

# Automatic Sequential Pattern Mining in Data Streams

Koki Kawabata\*

ISIR, Osaka University  
koki88@sanken.osaka-u.ac.jp

Yasuko Matsubara\*

ISIR, Osaka University  
yasuko@sanken.osaka-u.ac.jp

Yasushi Sakurai\*

ISIR, Osaka University  
yasushi@sanken.osaka-u.ac.jp

## ABSTRACT

Given a large volume of multi-dimensional data streams, such as that produced by IoT applications, finance and online web-click logs, how can we discover typical patterns and compress them into compact models? In addition, how can we incrementally distinguish multiple patterns while considering the information obtained from a pattern found in a streaming setting? In this paper, we propose a streaming algorithm, namely *STREAMSCOPE*, that is designed to find intuitive patterns efficiently from event streams evolving over time. Our proposed method has the following properties: (a) it is effective: it operates on semi-infinite collections of co-evolving streams and summarizes all the streams into a set of multiple discrete segments grouped by their similarities. (b) it is automatic: it automatically and incrementally recognizes such patterns and generates models for each of them if necessary; (c) it is scalable: the complexity of our method does not depend on the length of the data streams. Our extensive experiments on real data streams demonstrate that *STREAMSCOPE* can find meaningful patterns and achieve great improvements in terms of computational time and memory space over its full batch method competitors.

## CCS CONCEPTS

• Information systems → Data mining;

### ACM Reference Format:

Koki Kawabata, Yasuko Matsubara, and Yasushi Sakurai. 2019. Automatic Sequential Pattern Mining in Data Streams. In *The 28th ACM International Conference on Information and Knowledge Management (CIKM'19), November 3–7, 2019, Beijing, China*. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3357384.3358002>

## 1 INTRODUCTION

Pattern discovery in time series has attracted a considerable amount of interest because it has been widely used in numerous applications such as sensor monitoring [19, 26], online social activity [18], medical data analysis [16] and financial data analysis [27]. However, recently, the emergence of small intelligent devices has opened up new opportunities for various types of applications, such as automated factories, smart cities and connected healthcare, namely the “Internet of Things (IoT)” [10, 24]. Machines are getting in on

\*Artificial Intelligence Research Center, The Institute of Scientific and Industrial Research (ISIR), Osaka University. Mihogaoka 8-1, Ibaraki, Osaka, 567-0047, Japan

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or to publish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

CIKM '19, November 3–7, 2019, Beijing, China

© 2019 Association for Computing Machinery.

ACM ISBN 978-1-4503-6976-3/19/11...\$15.00

<https://doi.org/10.1145/3357384.3358002>

the act. This sudden explosion of intelligent connected devices imposes new requirements on data stream mining, which include (a) real-time modeling, and (b) automatic mining. The ideal method should capture the dynamics of time-series data, efficiently and automatically over IoT data streams. It should also enable us to perform model estimation and advanced analytics, such as anomaly detection, forecasting, and pattern recognition, without any parameter tuning steps.

For example, consider a large collection of observations continuously generated by sensors in an automated factory. The most important goal is to acquire knowledge of sequence patterns reflecting various aspects of the systems such as ordinary processing and signs of accidents, which may have not appeared in the historical/training data, and also to quickly utilize this knowledge for subsequent pattern recognition. In this situation, the data sources generate data with no end in sight, making it difficult to efficiently identify meaningful (i.e., frequent, periodic and abnormal) patterns over one or multiple attributes. Given such data streams, how can we adaptively recognize newly arrived patterns as well as already known ones? How do we efficiently and exactly identify such dynamical patterns?

This paper focuses on an important time series analysis task, namely, automatic and efficient pattern discovery that can reveal all meaningful patterns, which we refer to as a *regime*, consisting of a semi-infinite time-evolving data streams. Intuitively, we design a streaming algorithm, called *STREAMSCOPE*, thus allowing us to deal with the following informal problem.

**INFORMAL PROBLEM 1.** *Given a data stream  $X$ , which is composed of a series of multi-dimensional vectors, i.e.,  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_t\}$ , our challenging problem has the following components:*

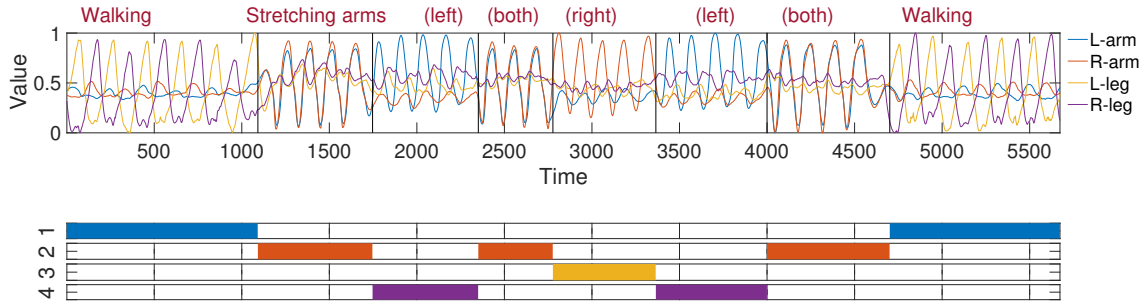
- *splitting  $X$  into a set of segments by the cut points of multiple meaningful patterns,*
- *estimating the number of similar segment groups, i.e., regimes,*
- *summarizing all such dynamical patterns into model parameters,*

*quickly and automatically in a streaming setting.*

### 1.1 Preview of our results

Figure 1 shows the online pattern discovery result of *STREAMSCOPE*. The top of Figure 1 shows the original motion capture stream corresponding to left/right legs and arms. The data stream consists of four intuitive motion patterns: walking, stretching left, right, and both arms.

As shown in the bottom of Figure 1, our method can incrementally find these four patterns (i.e., regimes) and in terms of the time-position of each transition (i.e., the colored rectangles in the figure). More specifically, our method first discovers the motion “walking” and summarizes the feature as regime #1. After the motion has been changed to “stretching arms” at time tick 1300, our method



**Figure 1: Online pattern discovery with STREAMSCOPE: the original motion sensor streams (top), and the output of our method (bottom). At every time tick, the method incrementally identifies intuitive motions (such as walking and stretching arms) and classifies them into groups, namely regimes. If necessary, it automatically generates a new regime with no prior knowledge.**

autonomously generates regime #2 for the new pattern because the two regimes give us better summarization for the data stream. Similarly, as the kinds of stretching arms increases, regime #2 becomes insufficient for capturing the features. Thus, it employs regime #3 and #4 at time tick 3600. If a pattern is observed again (i.e., stretching both arms and walking), then it can recognize the change with the models of already estimated regimes and assign the part of the stream as a segment into the most suitable regime.

As a consequence, our method can perform incremental segmentation and dynamic modeling simultaneously, without any prior knowledge and user intervention.

### 1.2 Contributions

In this paper, we present an online method, namely, STREAMSCOPE, for automatic pattern discovery in co-evolving data streams. In summary, the contributions of this paper are as follows:

- *Effective*: it efficiently and correctly realizes pattern discovery in a streaming fashion. Our method can find high-level patterns, i.e., the dynamics not only within a single pattern but also between patterns.
- *Automatic*: for each pattern it can automatically and incrementally estimate the number, transitional points and model parameters. All of these characteristics are unknown in advance.
- *Scalable*: thanks to our proposed incremental update strategy, it is fast and the computational time and memory space are independent of the entire length of the input streams.

We evaluate our approach on real datasets and show that it discovers meaningful patterns effectively and efficiently.

**Outline.** The remainder of this paper is organized as follows. We first introduce related work, followed by problem statement, our proposed model and algorithms, experiments and conclusions.

## 2 RELATED WORK

Time series pattern discovery has been studied especially in relation to database and data mining [3, 6, 28, 34]. Table 1 shows the relative advantages of our method. We loosely classify the related work into two groups.

**Modeling dynamics and segmentation.** Traditional approaches applied to time series modeling include hidden Markov models

**Table 1: Capabilities of approaches. Only STREAMSCOPE meets all requirements.**

	HMM/++	pHMM/TICC	AUTOPLAIT	CLUSSTREAM/++	MOTIF/++	ADS+	STREAMSCAN	STREAMSCOPE
Data stream	-	-	-	✓	✓	-	✓	✓
Online optimization	-	-	-	✓	✓	✓	-	✓
Data compression	✓	✓	✓	✓	-	-	-	✓
Segmentation	-	✓	✓	-	✓	✓	✓	✓
Regime identification	-	-	✓	-	-	-	-	✓
Parameter free	-	-	✓	-	-	-	-	✓

(HMM), autoregression (AR), linear dynamical systems (LDS), and their variants [5, 15, 17]. More recent advances are based on modeling non-linear behavior on co-evolving time series [21, 23]. However, they are incapable of multiple distinct patterns in large IoT sensor data. For modeling such patterns, there has also been study on segmenting time series by their similar dynamics. Wang et al. [31] presented a pattern-based hidden Markov model (pHMM), which converts a time series into a sequence of line segments, and learns a Markov model from the sequence. Hallac et al. [11] proposed Toeplitz inverse covariance-based clustering (TICC). It is based on a Markov random field (MRF), and can find interpretable clusters, each of which captures the interdependencies between multivariate observations in a typical subsequence of that cluster. However, these methods need parameter tuning for error thresholds, model structure, the number of clusters, etc.

To achieve automatic mining, the minimum description length (MDL) principle has been applied in various contexts [14, 20, 29]. In particular, AutoPlait [20] can automatically identify the intuitive number of high-level patterns (i.e., regimes) using hierarchical HMM-based model, but unfortunately, none of above are streaming methods. The ideal method should require no parameter settings and analyze data in an online fashion.

**Data mining on streams.** The processing and mining of data streams have also attracted significant interest [12, 32, 33]. Ghemoune et al. [7] provides a comprehensive overview of basic

clustering approaches for data streams, such as CluStream [1] and DenStream [4]. EDMStream [8] has recently proposed, which can adaptively change cluster centers to handle time-varying data distribution known as *concept drift* [24]. However, these algorithms process each data point individually and ignore temporal dependency between data, which is unsuitable for extracting the dynamics of time-evolving streams. Time series motif discovery is a well-studied research topic [13, 25, 30], which enables us to detect the frequently occurring subsequences in time series efficiently. ADS+ [34] provides a new adaptive indexing approach for big data series, where the index is built incrementally and adaptively. However, these methods use distance-based metrics (e.g., DTW and Euclidean distance), thus they cannot summarize the dynamics of the patterns into compact models. Although StreamScan [22] is able to identify variable length patterns with HMM on data streams, it requires query models for the pattern we want to find.

In conclusion, none of the above methods satisfy all the requirements, namely, anytime incremental operation, scalability on massive streams, adaptivity, and automation for intuitive pattern discovery.

### 3 PROBLEM FORMULATION

In this section, we formally describe the problem we want to solve for automatic stream mining of sequential patterns. Table 2 lists the symbols used in this paper. The data stream we consider is a collection of  $d$ -dimensional vectors  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_t\}$ , where  $\mathbf{x}_t = \{x_i\}_{i=1}^d$  arrives at each time tick  $t$ .

Given a data stream  $X$ , our goal is to partition it into a set of  $m$  non-overlapping segments denoted by  $\mathcal{S} = \{s_1, \dots, s_m\}$ , where  $s_i$  consists of the start and end position, i.e.,  $s = \{t_s, t_e\}$ , and  $m$  denotes the number of segments.

We also want to assign each segment  $s_i$  to any one of  $r$  groups of segments (namely *regimes*), where,  $r$  is the number of regimes. Here, the regime assignment of each segment is denoted by  $\mathcal{F} = \{f(1), \dots, f(m)\}$ , where  $f(i)$  is the regime index of  $i$ -th segment (i.e.,  $1 \leq f(i) \leq r$ ).

For example, in Figure 1, the data stream consists of  $m = 8$  segments, each of which belongs to one of the  $r = 4$  regimes, e.g.,  $f(1) = 1, f(2) = 2, f(3) = 4, f(4) = 2, \dots$ .

In this paper, we assume that data streams contains two levels of dynamics, namely, *within-regime* transitions and *across-regime* transitions.

To represent dynamic time-evolving pattern in a single regime (i.e., within-regime transitions), we propose using a hidden Markov model (HMM)<sup>1</sup> [2]. An HMM is a well-known stochastic model used in many applications such as speaker recognition and the analysis of biological sequences. It assumes the system to be a Markov process with discrete hidden states. An HMM with  $k$  states is composed of initial probabilities  $\boldsymbol{\pi} = \{\pi_i\}_{i=1}^k$ , transition probabilities  $\mathbf{A} = \{a_{ij}\}_{i,j=1}^k$ , and output probabilities  $\mathbf{B} = \{b_i(\mathbf{x}_t)\}_{i=1}^k$ . Consequently, a single regime dynamics can be described as a set of parameters:  $\boldsymbol{\theta} = \{\boldsymbol{\pi}, \mathbf{A}, \mathbf{B}\}$ .

<sup>1</sup>In our setting, we assume a Gaussian distribution for the output probability, which is able to handle multi-dimensional vectors at each time tick (i.e.,  $\mathbf{B} = \{\mathcal{N}(\mu_i, \sigma_i^2)\}_{i=1}^k$ ).

**Table 2: Symbols and definitions.**

Symbol	Definition
$d$	Number of dimensions
$\mathbf{x}_t$	$d$ -dimensional vector at time tick $t$
$X$	Co-evolving data stream, i.e., $X = \{\mathbf{x}_1, \dots, \mathbf{x}_t\}$
$m$	Number of segments
$\mathcal{S}$	Segment set, i.e., $\mathcal{S} = \{s_1, \dots, s_m\}$
$\mathcal{F}$	Segment-membership, i.e., $\mathcal{F} = \{f(1), \dots, f(m)\}$
$r$	Number of regimes
$\boldsymbol{\theta}_i$	Model parameters governing $i$ -th regime
$k_i$	Number of hidden states in $\boldsymbol{\theta}_i$
$\Phi$	Regime transitions, i.e., $\Phi = \{\phi_{ij}\}_{i,j=1}^r$
$\Theta$	Model set of $r$ regimes, i.e., $\Theta = \{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_r, \Phi\}$
$C$	Candidate solution, i.e., $C = \{m, r, \mathcal{S}, \Theta, \mathcal{F}\}$

For the upper level, we assume that there are across-regime transitions between multiple distinct regimes  $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_r$ . Here, we define a *regime transition matrix*  $\Phi = \{\phi_{ij}\}_{i,j=1}^r$ , which consists of the regime transition probabilities  $\phi_{ij}$  from the  $i$ -th regime  $\boldsymbol{\theta}_i$  to the  $j$ -th regime  $\boldsymbol{\theta}_j$ . Each probability is computed as  $\phi_{ij} = n_{ij} / \sum_{s \in \mathcal{S}_i} |s|$  where  $\sum_{s \in \mathcal{S}_i} |s|$  indicates the total length of the segments belonging to regime  $\boldsymbol{\theta}_i$ , and  $n_{ij}$  shows the regime-switch count from  $\boldsymbol{\theta}_i$  to  $\boldsymbol{\theta}_j$ . Finally, the entire model parameters can be defined as  $\Theta = \{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_r, \Phi\}$ .

Finally, we adapt all of the above components for dynamic summarization on data stream  $X$ . Specifically, letting  $C = \{m, r, \mathcal{S}, \Theta, \mathcal{F}\}$  be a compact representation of  $X$ , which we refer to as a candidate solution. The problem we want to solve is as follows.

**PROBLEM 1.** Given a  $d$ -dimensional data stream  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_t\}$ , where  $t$  is the most recent time tick, continuously optimize the candidate solution  $C$ , namely,

- the number of segments  $m$  and their switching positions,  $\mathcal{S} = \{s_1, \dots, s_m\}$ ,
- the number of regimes  $r$  and their model parameters,  $\Theta = \{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_r, \Phi\}$ , and
- the segment membership of each segment,  $\mathcal{F} = \{f(1), \dots, f(m)\}$ ,

in a streaming (i.e., online) setting.

## 4 AUTOMATIC PATTERN MINING

In this section, we first present how we can define a good summary of co-evolving event streams, and then show our optimization problem, which consists of dynamically estimating a compact description of data streams in an automatic way.

### 4.1 Offline Data Compression

Here, we assume that we can access and process whole input time series  $X$  (i.e., batch processing) to simplify the discussion. Intuitively, our goal is to keep the description of  $X$  compact and reasonable using  $C$ . This approach indicates that optimal model should minimize the total description cost:

$$Cost_T(X; C) = Cost_M(C) + Cost_C(X|C), \quad (1)$$

where  $Cost_M(C)$  is the description cost of  $C$ , and  $Cost_C(X|C)$  is the data coding cost given the model  $C$ . Thus, we apply the minimum description length (MDL) principle [9], which makes it possible to evaluate how well our model compresses  $X$ . We begin by defining the two components of the total cost more concretely.

**Model description cost.** The model description cost is the number of bits needed to describe the model. If we use more powerful model architecture, the total cost becomes higher. For our full parameter set  $C$ , the number of dimensions, segments and regimes requires  $\log^*(d) + \log^*(m) + \log^*(r)$  bits<sup>2</sup>. The segment membership of the regimes requires  $Cost_M(\mathcal{F}) = m \log(r)$  bits. The length of each segment  $s$  needs  $Cost_M(\mathcal{S}) = \sum_{i=1}^m \log^* |s_i|$  bits. The model parameter set  $\Theta$  needs:  $Cost_M(\Theta) = \sum_{i=1}^r Cost_M(\theta_i) + Cost_M(\Phi)$ , where a single model  $\theta$  requires the sum of  $\log^*(k)$  bits for the number of hidden states, and  $C_F(k + k^2 + 2kd)$  for a Gaussian HMM, i.e.,  $\{\pi, A, B\}$ . Here,  $C_F$  is the floating point cost<sup>3</sup>. The regime transition matrix also requires  $Cost_M(\Phi) = C_F \cdot r^2$ .

Finally, the total model description cost  $Cost_M(C)$  is the sum of the above terms:

$$Cost_M(C) = \log^*(d) + \log^*(m) + \log^*(r) + Cost_M(\mathcal{S}) + Cost_M(\Theta) + Cost_M(\mathcal{F}) \quad (2)$$

**Data encoding cost.** Given a full parameter set  $\Theta$ , data compression using Huffman coding assigns a number of bits to each value in  $X$ , which is a negative log-likelihood. The encoding cost of  $X$  given  $\theta$  is computed by:  $Cost_C(X|\theta) = -\log P(X|\theta)$ , where  $P(X|\theta)$  is the likelihood of  $X$ . Thus, the total encoding cost is:

$$Cost_C(X|C) = \sum_{i=1}^m Cost_C(X[s_i]|\Theta) = \sum_{i=1}^m -\log(\phi_{vu}(\phi_{uu})^{|s_i|-1} P(X[s_i]|\theta_u)), \quad (3)$$

where the  $i$ -th and  $(i-1)$ -th segments are governed by the  $u$ -th and  $v$ -th regimes, respectively.  $X[s_i]$  is a subsequence of segment  $s_i$ .

## 4.2 Online Optimization Objective

Next, we focus on how to evaluate the total cost  $Cost_T(X; C)$  efficiently. In a streaming setting, we cannot store and process all historical data because its size becomes extremely large as time progresses. To search for the optimal segments/regimes incrementally, we only consider the increase in the total cost when the set of new stream elements is added to  $X$ . We roughly refer to this as the most recent subsequence  $X_c = X[t' : t]$  (i.e.,  $1 \ll t' < t$ ). Assuming that the two numbers  $m$  and  $r$  have increased to  $m_+$  and  $r_+$ , respectively, given  $X_c$  at time tick  $t$ , we define the additional description cost in Equation (1) as follows:

$$\Delta Cost_T(X_c; C) = \log^*(m_+) - \log^*(m) + \log^*(r_+) - \log^*(r) + \Delta Cost_M(\mathcal{S}) + \Delta Cost_M(\Theta) + \Delta Cost_M(\mathcal{F}) + Cost_C(X_c|\Theta). \quad (4)$$

<sup>2</sup>Here,  $\log^*$  is the universal code length for integers.

<sup>3</sup>We used  $4 \times 8$  bits in our setting.

Here  $\Delta Cost_M(\cdot)$  represents the increase in description length for new entries in the set:

$$\Delta Cost_M(\mathcal{S}) = \sum_{i=m}^{m_+-1} \log^* |s_i|, \quad (5)$$

$$\Delta Cost_M(\Theta) = \sum_{i=r+1}^{r_+} Cost_M(\theta_i) + C_F \cdot r_+^2 - C_F \cdot r^2, \quad (6)$$

$$\Delta Cost_M(\mathcal{F}) = m_+ \log(r_+) - m \log(r). \quad (7)$$

A key advantage of applying our techniques to data streams is that it can jointly estimate the number of regimes while processing. The above updating strategy is designed to flexibly deal with the number of regimes in  $X_c$ . If we represent  $X_c$  only with existing regimes, i.e.,  $r_+ = r$ , then the additional cost  $\Delta Cost_M(\Theta) = 0$ , otherwise, further cost is required corresponding to the number and complexity of the new regimes  $\theta$ . This observation indicates that more complex or duplicate regimes requires a higher cost to employ. Our objective function is likely to keep the entire model structure reasonable.

## 5 STREAMING ALGORITHMS

We now present our proposed STREAMSCOPE for the fast and automatic pattern discovery of time-evolving data streams.

### 5.1 Overview

Although we have proposed our incremental coding scheme, i.e., Equation (4), it is still difficult to globally optimize the score because the number of all possible combinations of candidate segments/regimes in the whole input  $X$  becomes extremely rich. So, how efficiently can we monitor time-evolving streams while maintaining the quality of a full model set  $C$ ?

To tackle this important problem, we carefully design an online optimization algorithm, namely STREAMSCOPE, which utilizes two sub-algorithms and maintains the model parameter set  $C$ . In STREAMSCOPE, we define the most recent subsequence  $X_c$  as the latest pattern we found, i.e.,  $X_c = X[t' : t]$  where  $t'$  and  $t$  show the last regime switching point and current time tick, respectively. Given  $X_c$ , the decision for its model update can be decided via two procedures, namely,

- (1) SEGMENTASSIGNMENT: This quickly searches the most suitable model/regime to represent arriving data  $x_t$  without any modification or addition of regimes.
- (2) REGIMEGENERATION: This considers inner/outer extensions of regimes by estimating new model(s)  $\theta$  that capture dynamical patterns in  $X_c$ .

Most importantly, since these procedures are completely separate from each other, we can decide how the entire model  $C$  updates by applying either of the outputs minimizing Equation (4). If  $X_c$  includes different dynamics from existing regimes, the model structure should be adapted to the dynamics even if it incurs the extra costs of model extensions. We describe STREAMSCOPE in detail after presenting the two algorithms in steps.

### 5.2 SegmentAssignment

As the simplest case scenario, we first assume that we know several patterns/regimes and their model parameters. Our first goal is

to identify which of the patterns is arriving at any given moment. Given a  $d$ -dimensional vector  $\mathbf{x}_t$  at time tick  $t$  and a model parameter set of  $r$  regimes  $\Theta = \{\theta_1, \dots, \theta_r, \Phi\}$ , we want to find cut points that divide  $X$  into individual patterns (i.e., segments), so that we can minimize  $Cost_C(X|\Theta)$ .

Algorithm 1 shows the overall procedure for the incremental assignment, namely, SEGMENTASSIGNMENT. In the first step, we compute the likelihood value  $P(\mathbf{x}_t|\Theta)$  as follows:

$$P(\mathbf{x}_t|\Theta) = \max_{1 \leq i \leq r} \{ \max_{1 \leq j \leq k_i} \{ p_{t;i;j} \} \}, \quad (8)$$

$$p_{t;i;j} = \max \begin{cases} \phi_{ui} \cdot \max_v \{ p_{t-1;u;v} \} \cdot \pi_{i;j} \cdot b_{i;j}(\mathbf{x}_t) \\ \phi_{ii} \cdot \max_u \{ p_{t-1;i;u} \cdot a_{i;uj} \} \cdot b_{i;j}(\mathbf{x}_t) \end{cases} \quad (9)$$

where:  $p_{t;i;j}$  is the maximum probability of state  $j$  of regime  $\theta_i$  at time tick  $t$ .  $\pi_{i;j}$  is the initial probability of state  $j$  of  $\theta_i$ .  $b_{i;j}(\mathbf{x}_t)$  is the output probability of  $\mathbf{x}_t$  for state  $j$  of  $\theta_i$ , and  $a_{i;uj}$  is the transition probability from state  $u$  to state  $j$  in  $\theta_i$ . At time tick  $t = 1$ , the probability for each regime is given by  $p_{1;i;j} = \phi_{ii} \cdot \pi_{i;j} \cdot b_{i;j}(\mathbf{x}_1)$ . Next, we retain the cut points using a candidate set  $\mathcal{L} = \{L_1, L_2, \dots\}$  for each state of all the regimes, where  $L_i$  denotes the tuple of the location  $t_i$  and the regime index  $w_i$  of the  $i$ -th cut point, which it visited (i.e.,  $L_i = \{t_i, w_i\}$ ). The entire candidate set is:

$$\mathcal{L}_{t;i;j} = \begin{cases} \mathcal{L}_{t-1;u;v} \cup \{L\} & // \text{ if regime switches} \\ \mathcal{L}_{t-1;i;j} & // \text{ otherwise} \end{cases} \quad (10)$$

where  $\mathcal{L}_{t;i;j}$  shows the candidate set for state  $j$  of  $\theta_i$  at time tick  $t$ . We can incrementally update these likelihood and candidate cut-point sets using a dynamic programming approach that maximizes the probabilities solely from previous regimes and states.

**$\gamma$ -guarantee.** If the regime that gives the maximum likelihood is switched to another regime, the time tick  $t$  is a cut point, i.e., a starting point of the new segment. However, this determinant cannot be optimum due to noise and uncertainty about similar patterns. To avoid a misassignment due to the temporal maximum likelihood, we propose  $\gamma$ -guarantee, which provides the cut-point set with a guarantee of exactness by using a time interval  $\gamma$ . If we find the first candidate cut point switching to the  $i$ -th regime at time tick  $t$ , the determinant cut point is obtained at time tick  $t + \gamma$ . The interval  $\gamma$  is typically set to a half the average length of segments.

**LEMMA 5.1.** *For a given model  $\Theta = \{\theta_1, \dots, \theta_r, \Phi\}$  and stream  $X[1 : t + \gamma]$ , the exactness of the cut point set until time tick  $t$  is guaranteed at time tick  $t + \gamma$ .*

**PROOF.**  $P(\mathbf{x}_{t+\gamma}|\Theta)$  shows the maximum likelihood at time tick  $t + \gamma$ . SEGMENTASSIGNMENT computes the likelihood of every regime each with states for each time tick. At the same time, it also reports all candidate cut points without omission. Therefore, at time tick  $t + \gamma$ , the cut-point set giving the maximum probability  $P(\mathbf{x}_{t+\gamma}|\Theta)$  is equivalent to a Viterbi path, that is, the optimal cut-point set. The cut point set until time tick  $t$  is thus exactly optimal at time tick  $t + \gamma$ .  $\square$

---

**Algorithm 1** SEGMENTASSIGNMENT ( $\mathbf{x}_t, \Theta$ )

---

**Input:** (a) New vector  $\mathbf{x}_t$  at time tick  $t$   
(b) Regime parameter set  $\Theta = \{\theta_1, \dots, \theta_r, \Phi\}$   
**Output:** (a) Number of segments  $m$   
(b) Segment set  $\mathcal{S}$   
(c) Segment membership  $\mathcal{F}$

```

1:  $m \leftarrow 0$ ;  $\mathcal{S} \leftarrow \emptyset$ ;  $\mathcal{F} \leftarrow \emptyset$ ;
2: /* Find candidate cut-points with Equations (9) and (10) */
3: for  $i = 1 : r$  do
4:   Compute  $p_{t;i;j}$  for state  $j = 1$  to  $k_i$ ;
5:   Update  $\mathcal{L}_{t;i;j}$  for state  $j = 1$  to  $k_i$ ;
6: end for
7: if the first regime transition is found then
8:    $t_\gamma \leftarrow t + \gamma$ ;
9: end if
10: if  $t = t_\gamma$  then
11:   /* Add segments into the optimal regime */
12:    $\mathcal{L}_{best} \leftarrow \arg \max_{\mathcal{L}_{t;i;j} | 1 \leq i \leq r, 1 \leq j \leq k_i} P(\mathbf{x}_t|\Theta)$ ; // Equation (8)
13:   for each cut point  $L_i = \{t_i, w_i\} \in \mathcal{L}_{best}$  do
14:     Create a new segment  $s = \{t_s, t_i\}$ ;
15:     Add  $s$  into  $\mathcal{S}$ ;  $f(m+1) = w_i$ ;  $m \leftarrow m + 1$ ;
16:      $t_s \leftarrow t_i + 1$ ;
17:   end for
18:    $\mathcal{L}_t \leftarrow \emptyset$ ;
19: end if
20: return  $\{m, \mathcal{S}, \mathcal{F}\}$ ;

```

---

### 5.3 RegimeGeneration

In this step, our goal is to find the local optimal number of segments, regimes and their model parameters thus minimizing the cost  $\Delta Cost_T(X_c; C) = \Delta Cost_M(C) + Cost_C(X_c; C)$ . REGIMEGENERATION (Algorithm 2) initially regards  $X_c$  as a single regime, i.e.,  $m_0 = 1, \mathcal{S}_0 = \{1, |X_c|\}$ , and estimates the model parameter  $\theta_0$ . Given the initial parameter set  $\{m_0, \mathcal{S}_0, \theta_0\}$ , it greedily tries to divide  $X_c$  by iterating the following two steps until a stack  $Q$  for keeping candidate regimes becomes empty.

- (1) The algorithm first finds a new candidate pair  $\{\theta_1, \theta_2, \Phi\}$ . Given the candidate regimes  $\{\theta_1, \theta_2, \Phi\}$ , it iterates finding the cut points of segments,  $\{\mathcal{S}_1, \mathcal{S}_2\}$ , by computing the maximum likelihood (i.e., Equation (8)) and re-estimating  $\{\theta_1, \theta_2, \Phi\}$ , while  $\Delta Cost_T(X_c[\mathcal{S}_0]|\theta_1, \theta_2, \Phi)$  is improved.
- (2) If the total cost of applying the candidate pair is lower than when applying only  $\theta$ , the algorithm stores the new candidate set  $\{m_1, \mathcal{S}_1, \theta_1\}$  and  $\{m_2, \mathcal{S}_2, \theta_2\}$  into  $Q$  to further divide the candidate regimes. Otherwise, it updates the parameters in  $C$  based on  $\{m_0, \mathcal{S}_0, \theta_0\}$ .

Note that we use the Baum-Welch algorithm whenever we undertake the parameter estimation of a single regime  $\theta$  in the above procedure. Since the Baum-Welch algorithm needs the number of hidden states  $k$  in a model  $\theta$  to be specified, we vary  $k = 1, 2, \dots$  while the cost,  $\Delta Cost_T(X[\mathcal{S}]|\theta)$ , can be decreased. This approach prevents overfitting and provides a reasonable  $k$  without any prior information.

**Model initialization.** When we start the iterations with a candidate regime including segments  $\mathcal{S}$ , we uniformly take several sample segments from  $X_c[\mathcal{S}]$  and estimate the model parameters  $\theta_s$  for

**Algorithm 2** REGIMEGENERATION ( $X_c$ )

---

**Input:** Subsequence  $X_c$   
**Output:** Candidate solution  $C$

- 1:  $Q \leftarrow \emptyset$ ; // Stack for candidate regimes
- 2:  $m_0 \leftarrow 1$ ;  $S_0 \leftarrow \{1, |X_c|\}$ ;  $\theta_0 \leftarrow \text{MODELESTIMATION}(X_c)$ ;
- 3: Push an entry  $\{m_0, S_0, \theta_0\}$  into  $Q$ ;
- 4: **while** stack  $Q \neq \emptyset$  **do**
- 5:   Pop an entry  $\{m_0, S_0, \theta_0\}$  from  $Q$ ;
- 6:    $\{\theta_1, \theta_2\} \leftarrow \text{MODELINITIALIZATION}(X_c[S_0])$ ; // Equation (11)
- 7:   **repeat**
- 8:     **for**  $x_t \in X_c[S_0]$  **do**
- 9:        $\{m^*, S^*, F^*\} \leftarrow \text{SEGMENTASSIGNMENT}(x_t, \theta_1, \theta_2, \Phi_{2 \times 2})$ ;
- 10:     **end for**
- 11:      $S_1 \leftarrow \{s_i \in S^* | f(i) = 1, i = 1, \dots, m^*\}$ ;
- 12:      $S_2 \leftarrow \{s_i \in S^* | f(i) = 2, i = 1, \dots, m^*\}$ ;
- 13:      $\theta_1 \leftarrow \text{MODELESTIMATION}(X_c[S_1])$ ;
- 14:      $\theta_2 \leftarrow \text{MODELESTIMATION}(X_c[S_2])$ ;
- 15:     Update regime transition matrix  $\Phi_{2 \times 2}$ ;
- 16:     **until** convergence;
- 17:     **if**  $\Delta \text{Cost}_T(X_c; S_1, S_2, \theta_1, \theta_2) < \Delta \text{Cost}_T(X_c; S_0, \theta_0)$  **then**
- 18:       /\* Recursively try to split new regimes \*/;
- 19:       Push entries  $\{m_1, S_1, \theta_1\}, \{m_2, S_2, \theta_2\}$  into  $Q$ ;
- 20:     **else**
- 21:        $S \leftarrow S \cup S_0$ ;  $\Theta \leftarrow \Theta \cup \theta_0$ ;  $r \leftarrow r + 1$ ;
- 22:        $f(i) = r$ , for each  $i = m + 1, \dots, m + m_0$ ;  $m \leftarrow m + m_0$ ;
- 23:       Update  $\Phi_{r \times r}$ ;
- 24:     **end if**
- 25:   **end while**
- 26: **return**  $C = \{m, r, S, \Theta, F\}$ ;

---

each of them. The most appropriate pair is chosen by computing the coding cost of all possible pairs:

$$\{\theta_1, \theta_2\} = \arg \min_{\theta_{s_1}, \theta_{s_2} | s_1, s_2 \in \mathcal{X}} \text{Cost}_C(X_c[S] | \theta_{s_1}, \theta_{s_2}), \quad (11)$$

where  $\mathcal{X} = \{s_1, s_2, \dots\}$  is a set of samples taken from  $X_c[S]$ .

## 5.4 StreamScope

Now, we explain how to efficiently identify regimes in a large event stream. Our final goal is to answer the questions: (a) when should we estimate model parameter  $\theta$  for new regimes and confirm their necessity for an effective pattern mining? And (b) how can we obtain the compact description  $C$  at each time tick from the outputs of the two algorithms we have discussed?

Now, we address the first question. Algorithm 3 shows the overall procedure of STREAMSCOPE. Given a vector  $x_t$  and the candidate solution  $C$ , it first executes SEGMENTASSIGNMENT to detect regime transitions. Recall that it retains candidate transition points i.e., Equation (10). STREAMSCOPE tries to estimate new regimes by using REGIMEGENERATION, at the same time as the candidates are verified, which can reveal that the current regime is switching to new regime not to one of the existing regimes.

**Online update of full parameter set.** Here, we describe how to update  $C$  based on the output of SEGMENTASSIGNMENT (i.e.,  $C_S$ ) or REGIMEGENERATION (i.e.,  $C_R$ ). Given the two candidates, STREAMSCOPE compares the extra costs, i.e., Equation (4), then decides which candidate it applies to  $C$ .

**Algorithm 3** STREAMSCOPE ( $x_t, C$ )

---

**Input:** (a) New observation  $x_t$  at time tick  $t$   
(b) Previous candidate solution  $C = \{m, r, S, \Theta, F\}$   
**Output:** Updated candidate solution  $C$

- 1:  $X_c \leftarrow X_c \cup x_t$ ;
- 2: /\* (I) Assign  $x_t$  into an existing regime \*/
- 3:  $\{m_S, S_S, F_S\} \leftarrow \text{SEGMENTASSIGNMENT}(x_t, \Theta)$ ;
- 4:  $C_S = \{m_S, r, S_S, \Theta, F_S\}$ ;
- 5: **if**  $t = t_\gamma$  **then**
- 6:   /\* (II) Generate new regime(s) \*/
- 7:    $C_R = \{m_R, r_R, S_R, \Theta_R, F_R\} \leftarrow \text{REGIMEGENERATION}(X_c)$ ;
- 8:   /\* (III) Online update according to Equation (4) \*/
- 9:   **if**  $\Delta \text{Cost}_T(X_c; C_S) < \Delta \text{Cost}_T(X_c; C_R)$  **then**
- 10:     **for**  $i = 1 : m_S$  **do**
- 11:        $w \leftarrow f(i) \in F_S$ ;  $\{t_s, t_e\} \leftarrow s_i \in S_S$ ;
- 12:        $\theta_w \leftarrow \text{MODELUPDATE}(X_c[t_s : t_e], \theta_w)$ ;
- 13:     **end for**
- 14:      $S \leftarrow S \cup S_S$ ;  $F \leftarrow F \cup F_S$ ;  $m \leftarrow m + m_S$ ;
- 15:   **else**
- 16:      $w \leftarrow f(m) \in F$ ;  $S_w \leftarrow \{s_i \in S_R | w = f(i) \in F_R\}$ ;
- 17:      $\theta_w \leftarrow \text{MODELUPDATE}(X_c[S_w], \theta_w)$ ;
- 18:      $\Theta \leftarrow \Theta \cup \Theta_R$ ;  $r \leftarrow r + r_R - 1$ ;
- 19:      $S \leftarrow S \cup S_R$ ;  $F \leftarrow F \cup F_R$ ;  $m \leftarrow m + m_R - 1$ ;
- 20:   **end if**
- 21:   Update regime transition matrix  $\Phi_{r \times r}$ ;
- 22:    $X_c \leftarrow X[S_m]$ ;  $t_\gamma \leftarrow t + \gamma$ ;
- 23: **end if**
- 24: **return**  $C = \{m, r, S, \Theta, F\}$ ;

---

- If STREAMSCOPE uses  $C_S$ , the regime parameters  $\theta$  used in  $X_c$  are updated using their own new segments, more specifically, it incrementally computes the means, variances and probabilities for hidden states based on a series of predicted states with  $\theta$  for  $X_c$ . Then, it merges new segment information  $\{m_S, S_S, F_S\}$  with  $\{m, S, F\}$  in  $C$ .
- If  $C_R$  is chosen, it requires renovation with regard to  $\Theta$  in addition to the updates on new  $m_R$  segments. Since STREAMSCOPE monitors regime switching from the latest pattern to the next,  $X_c$  definitely includes a known dynamical pattern, which denoted by a regime  $\theta_w$  in  $C$ . Therefore, the algorithm online updates  $\theta_w$  using new segments  $S_w$  in  $X_c$ . Then, it employs the other regimes in  $\Theta_R$  as members of  $\Theta$ .

In both cases, the transition matrix  $\Phi$  is also updated using the regime switching count in  $X_c$ .

**5.4.1 Theoretical analysis.** Let  $n$  be the length of a subsequence  $X_c$  and  $k$  be the maximum number of hidden states in the regimes  $\Theta$ , i.e.,  $k = \max\{k_i\}_{i=1}^r$ .

**LEMMA 5.2.** *The time complexity of STREAMSCOPE is  $O(drk^2 + dr^2k)$  at per time tick.*

**PROOF.** STREAMSCOPE first runs SEGMENTASSIGNMENT, which updates the likelihood from all the  $k$  previous states in a regime and all the other  $r - 1$  regimes for each regime and dimension. The inner state transition in a regime and regime switch require  $O(dk^2)$  and  $O(drk)$ , respectively, thus the total time complexity is  $O(drk^2 + dr^2k)$ . If the algorithm estimates regime parameter  $\theta$ , it iteratively performs the Baum-Welch algorithm and calculation of

the coding cost in  $X_c$ . This process requires  $O(\#iter \times n \times dk^2)$ , where  $\#iter$  is the number of iterations for REGIMEGENERATION. Since  $\#iter$  and  $n$  are small values, which are negligible. Thus, the complexity of STREAMSCOPE is  $O(drk^2 + dr^2k)$ .  $\square$

LEMMA 5.3. *STREAMSCOPE requires  $O(drk + rk^2 + r^2)$  space per time tick.*

PROOF. To maintain the cumulative likelihood for all  $r$  regimes, STREAMSCOPE maintains two arrays of  $r \times k$  numbers for a single trellis structure. This computation also needs  $r \times (k + k^2 + 2kd) + r^2$  space to maintain the parameters in  $\Theta$ . STREAMSCOPE also requires  $n \times d$  space for  $X_c$  to estimate new regime(s)  $\theta$  but the space is considered a small constant. Thus, the total memory required for STREAMSCOPE is  $O(drk + rk^2 + r^2)$ .  $\square$

These lemmas show that STREAMSCOPE scales quadratically with respect to the numbers  $r$  and  $k$ . However, from the experiments in Section 6, we will learn that this bound is exaggerated. In practice the numbers stay relatively small because they are decided so that they compress the data well through the effect of our proposed coding scheme Equation (4). The runtime of STREAMSCOPE therefore stays approximately constant, and is independent of the length of the data stream  $X$ .

## 6 EXPERIMENTS

In this section, we describe the experiments we undertook to evaluate the performance of our proposed method, STREAMSCOPE, with real data streams. We designed the experiments to answer the following questions:

- (1) *Q1. Effectiveness:* How successful is our method in discovering typical patterns in given input streams?
- (2) *Q2. Accuracy:* How well does our method detect the cut points and regimes?
- (3) *Q3. Scalability:* How does our method scale in terms of computational time?

We conducted our experiments on an Intel Xeon E5 2.7GHz with 64GB memory, running Linux. To learn initial model  $\Theta$  of STREAMSCOPE, we applied REGIMEGENERATION to 10% of an input data stream. We used the following three datasets.

**#MoCap** was obtained from the CMU motion capture database<sup>4</sup>, which consists of 64-dimensional vectors. We used a set of 20 streams, containing about  $m = 100$  segments in total. Each stream includes from 4 to 8 patterns. In our setting, we chose four dimensions corresponding to left-right arms and legs.

**#Bicycle** is a set of measurements from a 3-axis accelerometer fixed to a bicycle handle. The values are 25Hz and we use the absolute values for the frontal (i.e., z-axis) measurements. This dataset contains 10 streams with about  $m = 200$  segments that we labeled with five driving patterns including going straight, turning and stopping on paved/unpaved roads.

**#Workout** consists of 7-dimensional streams of unified 3-axis acceleration and 4-dimensional EMG collected at 25Hz with an arm-band for smart gesture control, considering the task of distinguishing four kinds of calisthenics: arm curl, rowing, side raise and push

up, and intervals between them. We used a set of 10 streams, containing about  $m = 100$  segments.

The values of each dataset were normalized so that they had the same mean and variance (i.e., z-normalization).

### 6.1 Q1: Effectiveness

We first demonstrate the effectiveness of the proposed method in terms of capturing intuitive patterns using real datasets.

**Motion sensor stream analysis.** Here we show that how effectively STREAMSCOPE works on a real motion sensor stream. Figure 2 (a), (b) and (c) show outputs of our method for three streams: “Exercise”, “Boxing” and “Basketball”, respectively.

Figure 2 (a) consists of four kinds of regimes related to exercise such as walking and running. Initially, the algorithm employs regime #1 for walking and starts monitoring the data stream. If it is no longer able to capture  $X$  with only existing regime, and the dynamics have changed significantly (e.g., segment #1→#2 and segment #4→#5), it adopts the segmentation by REGIMEGENERATION and increases the number of regimes to maintain. Otherwise, it adopts SEGMENTASSIGNMENT to assign new observations into a suitable regime, then updates the regime parameters incrementally. Finally, STREAMSCOPE successfully found six segments (i.e.,  $m = 6$ ) and four regimes (i.e.,  $r = 4$ ), which corresponded to meaningful human activities. Our method also successfully divides the two event streams, Figure 2 (b) and (c), into their own unique motions such as “punching” and “dribble”. Most importantly, our method needs no prior knowledge of dynamical patterns in data streams.

**Driving pattern recognition.** Figure 3 shows the original sensor stream and our online segmentation result using the #Bicycle dataset. STREAMSCOPE estimates regime #1/#2 for stable driving and regime #3/#4 for bumpy driving on a gravel road. Our method can automatically and precisely identify such typical regime patterns, and also find the sudden event “stopping”. Note that our segmentation and pattern recognition are independent of the occurrence frequency of dynamical patterns. It is beneficial that our method can deal with unexpected situations such as the traffic accidents and heart attacks.

**Analysis of fitness tracking streams.** Lastly, Figure 4 shows the analytical result for the #Workout data stream. In this realistic scenario, users introduce their own intervals between four training motions, thus our approach can robustly recognize all the main motions (regime #1/#2/#4/#5) as well as extract the other subsequences assigned to regime #3. As shown in the figure, our method requires no user intervention to handle variable length patterns, which expands its applicability of our method.

### 6.2 Q2: Accuracy

We first evaluate the quality of our method in terms of segmentation and clustering accuracy. We compare our method with pHMM [31], AutoPlait [20] and TICC [11], which are offline algorithms but state-of-the-art for simultaneous segmentation and clustering of time series. Note that pHMM requires two parameters,  $\epsilon_f$  and  $\epsilon_c$ , which correspond to the fitting error threshold. Therefore, we varied the two parameters from 0.1 to 10.0 to tune the quality of its segments and clusters. For tuning TICC, we varied the smoothness penalty  $\beta = 0, 1, \dots, 2000$ , and set the sparsity level in the model

<sup>4</sup><http://mocap.cs.cmu.edu/>

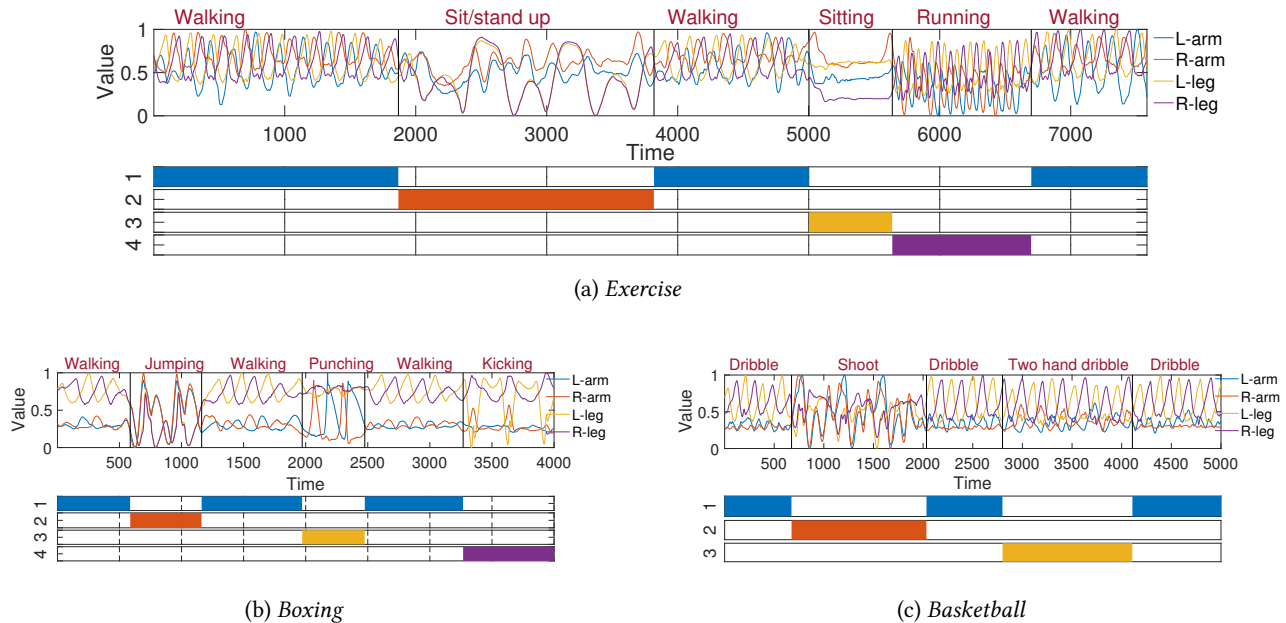


Figure 2: STREAMSCOPE is efficient: STREAMSCOPE can incrementally find regime switching (e.g., “stand up”→“walking”, “sitting”→“running” in “Exercise” data stream) for a wide variety of motion data streams. It retains the latest pattern and realises its dynamics has changed in real time. All the patterns it found are compressed into a compact model, which are updated and generated while online processing.

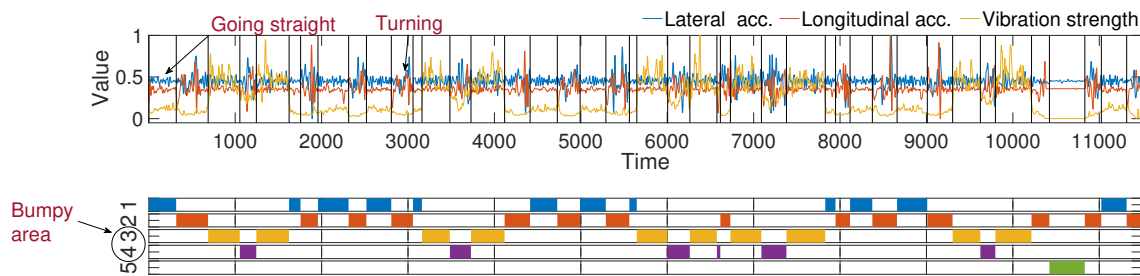


Figure 3: Driving pattern discovery of STREAMSCOPE for bicycles: Given the 3-dimensional sensor stream (top), our method finds “going straight” and “turning” on a paved/unpaved road (i.e., regimes #1/#3 and #2/#4, respectively). The approach incrementally generates regime #5 for the pattern “stopping”.

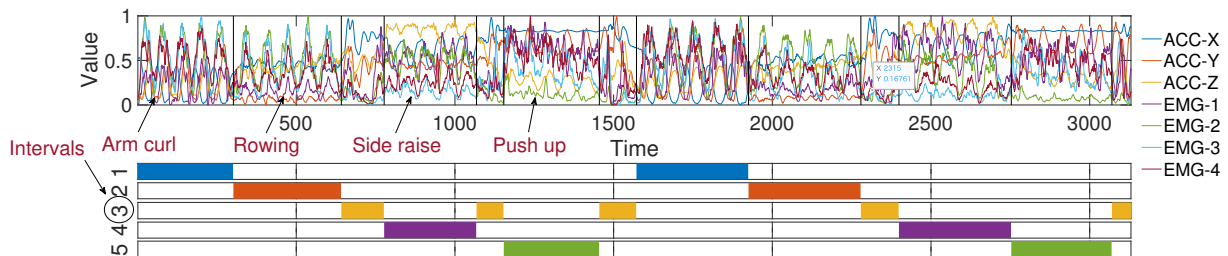


Figure 4: STREAMSCOPE can efficiently find important user activities in order, i.e., four kinds of training motions from a high dimensional data stream including 3-axis acceleration (i.e., ACC) and biological information (i.e., EMG). The algorithm also points out another pattern regime #3 (i.e., intervals) automatically.



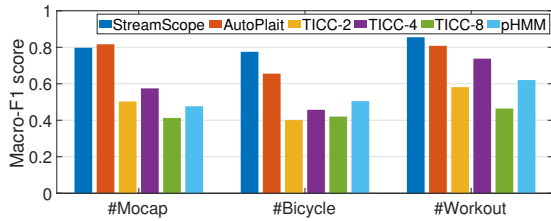


Figure 5: Segmentation accuracy with respect to macro- $F_1$  score: our online method performs as well as the other offline methods (higher is better).

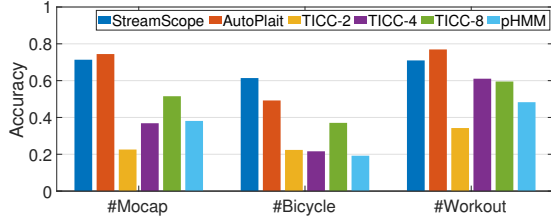


Figure 6: Clustering accuracy of STREAMSCOPE: the number of true regime labels correctly inferred with respect to the total number of observations (higher is better).

$\lambda = 10$  and window size  $w = 5$ . Moreover, since TICC needs to specify the number of clusters  $K$ , we also varied the numbers between 2, 4 and 8, namely TICC- $K$ , to test its parameter sensitivity. **Segmentation accuracy.** Figure 5 shows the segmentation accuracy with respect to the macro- $F_1$  score, which is the average of the  $F_1$  scores for all data streams. The  $F_1$  score is given by harmonic mean of the precision/recall defined as the ratio of reported correct cuts vs. the total number of reported/correct cuts, respectively. Overall, our proposed method outperforms the other offline methods because it can capture high-level patterns, i.e., regimes, regardless of its online approach thanks to our effective updating strategy for data stream summarization.

**Clustering accuracy.** Next, Figure 6 shows the clustering performance using the measure of accuracy,  $ACC$ , which is computed from the confusion matrix,  $CM$ , for prediction regime labels against true regime labels:  $ACC = \sum_{i=1} CM_{ii} / \sum_{i,j=1} CM_{ij}$ . The ideal confusion matrix will be diagonal, in which case  $ACC = 1$ . Given the incorrect number of clusters, the TICC algorithm misses important patterns, especially for *MoCap* data streams because of their variation in the number between streams. Thus, our automatic approach is effective to find multiple sequential patterns in the realistic situation where the number of patterns in streams is unknown and variable.

### 6.3 Q3: Scalability

We also compare the wall clock time of STREAMSCOPE with those of AutoPlait, TICC, and pHMM. The left of Figure 7 indicates the wall clock time of an experiment performed on a large *MoCap* dataset with four dimensions and various sequence lengths  $t$ . Note that the computation time of our method suddenly increases. This is because STREAMSCOPE updates the model parameter set as needed, however, it is significantly faster than its competitors. When our

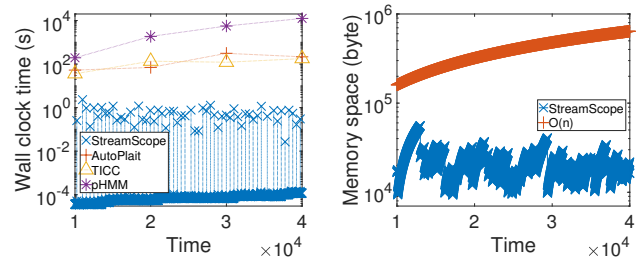


Figure 7: Scalability of STREAMSCOPE: (left) Wall clock time vs. sequence length  $t$ . it consistently surpasses its competitors. (right) Memory space consumption vs. sequence length  $t$ . it can continuously model and identify time-evolving patterns in data streams with a small space.

method does not need to update the parameters, it is also the fastest algorithm, (i.e., up to four orders of magnitude), since the likelihood computation scales constantly with the number of regimes and each state. Overall, our online algorithm is much faster than the full-batch algorithm.

The right of Figure 7 shows the space requirements of STREAMSCOPE. The requirement of our method depends on the number of regimes  $r$  and its inner states  $k$ . On the other hand, as the number of samples increases, the requirement of competitors begins to diverge to infinity; the complexity of the other competitors depends on data size  $n$ . As we expected, STREAMSCOPE only needs approximately constant space to keep track of the regime transitions, which constituted a large improvement.

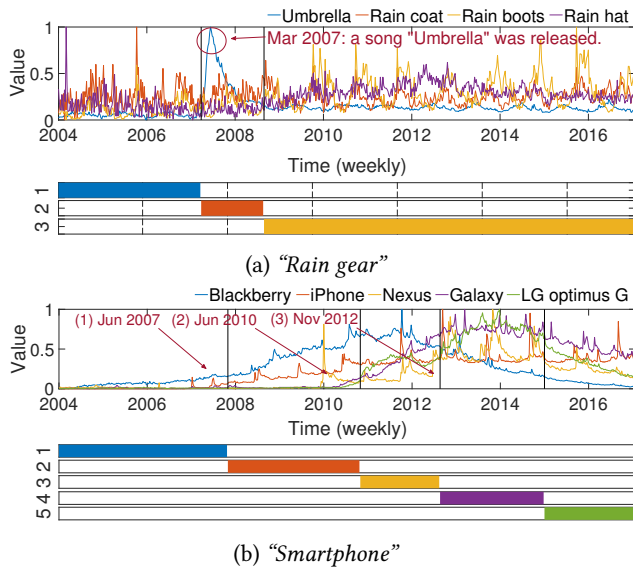
## 7 APPLICATIONS

In this section, we present an example of pattern discovery with our method as regards web-click activities, i.e., *GoogleTrends*<sup>5</sup>. The dataset we use is a multi-dimensional stream on related queries, and each dimension consists of the weekly search volumes for the query on Google over thirteen years.

**Event detection.** Figure 8 (a) shows a  $d = 4$  dimensional stream. Each dimension represents keywords on rain gear (i.e., umbrellas, raincoats, rain boots and rain hats). “Umbrella” is searched for throughout the year and the others exhibit yearly periodicity because of their seasonality. In March 2007, the search count for umbrella suddenly increased caused by an abnormal event, namely, that a song named “Umbrella” was released. STREAMSCOPE successfully detected this event as regime #2.

**Trend discovery.** Figure 8 (b) shows our result for another data stream related to some major smartphones (e.g., iPhone and Nexus). STREAMSCOPE captures the following key trends in relation to the smartphone industry caused by some latent events with the observed stream: (#1) Blackberry was used among businessmen as a pioneering smartphone, (#2) iPhone impacted and its market grew when it entered the market in June 2007, (#3) Android smartphones represented by the Nexus and Galaxy series entered the market from 2010, (#4) The competitive relationship between iOS and Android, and (#5) represents the recent global share of smartphones.

<sup>5</sup><https://trends.google.com/trends/>



**Figure 8: Application to real web data: STREAMSCOPE spots meaningful changeovers on GoogleTrends. (a) Rain gear-related keywords: it discovers one sudden event in 2007 (i.e., a hit song); (b) smartphone-related keywords: it finds five phases of smartphone market share (i.e., #1: Blackberry, #2: iPhone, #3: Android, #4: iOS vs Android, #5: Recent share).**

## 8 CONCLUSIONS

In this paper, we have proposed a streaming algorithm, STREAMSCOPE, for automatic pattern discovery in multi-dimensional streams. As the processing progresses, our algorithm partitions time-series streams into a set of segments at points where characteristics change, and then similar segments are grouped into a regime. The advantages of STREAMSCOPE are as follows:

- It is **Effective**: our experiments with several datasets show that STREAMSCOPE successfully finds segments/regimes that match human intuition.
- It is **Automatic**: it can handle arbitrary semi-infinite data streams, and automatically maintain the number  $r$  of regimes by exploiting our novel coding scheme.
- It is **Scalable**: the computational time does not depend on the length of the input streams.

## ACKNOWLEDGMENTS

The authors would like to thank the anonymous referees for their valuable comments and helpful suggestions. This work was supported by JSPS KAKENHI Grant-in-Aid for Scientific Research Number JP19J11125, JP18H03245, JP17H04681, JP16K12430, PRESTO JST.

## REFERENCES

[1] Charu C. Aggarwal, Jiawei Han, Jianyong Wang, and Philip S. Yu. 2003. A Framework for Clustering Evolving Data Streams (*VLDB*). 81–92.  
 [2] Christopher M. Bishop. 2006. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer.  
 [3] George E.P. Box, Gwilym M. Jenkins, and Gregory C. Reinsel. 1994. *Time Series Analysis: Forecasting and Control* (3rd ed.). Prentice Hall, Englewood Cliffs, NJ.

[4] Feng Cao, Martin Estert, Weining Qian, and Aoying Zhou. 2006. Density-based clustering over an evolving data stream with noise. In *SDM*. SIAM, 328–339.  
 [5] Joel Janek Dabrowski, Ashfaqur Rahman, Andrew George, Stuart Arnold, and John McCulloch. 2018. State Space Models for Forecasting Water Quality Variables: An Application in Aquaculture Prawn Farming. In *KDD*. 177–185.  
 [6] Ian Fox, Lynn Ang, Mamta Jaiswal, Rodica Pop-Busui, and Jenna Wiens. 2017. Contextual Motifs: Increasing the Utility of Motifs Using Contextual Data. In *KDD*. 155–164.  
 [7] Mohammed Gheshmoune, Mustapha Lebbah, and Hanene Azzag. 2016. State-of-the-art on clustering data streams. *Big Data Analytics* 1, 1 (2016), 13.  
 [8] Shufeng Gong, Yanfeng Zhang, and Ge Yu. 2017. Clustering Stream Data by Exploring the Evolution of Density Mountain. *PVLDB* 11, 4 (2017), 393–405.  
 [9] Peter D Grünwald, In Jae Myung, and Mark A Pitt. 2005. *Advances in minimum description length: Theory and applications*. MIT press.  
 [10] Jayavardhana Gubbi, Rajkumar Buyya, Slaven Marusic, and Marimuthu Palaniswami. 2013. Internet of Things (IoT): A Vision, Architectural Elements, and Future Directions. *Future Gener. Comput. Syst.* 29, 7 (2013), 1645–1660.  
 [11] David Hallac, Sagar Vare, Stephen Boyd, and Jure Leskovec. 2017. Toeplitz Inverse Covariance-Based Clustering of Multivariate Time Series Data. In *KDD*. 215–223.  
 [12] Ahsanul Haque, Zhuoyi Wang, Swarup Chandra, Bo Dong, Latifur Khan, and Kevin W. Hamlen. 2017. FUSION: An Online Method for Multistream Classification. In *CIKM*. 919–928.  
 [13] Eamonn J. Keogh, Selina Chu, David Hart, and Michael J. Pazzani. 2001. An Online Algorithm for Segmenting Time Series. In *ICDM*. 289–296.  
 [14] Jae-Gil Lee, Jiawei Han, and Kyu-Young Whang. 2007. Trajectory clustering: a partition-and-group framework. In *SIGMOD*. 593–604.  
 [15] Lei Li, B. Aditya Prakash, and Christos Faloutsos. 2010. Parsimonious Linear Fingerprinting for Time Series. *PVLDB* 3, 1 (2010), 385–396.  
 [16] Zitao Liu and Milos Hauskrecht. 2017. A Personalized Predictive Framework for Multivariate Clinical Time Series via Adaptive Model Selection. In *CIKM*. 1169–1177.  
 [17] Alysha M De Livera, Rob J Hyndman, and Ralph D Snyder. 2011. Forecasting time series with complex seasonal patterns using exponential smoothing. *J. Amer. Statist. Assoc.* 106, 496 (2011), 1513–1527.  
 [18] Michael Mathioudakis, Nick Koudas, and Peter Marbach. 2010. Early online identification of attention gathering items in social media. In *WSDM*. 301–310.  
 [19] Yasuko Matsubara and Yasushi Sakurai. 2019. Dynamic Modeling and Forecasting of Time-evolving Data Streams. In *KDD*. 458–468.  
 [20] Yasuko Matsubara, Yasushi Sakurai, and Christos Faloutsos. 2014. AutoPlait: Automatic Mining of Co-evolving Time Sequences. In *SIGMOD*.  
 [21] Yasuko Matsubara, Yasushi Sakurai, and Christos Faloutsos. 2015. The Web as a Jungle: Non-Linear Dynamical Systems for Co-evolving Online Activities. In *WWW*.  
 [22] Yasuko Matsubara, Yasushi Sakurai, Naonori Ueda, and Masatoshi Yoshikawa. 2014. Fast and Exact Monitoring of Co-Evolving Data Streams. In *ICDM 2014*. 390–399.  
 [23] Yasuko Matsubara, Yasushi Sakurai, Willem G. van Panhuis, and Christos Faloutsos. 2014. FUNNEL: automatic mining of spatially coevolving epidemics. In *KDD*. 105–114.  
 [24] Gianmarco De Francisci Morales, Albert Bifet, Latifur Khan, Joao Gama, and Wei Fan. 2016. IoT Big Data Stream Mining. In *KDD, Tutorial*. 2119–2120.  
 [25] Abdullah Mueen and Eamonn J. Keogh. 2010. Online discovery and maintenance of time series motifs. In *KDD*. 1089–1098.  
 [26] Spiros Papadimitriou and Philip S. Yu. 2006. Optimal multi-scale patterns in time series streams. In *SIGMOD*. 647–658.  
 [27] Tobias Preis, Helen Susannah Moat, and H. Eugene Stanley. 2013. Quantifying Trading Behavior in Financial Markets Using Google Trends. *Sci. Rep.* 3 (04 2013).  
 [28] Yasushi Sakurai, Yasuko Matsubara, and Christos Faloutsos. 2015. Mining and Forecasting of Big Time-series Data. In *SIGMOD, Tutorial*. 919–922.  
 [29] Mohammad Shokoochi-Yekta, Yanping Chen, Bilson Campana, Bing Hu, Jesin Zakaria, and Eamonn Keogh. 2015. Discovery of Meaningful Rules in Time Series. In *KDD*. 1085–1094.  
 [30] Machiko Toyoda, Yasushi Sakurai, and Yoshiharu Ishikawa. 2013. Pattern discovery in data streams under the time warping distance. *VLDB J.* 22, 3 (2013), 295–318.  
 [31] Peng Wang, Haixun Wang, and Wei Wang. 2011. Finding semantics in time series. In *SIGMOD*. 385–396.  
 [32] Byoung-Kee Yi, N.D. Sidiropoulos, Theodore Johnson, H.V. Jagadish, Christos Faloutsos, and Alexandros Biliris. 2000. Online Data Mining for Co-Evolving Time Sequences. *ICDE (2000)*, 13–22.  
 [33] Yunyue Zhu and Dennis Shasha. 2002. StatStream: Statistical Monitoring of Thousands of Data Streams in Real Time. In *VLDB*. 358–369.  
 [34] Kostas Zoumpatianos, Stratos Idreos, and Themis Palpanas. 2014. Indexing for interactive exploration of big data series. In *SIGMOD*. 1555–1566.