Fast and Multi-aspect Mining of Complex Time-stamped Event Streams

Kota Nakamura, Yasuko Matsubara, Koki Kawabata, Yuhei Umeda, Yuichiro Wada, Yasushi Sakurai





Complex Time-stamped Event Streams are Everywhere

□ A huge, online stream of time-stamped events with multiple attributes





Complex Time-stamped Event Streams are Everywhere

□ A huge, online stream of time-stamped events with multiple attributes

3 attributes (M=3)

	TimeStamp	Brand	Item category	Price
	2023-04-30-21:01	Tefal	Kettle	\$45
	2023-04-30-21:01	Bosch	Refrigerator	\$200
	2023-04-30-21:02	Samsung	TV	\$650
	2023-04-30-21:03	Sony	Portable audio	\$200
	2023-04-30-21:08	LG	TV	\$400
E-commerce	2023-04-30-21:11	Dell	Monitor	\$90
	2023-04-30-21:13	Philips	Headphones	\$190



Complex Time-stamped Event Streams are Everywhere

□ A huge, online stream of time-stamped events with multiple attributes

2 attributes (M=2)

	TimeStamp	Pick-up location	drop-up location
Image: constrained stateImage: constrained stateLocal mobility	2023-04-30-20:01	Museum C	Museum B
	2023-04-30-21:02	Cinema A	Street C
	2023-04-30-21:06	School D	Restaurant A
	2023-04-30-21:18	Office A	Station A
	2023-04-30-22:08	Street A	University D
	2023-05-01-09:11	Hotel B	Airport A
	2023-05-01-11:13	Station C	Street B



Limitations & Challenges

Complex time-stamped event streams ...

derail existing methods and even our interpretation



Because this is... **High-order tensor streams** (**) **High-dimensional** (**) **Sparse** (**) **Semi-infinite**

3rd -order tensor stream: each aspect indicates each attributes



Q. How can we summarize large, dynamic high-order tensor streams? Q. How can we see any hidden patterns, rules, and anomalies?



Q. How can we summarize large, dynamic high-order tensor streams? Q. How can we see any hidden patterns, rules, and anomalies?

Our answer is ... to focus on two types of patterns, **Regimes and Components**



Our Answer: Regimes and Components



□Summarize **semi-infinite** event stream into a handful number of segments



Our Answer: Regimes and Components

Components: Multi-aspect latent trends





Pick-up location Drop-off location

Centoral Park Grand Centoral Terminal Lower Manhattan

0.07 0.06 0.05 0.04 0.03 0.02 0.01 0.0

Timestamp (Pick-up time)

□Summarize **high-deimensional** and **sparse** events into major groups

Outline

Introduction



Algorithm

Experiments

Conclusion



Our Settings: Complex Time-stamped Event Streams

□ Event stream, which consist of {M attributes + Timestamp} → M+1th-order tensor stream $X \in \mathbb{N}^{U_1 \times \cdots \times U_M \times T}$

 \Box Continuously obtain **Current tensors** $\chi^C \in \mathbb{N}^{U_1 \times \cdots \times U_M \times \tau}$





Q1. What is the simplest mathmatical model for components?Q2. How can we represent regimes and summarize the whole stream?Q3. How can we formulate the summarization problem?

- G1. Multi-aspect component factorization
- **G2.** Compact description

G3. Problem formulation in a data compression paradigm



G1. Multi-aspect Component Factorization

Goal: to describe a high-dimensional and sparse tensor X^C as compact and interpretable model



Component matrices

Multi-aspect Component factorization

Model the generative process of events
 Assume that there are K major trends/components
 k-th component is defined by probability distribution w.r.t. M attributes and time

$$\mathbf{A}_{k}^{(m)} \in \mathbb{R}^{U_{m}}, \mathbf{B}_{t} \in \mathbb{R}^{K}$$

$$\mathbf{P}^{(m)} \sim \text{Dirichlet}(\boldsymbol{\alpha}^{(m)}), \mathbf{B}_{t} \sim \text{Dirichlet}(\boldsymbol{\beta})$$



G1. Multi-aspect Component Factorization

The generative process:



Summarize sparse activity into K components
 Mutli-aspect property: handle arbitrary-order tensors
 Online setting: capture temporal dependencies without storing tensors

15

G2. Compact description

Goal: to represent the whole stream X, containing distinct dynamical patterns



Compact description: $C = \{R, \Theta, G, S\}$ \Box the number of regimes *R* and the regime set Θ

 \Box the number of segments G and the assignments S



G3. Problem Formulation: Data Compression Paradigm

What is good summarization?

□ Minimum Description Length (MDL) principle:

"the more we can compress the data,

the more we can learn about their underlying patterns"

□ Evaluate the total encoding cost,

which is used to losslessly compress the original data streams

Summarization Problem

Find the compact description \mathcal{C} , which minimizes the total encoding cost

$$< X; C > = < C > + < X | C >$$

ModelDatacoding costcoding cost



G3. Problem Formulation: Data Compression Paradigm

□ Model Coding Cost: the number of bits needed to describe the model C□ Data Coding Cost: the coding cost of data X given the model C



Outline

Introduction

Model



Experiments

Conclusion





Given: Complex time-stamped event streams



CubeScope Finds Components (Multi-aspect latent trends/groups) Regimes (Distinct time-evolving patterns) Detects anomalies and their types





Our **CubeScope** consists of two sub-algorithms:



C-Decomposer:

 \Box incrementally monitors X^C

 \Box estimates a candidate regime θ_c

C-Compressor:

Updates the compact description C
 Measures the anomalousness of X^C



C-Decomposer





C-Decomposer is **Efficient**

□ Independ on dimensionality, i.e., it takes O(N), N: the number of events □ Conventional algorithms (e.g., ALS) are expensive for high-order tensor these scale w.r.t. all the attributes, i.e., take $O(\prod_{m=1}^{M} U_m)$



WWW'23

© 2023 Kota Nakamura et.al

Insertion-based algorithm:

Maintains a resonable description C for X and generates new regime if necessary







Compression-based anomaly detection

 \Box Higher compression cost \rightarrow higher anomalousness score

$$norm = \arg \max_{r \in R} |\mathcal{S}_{r}^{-1}|,$$

$$score(\mathcal{X}^{C}) = \langle \mathcal{X}^{C} | \theta_{norm} \rangle,$$

C-Compressor is Adaptive

- □ The concept of **normal changes** over time
 - \rightarrow Adaptively change the baseline to judge incoming tensors
- Data streams contain multiple anomalies over time
 - \rightarrow Discard anomalies from the baseline



Outline

Introduction

Model

Algorithm



Conclusion



Experimental Questions

We aim to evaluate that *CubeScope* has ...

Q1. Effectiveness:

How successfully does it discover meaningful patterns?

Q2. <u>Accuracy</u>:

How accurately does it achieve modeling, clustering, and anomaly detection?

Q3. <u>Scalability</u>:

How does it scale in terms of computational time?





Experimental Setup

12 datasets

(8 real-world datasets + 4 synthetics)

•	Dataset	The form of entry	Order				
A	Local Mobility: Ride information attributes & timestamp \rightarrow #rides						
	#1 NYC-Taxi [8] #2 Bike-Share [2]	(Pick-up/Drop-off location ID, Time) (User's age, Start/End station ID, Time)	3 4				
	E-commerce: Purchase information attributes & timestamp \rightarrow #purchases						
	#3 Jewelry [4] #4 Electronics [3]	(Price, Brand, Gem, Accessory type, Time) (Brand, Item category, Time)	4 3				
	Network traffic/intrusion: Access detail attributes & timestamp \rightarrow #accesses						
-	#5 AirForce [5]	(Protocol type, Service, Flag, Land, Duration Src/Dst bytes, Wrong fragment, Urgent, Time)	10				
	#6 External [1]	(Proto, Src/Dst IP Addr, Src/Dst Pt, Flags,Duration,Packets,Bytes, Time)	10				
AH (V)	#7 OpenStack [1]	» "	10				
·	#8 Kyoto [9]	(Src/Dst bytes, Count, Same srv/Serror/Srv serror rate, Dst host serror rate/same src port rate/srv serrors rate, Dst host count/srv count, Duration,Service,Flag,Time)	15				



Probabilistic generative models

Clustering approaches for time series, tensor, and data streams

Unsupervised anomaly detection methods



© 2023 Kota Nakamura et.al



Jewerly Dataset: 4rd-order tensor stream {*Time, Price, Brand, Gem, Accessory type*}





Q1. Effectiveness:

Jewerly Dataset: 4rd-order {*Time, Price, Brand, Gem,*





Regimes: Distinct dynamical patterns Changes in Purchase behavior

WWW'23

Components: multi-aspect latent trends

User preferences

Online Marketing Analytics



tensor stream Accessory type} \$1Kcorundum_synthetic **brandB**brandF \$700 citrine brandEbrandC brandD \$350









AirForce Dataset: 10th-order tensor stream

{Time, Protocol type, Service, Flag, Land, Duration, Src/Dst bytes, Wrong fragment, Urgent}



found Regimes that most corresponded to actual intrusions These intrusions arise over time and thus their numbers, durations, and features are unknown in advance.

their numbers, durations, and features are unknown in advance



Q2. ACCURACY: Modeling, Clustering, Anomaly Detection

"How does *CubeScope* achieve modeling, clustering, and anomaly detection?"

#5 Airforce

of components

16

#6 External

8

of components

16

60











[Anomaly Detection] AUC score: higher is better

60

40

#7 OpenStack

of components

#8 Kyoto

of components

32

60

40



CubeScope consistently outperforms its baselines

Method

NTM

TriMine

CubeScope

WWW'23

© 2023 Kota Nakamura et.al

Q3. Scalability

"How does CubeScope scale in terms of computational time?"



CubeScope is up to 312,000x faster than baselines and scales linearly



33

Outline

- Introduction
- Model
- Algorithm
- Experiments





Conclusion

Effective

- □ Introduce regimes and components
- □ Formulate the summarization problem for capturing these patterns
- Design *CubeScope* to solve the summarization problem

<u>General</u>

- □ Perform data compression, pattern discovery, and anomaly detection
- Practical in multiple domains, such as local mobility, online market analytics, and cybersecurity

Scalable

□ Fast and constant computational time w.r.t. the entire stream length and its dimensionality



Thank you!







