic data streams," in *Proceedings of the 2005 ACM Symposium on Applied Computing*, H. Haddad, L. M. Liebrock, A. Omicini, and R. L. Wainwright, Eds. ACM Press, March 2005, pp. 573–577.
- [22] G. Hulten, L. Spencer, and P. Domingos, "Mining time-changing data streams," in *Proceedings of the Seventh ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM Press, 2001, pp. 97–106.
- [23] W. Hoeffding, "Probability inequalities for sums of bounded random variables," *Journal of the American Statistical Association*, vol. 58, no. 301, pp. 13–30, 1963.
- [24] K. Pearson, "Regression, heredity and panmixia," *Philosophical Transactions of the Royal Society*, vol. 187, pp. 253–318, 1896.
- [25] L. Leydesdorff, "Similarity measures, author cocitation analysis, and information theory," *Journal of the American Society for Information Science and Technology*, vol. 56, no. 7, pp. 769–772, 2005.
- [26] M. Wang and X. S. Wang, "Efficient evaluation of composite correlations for streaming time series," in *Advances in Web-Age Information Management - WAIM 2003*. Springer Verlag, 2003, pp. 369–380.
- [27] E. Fowlkes and C. Mallows, "A method for comparing two hierarchical clusterings," *Journal of American Statistical Association*, vol. 78, no. 383, pp. 553–569, 1983.

- [28] L. Hubert and J. Schultz, "Quadratic assignment as a general data-analysis strategy," *British Journal of Mathematical and Statistical Psychology*, vol. 29, pp. 190–241, 1975.
- [29] J. Dunn, "Well separated clusters and optimal fuzzy partitions," *Journal of Cybernetics*, vol. 4, no. 1, pp. 95–104, 1974.
- [30] R Development Core Team, *R: A language and environment for statistical computing*, R Foundation for Statistical Computing, Vienna, Austria, 2005. [Online]. Available: <http://www.R-project.org>
- [31] L. A. Zadeh, "Fuzzy sets," *Information and Control*, vol. 8, no. 3, pp. 338–353, 1965.
- [32] A. Bagnall and G. Janacek, "Clustering time series with clipped data," *Machine Learning*, vol. 58, no. 2–3, pp. 151–178, 2005.



João Gama received the PhD degree in Computer Science from the University of Porto, Portugal, in 2000. He is currently an Assistant Professor at the Faculty of Economics of the University of Porto, and a researcher at the Artificial Intelligence and Decision Support Laboratory. His main research interests are online learning from data streams, combining classifiers and probabilistic reasoning.



Pedro Pereira Rodrigues received the BSc and the MSc degrees in Computer Science from the Faculty of Sciences of the University of Porto, Portugal, in 2003 and 2005, respectively. He is currently a PhD candidate at the same university, working as a researcher at the Artificial Intelligence and Decision Support Laboratory, on distributed clustering of streaming data from sensor networks. His research interests concentrate on machine learning and data mining from distributed data streams and reliability of predictive and clustering analysis in streaming environments, with applications on industry- and health-related domains.



João Pedro Pedroso received the BSc in Chemical Engineering at the Faculty of Engineering of the University of Porto, Portugal, in 1989, and the PhD degree in Applied Sciences, group of Engineering Mathematics, at the Université Catholique de Louvain, Belgium, in 1996. He is currently an Assistant Professor at the Faculty of Sciences of the University of Porto. His research interests concentrate on optimization, mathematical programming and operations research.