

CatchCore: Catching Hierarchical Dense Subtensor

Wenjie Feng(✉), Shenghua Liu, and Xueqi Cheng

CAS Key Laboratory of Network Data Science & Technology,
Institute of Computing Technology, Chinese Academy of Sciences, 100190
University of Chinese Academy of Sciences (UCAS), Beijing 100049, China.
fengwenjie@ict.ac.cn, liu.shengh@gmail.com, cxq@ict.ac.cn

Abstract. Dense subtensor detection gains remarkable success in spotting anomaly and fraudulent behaviors for the multi-aspect data (i.e., tensors), like in social media and event streams. Existing methods detect the densest subtensors flatly and separately, with an underlying assumption that these subtensors are exclusive. However, many real-world scenario usually present hierarchical properties, e.g., the core-periphery structure or dynamic communities in networks. In this paper, we propose CATCHCORE, a novel framework to effectively find the hierarchical dense subtensors. We first design a unified metric for dense subtensor detection, which can be optimized with gradient-based methods. With the proposed metric, CATCHCORE detects hierarchical dense subtensors through the hierarchy-wise alternative optimization. Finally, we utilize the minimum description length principle to measure the quality of detection result and select the optimal hierarchical dense subtensors. Extensive experiments on synthetic and real-world datasets demonstrate that CATCHCORE outperforms the top competitors in accuracy for detecting dense subtensors and anomaly patterns. Additionally, CATCHCORE identified a hierarchical researcher co-authorship group with intense interactions in DBLP dataset. Also CATCHCORE scales linearly with all aspects of tensors. Code of this paper is available at: <http://github.com/wenchieh/catchcore>.

1 Introduction

Dense subgraph and subtensor detection have been successfully used in a variety of application domains, like detecting the anomaly or fraudulent patterns (e.g., lockstep behavior, boost ratings) in social media or review sites [10, 11], identifying malicious attacks in network traffic logs or stream data [20, 22], and spotting changing gene-communities in biological networks, etc.

Several algorithms detect the densest subtensors or blocks in a flat manner [11, 20–22], i.e., remove-and-redetect one-by-one, with an underlying assumption that these subtensors are exclusive and separate. However, many real-world tensors usually present hierarchical properties, like the core-peripheral structure in networks and dynamic communities in social media. So it will be difficult to identify subtle structures (like multi-layer core) within the dense block and

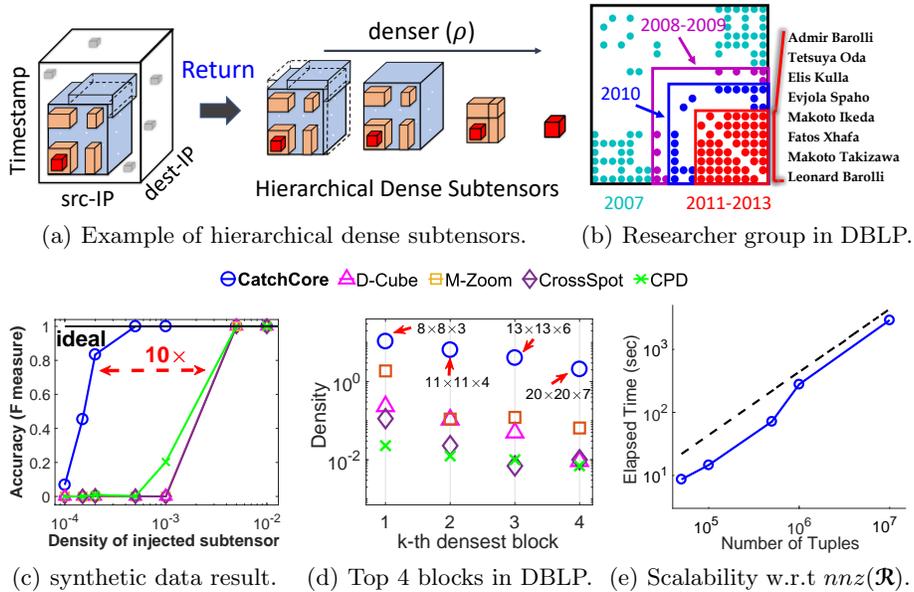


Fig. 1. Examples and CATCHCORE Performance Overview. (a) Example and workflow of hierarchical dense subtensors detection. (b) shows the detected dense co-authorship researcher group (a multi-layer core) of 20 users in DBLP. The densest block (red) last 3 years (2011-2013) containing 8 authors as the list shows, the outer hierarchies (with different colors) include other researchers and exist in various time ranges (text marked). (c) CATCHCORE outperforms competitors for detecting injected blocks in synthetic data, it achieves lower detection bound than others. (d) CATCHCORE detects dense subtensors with higher density compared with baselines for the top four densest blocks in DBLP. These blocks correspond to a hierarchical group as (b) shows. (e) CATCHCORE is linearly scalable w.r.t the number of tuples in tensor.

the relations (e.g., overlapping or inclusion) among the different blocks. Meanwhile, other methods for community detection [4, 6, 25] and dense subgraph detection [10, 19, 26] only concentrate on the plain graph.

One challenging problem is how to efficiently detect the hierarchical dense subtensors in the multi-aspect data, and Fig. 1(a) illustrates an example for the TCP dumps scenario. The network intrusion attacks dynamically changed in interacting-intensity at different stages along the time and among various hosts, resulting in a multi-layer and high-density core. So, hierarchical dense subtensor detection can help to understand the process and spot such anomaly patterns.

We propose CATCHCORE, a novel framework to detect hierarchical dense cores in multi-aspect data (i.e. tensors). Our main contributions are as follows:

- **Unified metric and algorithm:** We design a unified metric can be optimized with the gradient methods to detect dense blocks, and propose CATCHCORE for hierarchical dense core detection with theoretical guarantee and MDL based measurement.

- **Accuracy:** CATCHCORE outperforms the state-of-the-art methods in accurately detecting densest blocks and hierarchical dense subtensors in both synthetic and real-world datasets (Fig. 1(c), 1(d)).
- **Effectiveness:** CATCHCORE successfully spots anomaly patterns, including suspicious friend-connections, periodical network attacks, etc., and found a hierarchical researcher co-authorship group with heavy interactions (Fig. 1(b)).
- **Scalability:** CATCHCORE is scalable, with linear time (Fig. 1(e)) and space complexity with all aspects of tensors (Theorem 7, 8).

Reproducibility: Our open-sourced code and the data we used is available at <http://github.com/wenchieh/catchcore> (Supplement is also contained).

2 Notions and Concepts

Throughout the paper, vectors are denoted by boldface lowercases (e.g. \mathbf{x}), scalars are denoted by lowercase letters (e.g. c), and $[x] \equiv \{1, \dots, x\}$ for brevity.

Let $\mathcal{R}(A_1, \dots, A_N, C)$ be a relation consisting of N dimension attributes denoted by $\{A_1, \dots, A_N\}$, and the non-negative measure attribute $C \in \mathbb{N}^{\geq 0}$, (see the running example in supplement). We use \mathcal{R}_n to denote the set of distinct values of A_n , whose element is $a_k \in \mathcal{R}_n$. For each entry (tuple) $t \in \mathcal{R}$ and for each $n \in [N]$, we use $t[A_n]$ and $t[C]$ to denote the values of A_n and C respectively in t , i.e. $t[A_n] = a_n$ and $t[C] = c$. Thus, the relation \mathcal{R} is actually represented as an N -way tensor of size $|\mathcal{R}_1| \times \dots \times |\mathcal{R}_N|$, and the value of each entry in the tensor is $t[C]$. Let $\mathcal{R}(n, a_n) = \{t \in \mathcal{R}; t[A_n] = a_n\}$ denote all the entries of \mathcal{R} where its attribute A_n is fixed to be a_n . We define the mass of \mathcal{R} as $M_{\mathcal{B}} = \sum_{t \in \mathcal{R}} t[C]$, the volume of \mathcal{R} as $V_{\mathcal{R}} = \prod_{n=1}^N |\mathcal{R}_n|$, and the cardinality of \mathcal{R} as $D_{\mathcal{R}} = \sum_{n=1}^N |\mathcal{R}_n|$.

For a subtensor \mathcal{B} , which is composed of the subset of attributes in \mathcal{R} , is defined as $\mathcal{B} = \{t \in \mathcal{R}; t[A_n] \in \mathcal{B}_n, \forall n \in [N]\}$, i.e. the set of tuples where each attribute A_n has a value in \mathcal{B}_n . We use $\mathcal{B} \preceq \mathcal{R}$ to describe that \mathcal{B} is the subtensor of \mathcal{R} . Mathematically, for any $n \in [N]$, we can use a indicator vector $\mathbf{x} \in \{0, 1\}^{|\mathcal{R}_n|}$ to denote whether any $a_n \in \mathcal{R}_n$ belongs to \mathcal{B}_n , and $\mathbf{x}[a_n] = 1$ iff $\mathcal{B}(n, a_n) \subseteq \mathcal{R}(n, a_n)$. Thus the inclusion relationship between \mathcal{B} and \mathcal{R} can be represented with an indicator vectors collection $\mathbf{X}_{\mathcal{B}} = \{\mathbf{x}_n \in \{0, 1\}^{|\mathcal{R}_n|}; \forall n \in [N]\}$. Specially, $\mathbf{X}_{\mathbf{0}} = \{\{0\}^{|\mathcal{R}_n|}; \forall n \in [N]\}$ corresponds to **NULL** tensor (\emptyset), and $\mathbf{X}_{\mathbf{1}} = \{\{1\}^{|\mathcal{R}_n|}; \forall n \in [N]\}$ corresponds to \mathcal{R} .

Given an indicator vector $\mathbf{x} \in \{0, 1\}^{|\mathcal{R}_n|}$ for tensor \mathcal{R} , the subtensor whose n -th dimension consists of $\{a; \mathbf{x}[a] = 1, a \in \mathcal{R}_n\}$ can be denoted as $\mathcal{R} \times_n \mathbf{x}$, where “ \times_n ” is the n -mode product for a tensor and a vector¹.

3 Framework and Formulation

In this section, we propose a unified metric, which can be optimized with the gradient methods, to detect dense blocks, then we give the formal definition of the hierarchical dense subtensors detection problem.

¹ Entrywise, the n -mode product between the tensor \mathcal{R} and vector \mathbf{x} can be denoted as: $(\mathcal{R} \times_n \mathbf{x})_{i_1 \dots i_{n-1} i_{n+1} \dots i_N} = \sum_{i_n=1}^{|\mathcal{R}_n|} t(i_1, \dots, i_N, c) x_{i_n}$.

3.1 Densest Subtensor Detection Framework.

Let \mathcal{R} is an N -way tensor, and \mathcal{B} is the subtensor of \mathcal{R} defined by the indicator vectors collection $\mathbf{X}_{\mathcal{B}}$. Then the mass $M_{\mathcal{B}}$ can be represented as $M_{\mathcal{B}} = \mathcal{R} \bar{\times} \mathbf{X}_{\mathcal{B}} = \mathcal{R} \times_1 \mathbf{x}_1 \cdots \times_N \mathbf{x}_N$, where the *full-mode product* $\bar{\times}$ applies the n -mode tensor-vector product \times_n to indicator vectors along the corresponding dimension². We propose the following unified metric inspired by [24],

Definition 1 (Entry-Plenum). *Assume \mathbf{X} is an indicator vectors collection for some subtensor in \mathcal{R} , and $\phi > 0$ is a constant. Given any two strictly increasing functions g and h , the entry-plenum is defined as:*

$$f_{\phi}(\mathbf{X}) = \begin{cases} 0 & \mathbf{X} = \mathbf{X}_0, \\ g(M_{\mathbf{X}}) - \phi \cdot h(S_{\mathbf{X}}) & \text{otherwise.} \end{cases} \quad (1)$$

where $M_{\mathbf{X}}$ is the mass and $S_{\mathbf{X}}$ is the size of subtensor defined by \mathbf{X} in \mathcal{R} , and $S_{\mathbf{X}}$ can be $V_{\mathbf{X}}$, $D_{\mathbf{X}}$ or other forms.

Most popular existing subtensor density measures [11, 20, 21] can be subsumed into the above definition as

- Let $g(x) = \log x$, $h(x) = \log \frac{x}{N}$, $\phi = 1$ and $S_{\mathbf{X}} = D_{\mathbf{X}}$, $f_{\phi}(\mathcal{B})$ is equal to the *arithmetic average mass* $\rho_{ari}(\mathcal{B}) = M_{\mathcal{B}}/(D_{\mathcal{B}}/N)$.
- Let $g(x) = h(x) = \log x$, $S_{\mathbf{X}} = V_{\mathbf{X}}$, if $\phi = 1$, then $f_{\phi}(\mathcal{B})$ corresponds to *volume density* $\rho_{vol}(\mathcal{B}) = M_{\mathcal{B}}/V_{\mathcal{B}}$; and if set $\phi = \frac{1}{N}$, the $f_{\phi}(\mathcal{B})$ comes down to the *geometric average mass* $\rho_{geo}(\mathcal{B}) = M_{\mathcal{B}}/V_{\mathcal{B}}^{1/N}$.

In principle, for the entry-plenum definition, the first term $g(M_{\mathbf{X}})$ favors subtensors with the large mass, whereas the second term $-\phi \cdot h(S_{\mathbf{X}})$ acts as regularization to penalize large-size block. Thus, detecting the densest subtensor can be rewritten as the following problem under the entry-plenum metric.

Problem 2 (Densest (g, h, ϕ) -entry-plenum Subtensor). Given an N -way tensor \mathcal{R} , a constant $\phi > 0$, and a pair of increasing functions g and h , find an indicator vectors collection \mathbf{X}^* such that $f_{\phi}(\mathbf{X}^*) \geq f_{\phi}(\mathbf{X})$ for all feasible $\mathbf{X} = \{\mathbf{x}_n \in \{0, 1\}^{|\mathcal{R}_n|} : \forall n \in [N]\}$. The subtensor derived from \mathbf{X}^* is referred to be as the **Densest (g, h, ϕ) -entry-plenum Subtensor** of the tensor \mathcal{R} .

In general, finding the densest block in terms of some measure is NP-hard [2, 20], infeasible for the large dataset. Existing methods [10, 20, 21] resort to greedy approximation algorithm, which iteratively selects the local optimal subtensor from candidates based on some density measure defined in the ratio form, for scalability. Instead, our framework formulates the densest subtensor detection problem in an optimization perspective as follows, it utilizes the indicator vectors collection \mathbf{X} , which can be treated as a block-variable, to make the above problem can be solved through *block-nonlinear Gauss-Seidel (GS) method* [9] with convergence guarantee by introducing relaxation. $M_{\mathbf{X}}$ and $S_{\mathbf{X}}$ are derivable to each indicator vector under this condition, and we can use gradient-based optimization strategy for updating as long as the g and h are differentiable. Moreover, this process is linearly scalable as our proposed CATCHCORE does in Section 4.5.

² We use $\bar{\times}_{(-n)}$ to denote conducting full-mode product except the n -th mode.

3.2 Hierarchical Dense Subtensor Detection.

In practice, different dense blocks in \mathcal{R} may be overlapping or even inclusive rather than being separate or flat as many dense-block detection methods [11, 20, 21] assumed, whereas they can be described by the *hierarchical dense subtensors*. We present the following example and definition to manifest the hierarchical structure in multi-aspect data, and Fig. 1(a) gives a pictorial illustration.

Example 3 (Network Intrusion). The DARPA dataset contains 60% records labeled as anomalous belonging to (attacks), which mostly occurred in infrequent bursts, and dominant by 9 types of attacks, like **neptune**, **smurf**, and **satan** etc.. These different attacks had various intrusion intensity and burst activities among hosts, leading to discriminative dense patterns at various time ranges.

Given an N -way tensor \mathcal{R} , we want to find the dense subtensors, comprising a set of different entries, in several hierarchies. We use $\rho(\mathcal{B})$ to denote the density of subtensor \mathcal{B} , and \mathcal{B}^k as the k^{th} -hierarchy dense subtensor in \mathcal{R} . In order to find some meaningful patterns and to avoid getting identical subtensors across distinct hierarchies, we have following definition.

Definition 4 (Hierarchical Dense Subtensors (HDS-tensors)). *Given the tensor $\mathcal{B}^0 \leftarrow \mathcal{R}$ and a constant $K \in \mathbb{N}^+$. For any $k \in [K]$, the \mathcal{B}^{k-1} and \mathcal{B}^k are subtensors in two adjacent hierarchies, it is required that*

- i) **density:** the densities should be significantly different from each other, that is, for some $\eta > 1$, $\rho(\mathcal{B}^k) \geq \eta\rho(\mathcal{B}^{k-1})$.*
- ii) **structure:** subtensors in higher hierarchies are more “close-knit” (multi-layer dense core) $\mathcal{B}^k \preceq \mathcal{B}^{k-1}$, i.e., $\mathcal{B}_n^k \subseteq \mathcal{B}_n^{k-1}, \forall n \in [N]$.*

*Thus, all subtensors in K hierarchies consist of **Hierarchical Dense Subtensors**.*

Noteworthy is the fact that it is *not feasible* to recursively apply off-the-shelf dense subtensor detection methods to find HDS-tensors since they do not consider the relationship among different blocks, even if possible, it might return trivial results (e.g. identical dense subtensors across distinct hierarchies); and how to design the overall objective function to be optimized by the recursive heuristic is also not clear.

Formally, with the indicator vectors collection \mathbf{X}^k denoting the dense subtensor \mathcal{B}^k , the HDS-tensors detection problem is defined as follows.

Problem 5 (HDS-tensors Detection). **Given:** (1) the input N -way tensor \mathcal{R} , (2) the expected density ratio between two adjacent hierarchies η^3 , (3) the maximum number of hierarchies K . **Find:** the indicator vectors collections $\{\mathbf{X}^1, \dots, \mathbf{X}^r\}$, $r \leq K$ for hierarchical dense subtensors in the r significant hierarchies.

We require that $\rho(\mathbf{X}^r) \geq \eta\rho(\mathbf{X}^{r-1}) \geq \dots \geq \eta^{r-1}\rho(\mathbf{X}^1)$, and $\mathbf{X}^r \preceq \dots \preceq \mathbf{X}^1$.

In addition, we define a three-level coordinate (k, n, i) to index the indicator vectors collections, i.e., $\mathbf{X}_{(k,n,i)}$ denotes the i -th scalar element x_i of the n -th indicator vector \mathbf{x}_n in \mathbf{X}^k . Also, $\mathbf{X}_{(k,\cdot,\cdot)}$ and $\mathbf{X}_{(k,n,\cdot)}$ represent \mathbf{X}^k and indicator vector \mathbf{x}_n of \mathbf{X}^k respectively.

³ More generally, we can also set different density ratios between hierarchies rather than the fixed one parameter for specific concern.

4 Proposed Method

In this section, we propose an optimization based algorithm, CATCHCORE, to detect HDS-tensors, and provide analysis for the properties of CATCHCORE.

4.1 Optimization Based Dense Subtensor Detection.

Here, we provide an **interpretable instantiation** for the entry-plenum metric based on the volume density. With the indicator vectors collection $\mathbf{X}_{\mathcal{B}}$ of subtensor \mathcal{B} , the density is represented as $\rho_{\mathcal{B}} = \frac{M_{\mathcal{B}}}{V_{\mathcal{B}}} = \frac{M_{\mathcal{B}}}{\prod_{\mathbf{x} \in \mathbf{X}_{\mathcal{B}}} \|\mathbf{x}\|_1}$, where $V_{\mathcal{B}}$ is the product of the size of indicator vector for each mode, i.e., $\prod_{\mathbf{x} \in \mathbf{X}_{\mathcal{B}}} \|\mathbf{x}\|_1$, which equals to the total number of possible entries (including zeros).

To find the dense subtensor, if we directly maximize the density measure $\rho_{\mathcal{B}}$, however, it leads to some trivial solutions (the entries with maximum measure value, or any single entry in binary-valued tensor); while maximize the vector-based subtensor mass $M_{\mathcal{B}}$ by optimizing $\mathbf{X}_{\mathcal{B}}$ will also engender another trivial result — the \mathcal{R} itself, since no any size or volume limitation is taken into account.

Intuitively, we need to **maximize** the mass of entries while **minimize** the mass of missing entries in the block. So we proposed the optimization goal as

$$\max_{\mathbf{X}} \mathcal{F}(\mathbf{X}) = (1+p)\mathcal{R} \bar{\times} \mathbf{X} - p \prod_{\mathbf{x}_n \in \mathbf{X}} \|\mathbf{x}_n\|_1; \text{ s.t. } \mathbf{x}_n \in \{0, 1\}^{|\mathcal{R}_n|}, \forall n \in [N]. \quad (2)$$

where $p > 0$ is the penalty parameter, and $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$.

The rationale behind above definition is that each existing entry t in the resultant subtensor contributes $t[C]$ as itself to $\mathcal{F}(\mathbf{X})$, while each missing one \tilde{t} is penalized by p (i.e. $\tilde{t}[C] = -p$). In this way, the objective function maximize the total mass in the resultant subtensor while minimizing the total penalty of the missing entries. Moreover, it is also an instantiation of densest (g, h, ϕ) -entry-plenum subtensor by setting $g(x) = h(x) = x$, $\phi = p/(1+p)$, and $S_{\mathbf{X}} = V_{\mathbf{X}}$.

The optimization of the objective function $\mathcal{F}(\cdot)$ is an *NP-hard* problem due to the combinatorial nature stemming from the binary constraints of \mathbf{x}_n . So, we relax these constraints from a 0-1 integer programming to a polynomial-time solvable linear programming problem, i.e., $\mathbf{0}^{|\mathcal{R}_n|} \leq \mathbf{x}_n \leq \mathbf{1}^{|\mathcal{R}_n|}$, which represents the probability that the slice $\mathcal{R}(n, a_n)$ belonging to the resultant dense block. Finally, only the attribute value with probability exactly 1 will be selected.

4.2 Hierarchical Dense Subtensors Detection.

Based on the optimization formulation for finding dense subtensor in one hierarchy, intuitively, we maximize the objective function in Eq.(2) for each hierarchy to detect K hierarchical dense subtensors, i.e. to maximize $\sum_{k=1}^K \mathcal{F}(\mathbf{X}^k)$, and also consider aforementioned prerequisites of HDS-tensors. The *density* constraint is represented as $\rho_{\mathbf{X}^{k+1}} \geq \eta \rho_{\mathbf{X}^k}$ for the k^{th} hierarchy with density increase

ratio $\eta > 1$; and for the *structure* requirement ($\mathcal{B}^{k+1} \preceq \mathcal{B}^k \preceq \mathcal{B}^{k-1}$), we impose additional constraints on indicator vectors to prevent identical results, as $\mathbf{X}_{(k+1,n,\cdot)} \leq \mathbf{X}_{(k,n,\cdot)} \leq \mathbf{X}_{(k-1,n,\cdot)}, \forall n \in [N]$. We assume $\mathbf{X}^0 = \mathbf{X}_1, \mathbf{X}^{K+1} = \mathbf{X}_0$. Consequently, the overall optimization formulation is defined as follows,

$$\begin{aligned} \max_{\mathbf{X}^1, \dots, \mathbf{X}^K} \quad & \sum_{k=1}^K \mathcal{F}(\mathbf{X}^k) \\ \text{s.t.} \quad & \rho_{\mathbf{X}^{h+1}} \geq \eta \rho_{\mathbf{X}^h}, \mathbf{X}_{(h+1,n,\cdot)} \leq \mathbf{X}_{(h,n,\cdot)} \leq \mathbf{X}_{(h-1,n,\cdot)}. \quad \forall h \in [K]; n \in [N]. \end{aligned}$$

Obviously, this *bound-constrained multi-variable nonlinear programming (BMV-NLP)* optimization problem is non-convex and numerically intractable to solve. So we take following relaxation for the constraint to make it to be a regularization term in the objective function. Let $d^{k-1} = \eta \rho_{\mathbf{X}^{k-1}}$ which is a constant w.r.t \mathbf{X}^k , so the regularization term can be written as (also with entry-plenum form) $\mathcal{G}(\mathbf{X}^k) = \mathcal{R} \bar{\times} \mathbf{X}^k - d^{k-1} \prod_{n=1}^N \|\mathbf{X}_{(k,n,\cdot)}\|_1$. Thus, the resultant objective function with relaxation constraints for HDS-tensors detection is given by

$$\begin{aligned} \max_{\mathbf{X}^1, \dots, \mathbf{X}^K} \quad & \sum_{k=1}^K \mathcal{F}(\mathbf{X}^k) + \lambda \sum_{j=2}^K \mathcal{G}(\mathbf{X}^j) \\ \text{s.t.} \quad & \mathbf{X}_{(h+1,n,\cdot)} \leq \mathbf{X}_{(h,n,\cdot)} \leq \mathbf{X}_{(h-1,n,\cdot)}. \quad \forall h \in [K]; n \in [N]. \end{aligned} \quad (3)$$

where the parameter λ controls the importance of the regularization term.

4.3 Optimization Algorithms.

In this section, we first explain the optimization techniques, and then present the algorithm CATCHCORE (as Algorithm 1 summarized) to solve the problem.

Using the programming methods to solve the BMV-NLP optimization problem, the objective in Eq.(3) is a non-convex and higher-order bounded function with respect to each indicator vectors collection $\mathbf{X}_{(k,\cdot,\cdot)}$, which allows us to apply the alternative update method where we fix all variables of other hierarchies as constants except the current collections in each iteration; the similar strategy is also used to update each indicator vector $\mathbf{X}_{(k,n,\cdot)}$ alternatively.

Based on the *structure* constraints, for any dimension $n \in [N]$, the feasible solution $\mathbf{X}_{(k,n,\cdot)}$ is bounded by the indicator vectors $\mathbf{X}_{(k-1,n,\cdot)}$ and $\mathbf{X}_{(k+1,n,\cdot)}$ of two adjacent hierarchies in the high-dimensional space. we relax these constraints to $\mathbf{0}^{|\mathcal{R}_n|} \leq \mathbf{X}_{(k,n,\cdot)} \leq \mathbf{X}_{(k-1,n,\cdot)}$ for optimizing, thus we can obtain $\mathbf{X}^1, \mathbf{X}^2, \dots, \mathbf{X}^K$ in order. That is, we first get \mathbf{X}^1 with the constraints $\{\mathbf{0}^{|\mathcal{R}_n|} \leq \mathbf{X}_{(1,n,\cdot)} \leq \mathbf{1}^{|\mathcal{R}_n|}, \forall n \in [N]\}$ by ignoring the constraints of other variables in other \mathbf{X}^k s, then we obtain \mathbf{X}^2 based on the result at the first step under the constraints $\{\mathbf{0}^{|\mathcal{R}_n|} \leq \mathbf{X}_{(2,n,\cdot)} \leq \mathbf{X}_{(1,n,\cdot)}, \forall n \in [N]\}$ and also ignore other constraints. In this way, we can solve the K dense subtensors detection subproblems hierarchy-by-hierarchy. Technically, we adopt trust-region approach[5] to solve

Algorithm 1 CATCHCORE for HDS-tensors detection

Input: (1) the N -way tensor: \mathcal{R} (2) the maximum number of hierarchies: K
(3) the penalty value for each missing entry: p
(4) the density ratio between two adjacent hierarchies: η
(5) the regularization parameter: λ (6) maximum number of iterations: t_{\max}

Output: The dense subtensors indicator vector collections: $\{\mathbf{X}^1, \dots, \mathbf{X}^r\}$.

- 1: initialize $\mathbf{X}^1, \dots, \mathbf{X}^K$ as $\mathbf{X}_{(k,n,\cdot)}^{\text{init}}$, and $t \leftarrow 1, r \leftarrow 1$
- 2: **while** $t \leq t_{\max}$ **and** Eq. (5) is not satisfied **do** \triangleright Gauss-Seidel method
- 3: **for** $k \leftarrow 1 \dots K$ **do** \triangleright for the k^{th} hierarchy
- 4: **for** $n \leftarrow 1 \dots N$ **do** \triangleright for the n^{th} dimension
- 5: $\mathbf{x}_n^k \leftarrow \text{ONEWAYOPT}(\mathbf{x}_n^k)$
- 6: update \mathbf{X}^k
- 7: $t \leftarrow t + 1$
- 8: **while** $r \leq K$ **do** \triangleright select significant subtensors
- 9: $\mathcal{S} = \{\mathbf{X}_{(r,n,\cdot)}; \max \mathbf{X}_{(r,n,\cdot)} < 1, \forall n \in [N]\}$
- 10: **if** $\mathcal{S} \neq \emptyset$ **then**
- 11: **break** \triangleright no significant subtensors for hierarchies $> r$
- 12: **else:**
- 13: $r \leftarrow r + 1$
- 14: **return** the resultant r indicator vector collections $\{\mathbf{X}^1, \dots, \mathbf{X}^r\}$.

each nonlinear box-constrained programming subproblem. We rewrite the optimization problem in Eq.(3) as,

$$\begin{aligned}
\min_{\mathbf{X}^1, \dots, \mathbf{X}^K} f(\mathbf{X}^1, \dots, \mathbf{X}^K) &= -(1+p)\mathcal{R} \bar{\times} \mathbf{X}_{(1,\cdot,\cdot)} + p \prod_{n=1}^N \|\mathbf{X}_{(1,n,\cdot)}\|_1 \\
&\quad - (1+p+\lambda) \sum_{k=2}^K \mathcal{R} \bar{\times} \mathbf{X}_{(k,\cdot,\cdot)} + \sum_{k=2}^K (p+\lambda d^{k-1}) \prod_{n=1}^N \|\mathbf{X}_{(k,n,\cdot)}\|_1 \quad (4) \\
\text{s.t. } &\mathbf{X}_{(h+1,n,\cdot)} \leq \mathbf{X}_{(h,n,\cdot)} \leq \mathbf{X}_{(h-1,n,\cdot)}, \quad \forall h \in [K]; n \in [N].
\end{aligned}$$

We use the alternative projected gradient descent[15, 16] method that is simple and efficient to solve the optimization problem. For any dimension n , we denote the gradient of subtensor mass $M_{\mathcal{B}}$ w.r.t \mathbf{x}_n as $\nabla_{\mathbf{x}_n} M_{\mathcal{B}} = \mathcal{R} \bar{\times}_{(-n)} \mathbf{X}_{\mathcal{B}}$, so the gradient of $f(\cdot)$ w.r.t \mathbf{x}_n^1 ($\mathbf{X}_{(1,n,\cdot)}$) and \mathbf{x}_n^k ($\mathbf{X}_{(k,n,\cdot)}, k \geq 2$) are

$$\begin{aligned}
\nabla_{\mathbf{x}_n^1} f &= -(1+p)\nabla_{\mathbf{x}_n^1} M_{\mathbf{X}^1} + p \prod_{\mathbf{x}_n \in \mathbf{X}^1 / \{\mathbf{x}_n^1\}} \|\mathbf{x}_n\|_1 \mathbf{1}, \\
\nabla_{\mathbf{x}_n^k} f &= -(1+p+\lambda)\nabla_{\mathbf{x}_n^k} M_{\mathbf{X}^k} + (p+\lambda d^{k-1}) \prod_{\mathbf{x}_n \in \mathbf{X}^k / \{\mathbf{x}_n^k\}} \|\mathbf{x}_n\|_1 \mathbf{1}
\end{aligned}$$

where $\mathbf{1}$ is a $|\mathcal{R}_n|$ -dimensional all-ones vector. Let \mathbf{x}_n be the current iterator vector of any k^{th} hierarchy in the projected gradient approach, the new iterator is given by $\tilde{\mathbf{x}}_n = P(\mathbf{x}_n - \alpha \nabla_{\mathbf{x}_n} f)$ update rule. Here, the operator $P(\cdot)$ projects the vector back to the bounded feasible region, $-\nabla_{\mathbf{x}_n} f$ is the gradient-related

search direction and α is a step size computed e.g., by means of an Armijo step size strategy. In the case of an Armijo's rule line search, a good step size α is chosen until $f(\cdot, \mathbf{X}_{\mathbf{x}_n \rightarrow \tilde{\mathbf{x}}_n}^k, \cdot) - f(\cdot, \mathbf{X}^k, \cdot) \leq \sigma (\nabla_{\mathbf{x}_n} f)^T (\tilde{\mathbf{x}}_n - \mathbf{x}_n)$ is satisfied, where $\mathbf{X}_{\mathbf{x}_n \rightarrow \tilde{\mathbf{x}}_n}^k$ means replacing the indicator vector \mathbf{x}_n with the updated version $\tilde{\mathbf{x}}_n$ in \mathbf{X}^k , and a common choice of the parameter σ ($0 < \sigma < 1$) is 0.01. Thus we can alternatively update each indicator vector \mathbf{x}_n for the current k^{th} -hierarchy, and the details are listed in the ONEWAYOPT algorithm in supplement.

We propose CATCHCORE to solve the optimization problem in Eq.(4). First, we initialize the indicator vectors collection with rules that the probabilities of selecting the slices $\mathcal{R}(n, a_n)$ (i.e. $\mathbf{X}_{(k,n,i)}^{\text{init}}$) are 0.5 in the 1st hierarchy and 0.01 in other hierarchies. In this way, we can fairly avoid some trivial results. To make the solution to be close to the stationary point regarding convergence, we apply the following common condition as a stop criteria for the bounded-constrained optimization method besides the limitation for total iterations t_{\max} .

$$\left\| \left\{ \nabla_{\mathbf{x}_n^k}^P f; \forall n, k \right\} \right\|_2 \leq \epsilon \left\| \left\{ \nabla_{\mathbf{X}_{(k,n,\cdot)}^{\text{init}}}^P f; \forall n, k \right\} \right\|_2, \quad (5)$$

where $\nabla_{\mathbf{x}_n^k}^P f$ is the *elementwise projected gradient* defined as (the i -th element)

$$(\nabla_{\mathbf{x}_n^k}^P f)_i = \begin{cases} \min(0, (\nabla_{\mathbf{x}_n^k}^P f)_i) & \text{if } \mathbf{X}_{(k,n,i)} = \mathbf{X}_{(k+1,n,i)}, \\ (\nabla_{\mathbf{x}_n^k}^P f)_i & \text{if } \mathbf{X}_{(k+1,n,i)} < \mathbf{X}_{(k,n,i)} < \mathbf{X}_{(k-1,n,i)}, \\ \max(0, (\nabla_{\mathbf{x}_n^k}^P f)_i) & \text{if } \mathbf{X}_{(k,n,i)} = \mathbf{X}_{(k-1,n,i)}. \end{cases}$$

Then CATCHCORE calls ONEWAYOPT to alternatively update all the indicator vector for each dimension and hierarchy iteratively. In final, we only select these significant subtensors (the top r of K hierarchies) to return (Line 8 – 13).

4.4 Parameters Evaluation using MDL.

The penalty value p for missing entries controls the resultant lowest density, and the ratio parameter η affects density-diversity and the number of hierarchies in final. Thus, it is a challenging problem to set them appropriately or evaluate the quality of detection result under some parameter configuration, especially in the un-supervised application. We propose to measure the result w.r.t different parameter settings based on the Minimum Description Length (MDL). In the principle manner, we compute the number of bits needed to encode the tensor \mathcal{R} with detected hierarchical dense subtensors for selecting the best model (parameters), achieving the shortest code length. Intuitively, the less missing entries and the more accurate of detecting hierarchies will lead to fewer bits needed to encode \mathcal{R} in a lossless compression employing the characterization.

For the indicator vector $\mathbf{X}_{(k,n,\cdot)}$ ($k \in [K]$, $n \in [N]$), we can adopt Huffman or arithmetic coding to encode the binary string, which formally can be viewed as a sequence of realizations of a binomial random variable X . Due to $\mathbf{X}_{(k,n,\cdot)} \leq \mathbf{X}_{(k-1,n,\cdot)}$, we only consider the overlapping part $\tilde{\mathbf{x}}_n^k = \{\mathbf{X}_{(k,n,i)}; \mathbf{X}_{(k-1,n,i)} = 1, \forall i \in [|\mathcal{R}|_n]\}$ to avoid redundant encoding of 0s. We denote the entropy of

indicator vector \mathbf{x} as: $H(\mathbf{x}) = -\sum_{q \in \{0,1\}} P(X=q) \log P(X=q)$, where $P(X=q) = n_q / \|\mathbf{x}\|_1$ and n_q is the number of q in \mathbf{x} . The description length of \mathbf{X}^k is defined as $L(\mathbf{X}^k) = \sum_{n=1}^N (\log^* \|\mathbf{X}_{(k,n,\cdot)}\|_1 + \|\mathbf{X}_{(k-1,n,\cdot)}\|_1 \cdot H(\bar{\mathbf{x}}_n^k))$.⁴

Assume that $\mathbf{X}^{K+1} = \mathbf{X}_0$, For the dense subtensor \mathcal{B}^k defined by \mathbf{X}^k , we only need to encode the entries in $\bar{\mathcal{B}}^k = \mathcal{B}^k - \mathcal{B}^{k+1}$ due to $\mathcal{B}^{k+1} \preceq \mathcal{B}^k$, based on some probability distribution. For the entry $t \in \bar{\mathcal{B}}^k$, specifically, if $t[C] \in \{0,1\}$, t is sampled from binomial distribution; and if $t[C] \in \mathbb{N}^{\geq 0}$, we instead model the data by using the *Poisson* distribution [11] parameterized by the density of $\bar{\mathcal{B}}^k$, i.e. $\rho_{\bar{\mathcal{B}}^k}$. Therefore the code length for encoding $\bar{\mathcal{B}}^k$ is

$$L(\bar{\mathcal{B}}^k) = - \sum_{q \in \{t[C]; t \in \bar{\mathcal{B}}^k\}} n_q \cdot \log P(X=q) + C_{para},$$

where $P(X=q)$ is the probability of q in the probability distribution function \mathbf{P} , and C_{para} is for encoding the parameters of \mathbf{P} (like the mean in Poisson).

As for the residual part $\mathcal{R} = \mathcal{R} - \mathcal{B}^1$, we use Huffman coding to encode its entries considering the sparsity and discrete properties, the code length is denoted as L_e .

Putting all together, we can write the total code length for representing the tensor \mathcal{R} with the resultant K hierarchies indicator vectors collections as:

$$L(\mathcal{R}; \mathbf{X}^1, \dots, \mathbf{X}^K) = \log^* K + \sum_{n=1}^N \log^* |\mathcal{R}_n| + \sum_{k=1}^K L(\mathbf{X}^k) + L(\bar{\mathcal{B}}^k) + L_e. \quad (6)$$

To get the optimal parameters, we can heuristically conduct a grid search over possible values and pick the configuration that minimizes MDL. We demonstrate that the parameters according to the MDL principle results in optimal quality of detecting HDS-tensors, and the search space is limited.

4.5 Algorithm Analysis.

In this section, we provide the analysis for the convergence, the time and space complexity of CATCHCORE. The details of proofs refer to the supplement.

Lemma 6 states the convergence properties of the gradient method for our CATCHCORE.

Lemma 6. *CATCHCORE algorithm converges to a critical point.*

Theorem 7 states the time complexity of CATCHCORE algorithm, which is linear with K , N , and $nnz(\mathcal{R})$ — the number of non-zero entries in \mathcal{R} . And the space complexity is given in Theorem 8.

Theorem 7 (Worst-case Time Complexity). *Let t_{als} be the number of iterations for Armoji's line search used in the ONEWAYOPT Algorithm for updating any indicator vector, the worst-case time complexity of the CATCHCORE Algorithm 1 is $O(K \cdot N \cdot t_{max} \cdot t_{als} \cdot (nnz(\mathcal{R}) + c \cdot D_{\mathcal{R}}))$.*

⁴ $\log^* x$ is the universal code length for an integer x . [18].

Table 1. Summary of real-world datasets.

Name	Size	$D_{\mathcal{R}}$	$nnz(\mathcal{R})$	Category
Android	$1.32\text{M} \times 61.3\text{K} \times 1.28\text{K} \times 5$	1.38M	2.23M	Rating
BeerAdvocate	$26.5\text{K} \times 50.8\text{K} \times 1,472 \times 1$	78.7K	1.07M	
StackOverflow	$545\text{K} \times 96.7\text{K} \times 1,154 \times 1$	643K	1.30M	
DBLP	$1.31\text{M} \times 1.31\text{M} \times 72$	2.63M	18.9M	Social network
Youtube	$3.22\text{M} \times 3.22\text{M} \times 203$	6.45M	18.5M	
DARPA	$9.48\text{K} \times 23.4\text{K} \times 46.6\text{K}$	79.5K	522K	TCP dumps
AirForce	$3 \times 70 \times 11 \times 7.20\text{K} \times 21.5\text{K} \times 512 \times 512$	39.7K	863K	

Theorem 8 (Memory Requirements). *The amount of memory space required by CATCHCORE is $O(nnz(\mathcal{R}) + 2K \cdot D_{\mathcal{R}})$.*

Parameter Analysis: For the maximum significant hierarchies $K_{\max} = \log_{\eta}(\frac{\max(\mathcal{R})}{\rho_{\mathcal{R}}})$, where $\max(\mathcal{R})$ is the maximum value of measure attributes of \mathcal{R} . In practice, we have following observations to ensure the efficiency of CATCHCORE,

- $nnz(\mathcal{R}) \gg D_{\mathcal{R}}$, $K \ll K_{\max}$, i.e., there is only few significant hierarchies;
- $t < t_{\max}$, i.e., early stopping for iterations;
- a small t_{als} , i.e., few iterations for searching the step size.

and the dimension-update (Line 4) could be solved separately, a situation suitable for parallel environment.

5 Experiments

We design experiments to answer the following questions:

- **Q1. Accuracy:** How accurately does CATCHCORE detect HDS-tensors? Does the MDL evaluation select the optimal parameters?
- **Q2. Pattern and anomaly detection:** What patterns does CATCHCORE detect in real-world data? What is behavior of the detected anomalies?
- **Q3. Scalability:** Does CATCHCORE scale linearly with all aspects of data?

5.1 Experimental settings.

Baselines: We selected several state-of-the-art methods for dense-block detection as the baselines, including D-CUBE [21], M-ZOOM [20], CROSSSPOT [11], and CP Decomposition (CPD) [12]. In all experiments, a sparse tensor format was used for efficient memory usage, and the ρ_{ari} and ρ_{geo} were used for D-CUBE and M-ZOOM; we used a variant of CROSSSPOT which maximizes the same density metrics and used the CPD result for seed selection as did in [20].

Data: Table 1 summarizes the real-world datasets in our experiments. In *Rating* category, data are 4-way tensors (*user, item, timestamp, rating*), where entry values are the number of reviews. In *Social network*, data are 3-way tensors (*user, user, timestamps*), where entry values are the number of interactions

Table 2. Hierarchical subtensors detection results for BeerAdvocate dataset.

K	Injected Density	H1				H2					H3				
		CC	D/M	CS	CPD	CC	D	M	CS	CPD	CC	D	M	CS	CPD
2	0.05 + 0.01	1	1	1	0.99	1	0.12	0.04	0	0.08					
	0.2 + 0.1	1	0	0.24	0.24	1	1	1	1	0.99					
	0.7 + 0.1	1	0	0.33	0.33	1	1	1	1	0.52					
	1 + 0.05	1	0	0.78	0.78	1	1	1	1	0.99					
	1 + 0.2	1	0	0.79	0.79	1	1	1	1	0.99					
3	0.2 + 0.1 + 0.05	1	0	0.26	0.26	1	0	0	0	0	1	1	1	1	0.98
	1 + 0.2 + 0.01	1	0	0	0	1	0.99	1.0	1.0	0.22	1	0.44	0.43	0.75	0.85
	1 + 0.7 + 0.2	1	0	0	0	1	0	0	0.98	0.98	1	1	1	1	0.99

* The abbreviations mean that CC: CATCHCORE, D: D-CUBE, M: M-ZOOM, CS: CROSSPOT.

* The injected shape w.r.t density is: 0.01: $1K \times 800 \times 15$, 0.05: $800 \times 600 \times 10$, 0.1: $500 \times 500 \times 5$, 0.2: $300 \times 300 \times 5$, 0.7: $200 \times 100 \times 2$, 1: $100 \times 80 \times 1$.

(co-authorship / favorite). *DARPA* is the TCP dumps represented with a 3-way tensor (*source IP*, *destination IP*, *timestamps*), where entry values are the number of connections between two IP addresses (hosts). *AirForce* is also a network intrusion dataset, which is represented as a 7-way tensor (*protocol*, *service*, *src_bytes*, *dst_bytes*, *flag*, *host_count*, *src_count*, *#connections*). Timestamps are in minutes for *DARPA*, in dates for *Ratings* and *Youtube*, and in years for *DBLP*.

5.2 Q1. Accuracy of CatchCore

We compare how accurately each method detects injected dense subtensors in the synthetic and real-world datasets.

We randomly and uniformly generated a $5K \times 5K \times 2K$ 3-way tensor \mathcal{R} with a density of $3 \cdot 10^{-6}$. Into \mathcal{R} , one $200 \times 300 \times 100$ block is injected with distinct density, for testing the detection bound of each method. Fig. 1(c) demonstrated that CATCHCORE effectively detects block as low as a tenth of the density that the best baselines detect, which means that our method can spot such fraudsters with more adversarial effort.

We then injected K subtensors with different size and density into BeerAdvocate and synthetic tensor in a hierarchical manner. Table 2 lists the result for former (the result for synthetic data is listed in the supplement), where H1 is the densest or the first subtensor detected by methods, and the density (order) decreases (increases) from H1 to H3, and the information of injected blocks is listed at the bottom. CATCHCORE can accurately detect all injected subtensors in various hierarchies, size and density diversity, and consistently outperforms other baselines which fail to accurately detect or even miss at least one block. D-CUBE and M-ZOOM have similar accuracy (except some special cases as highlighted), they can not identify the structure of dense blocks, leading to some of the spars-est or densest injected blocks are missed or overwhelmed by large-volume result.

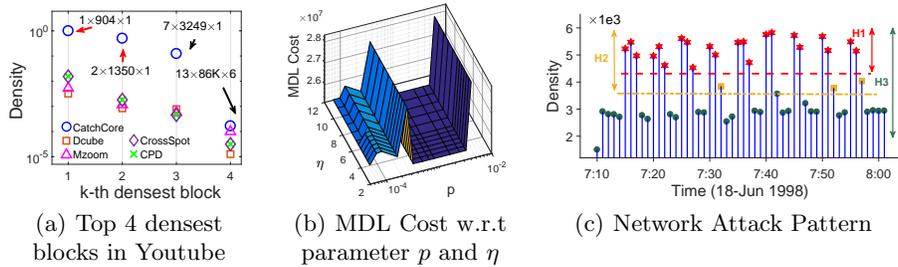


Fig. 2. (a) CATCHCORE detects higher-density blocks containing suspicious behaviors. (d) The optimal hierarchies achieve the lowest MDL cost w.r.t p and η . CATCHCORE can obtain optimal hierarchical dense subtensors for a wide parameters range. (c) Detected hierarchical network intrusion attack happened on June 18, 1998 in DARPA dataset.

CROSSSPOT and CPD also do not find hierarchical dense subtensors accurately. Similar conclusions can be drawn for the synthetic dataset.

Dense blocks in real data: We apply CATCHCORE for various real-world datasets, and measure the density instead of the mass to avoid the trivial results since that the blocks with higher density contain interesting patterns w.h.p. Fig. 1(d), 2(a) only show the densities of top four densest blocks found by the methods for DBLP and Youtube dataset. CATCHCORE spots denser blocks for each data, where it consistently outperforms the competitors for all blocks.

MDL-based Evaluation We evaluate the utility of our MDL criterion for measuring the quality of hierarchies returned by CATCHCORE under different parameters. The BeerAdvocate data with 3 hierarchical injected dense blocks (as the $K = 3$ case in Table 2) is used, we computed the MDL cost for the detection result by varying a pair of parameters p and η , the result is shown in Fig. 2. The optimal hierarchies achieve the lowest MDL cost w.r.t p and η . In addition, our model can obtain optimal results for a wide range of parameters.

5.3 Q2. Pattern and Anomaly detection.

Anomaly detection. CATCHCORE detected network intrusion in TCP dumps with high accuracy and identified attack patterns for DARPA and AirForce dataset, where each intrusion connection are labeled. Table 3 compares the accuracy of each method. CATCHCORE outperformed competitors for DARPA data, and also spotted the hierarchical behavior pattern of Neptune attack in H1 - H3, which are composed of the connections in different time for a pair of IPs. Fig. 2(c) shows the attack pattern snippet occurred during 7am - 8am on June 18, 1998. The densities (attack intensity) vary greatly over different hierarchies, i.e. the density in H1 is about 5K, while it is only about 3K for remain parts in H3. And the intense attacks represented cyclic patterns in 5 minutes. Although, the hierarchical structure of all subtensors include almost malicious connections (with recall = 98%) with the cost of containing some false positive samples, CATCHCORE achieves comparable performance for AirForce dataset.

Table 3. CATCHCORE identifies network attacks from the real-world TCP Dumps dataset with best or comparable performance in F-measure (Left); it spots different hierarchical dense blocks with interesting patterns for DARPA (Right).

	DARPA	AirForce	H	Subtensor Shape	Anomaly Ratio
CPD	0.167	0.785	1	$1 \times 1 \times 96$	100%
CROSSSPOT	0.822	0.785	2	$1 \times 1 \times 100$	100%
M-ZOOM	0.826	0.906	3	$1 \times 1 \times 274$	100%
D-CUBE	0.856	0.940	4	$16 \times 5 \times 24.7K$	87.0%
CATCHCORE	0.877	0.902	5	$171 \times 15 \times 29.2K$	85.4%

CATCHCORE also discerned denser and more anomaly subtensors. For the Youtube in Fig. 2(a), these dense blocks are missed by other competitors. Especially, the block with highest-density (H1) represents one user became friend with 904 users in a day, the other user in H2 also created connections with 475 users at the same time. So, they are more likely fraudulent. The densest block in StackOverflow shows one user marked 257 posts as the favorite in a day, which is too much than the normality. Although some of the densest blocks found may be rank-1 in all but one of the dimensions, CATCHCORE detects holistically optimal multi-layers dense subtensors and the densest one is only part of it rather than our direct target. The volume density tends to non-empty blocks and may result in some locally 1D slices (may not the densest slices within the whole tensor) in the highest-density layer. Other density metrics could eliminate this issue.

Evolving community pattern. As the Fig. 1(d), 1(b) show the evolving co-authorship structure of dense subtensors in the top 4 densest hierarchies for DBLP dataset, corresponding to the interaction between 20 researchers during 2007 to 2013, and Fig. 1(d) also presents their densities. The block in H1 with the size $8 \times 8 \times 3$, consists of research cooperation between **Leonard Barolli**, **Makoto Takizawa**, and **Fatos Xhafa**, etc. during 2011 to 2013 in ‘Algorithm and Distributed System’ field and the average connection between them is more than 10.7 in each year, forming a high-density clique. Also, the subtensors in other hierarchies are extended with their other common co-authors and years, and contain relatively less connections than H1, but the density of blocks in H4 is also more than 2. Therefore, CATCHCORE can cater to detect evolving community structures at different scales in a hierarchical fashion.

5.4 Q3. Scalability of CatchCore.

Empirically, we show that CATCHCORE scales (sub-) linearly with every aspect of input, i.e., the number of tuples, the number of dimension attributes, and the cardinality of dimension attributes we aim to find. To measure the scalability with each factor, we started with finding the injected subtensor with two hierarchies, which are $100 \times 100 \times 2$ with density 0.2 and $50 \times 50 \times 1$ with density 1.0, in a randomly generated tensor \mathcal{R} which contains 1 millions tuples with three

attributes whose total cardinality is 100K. Then, we measured the running time as we changed one factor at a time. As seen in Fig. 1(e) and result in supplement, CATCHCORE scales linearly with the number of tuples and the number of attributes, it also scales sub-linearly with the cardinalities of attributes, which illustrates the complexity of CATCHCORE in Theorem 7 is too pessimistic.

6 Related work

Dense Subgraph / Subtensor Detection. The detection of dense-subgraph has been extensively studied [2, 7], and many algorithms for the NP-hard problem are applied to detect community structure [3, 4, 26] and anomaly [1, 10], or extend to multi-aspect data [20–22]. CROSSSPOT [11] finds suspicious dense blocks by adjusting the seed in a greedy way until converges to a local optimum. Tensor decomposition such as HOSVD and CP decomposition [12] are also used to spot dense subtensors. M-ZOOM [20] and D-CUBE [21] adopt greedy approximation algorithm with quality guarantees to detect dense-subtensors for large-scale tensors. [22] spots dense subtensors for tensor stream with an incremental algorithm. None of them consider the relationship and structures of different blocks, and can not trace the evolving of dense subtensors or the hierarchical patterns.

Hierarchical Patterns Mining. Communities exit ubiquitously in various graphs [4, 14], their evolving behavior and hierarchical structure also have been explored in different scenes [13, 17, 25]. [8] proposed a framework for joint learning the overlapping structure and activity period of communities. [23] detected video event with hierarchical temporal association mining mechanism for multimedia applications. HIDDEN [26] detects hierarchical dense patterns on graph and also finds financial frauds. [19] uses k-core decomposition to compute the hierarchy of dense subgraphs given by peeling algorithm. Our method can deal with multi-aspect data, provide an information-theoretical measurement for the result, and advanced analyze for the performance.

Anomaly and Fraudulent Detection. The survey [1] provides a structured overview and summary for the methods of detection anomaly in graphs. The dense-subgraphs or dense-subtensors usually contain suspicious patterns, such as fraudsters in social network [10, 26], port-scanning activities in network analysis [11, 21], and lockstep behaviors or vandalism [11, 21, 22].

7 Conclusions

In this work, we propose CATCHCORE algorithm to detect the hierarchical dense subtensors with gradient optimization strategy, based on an novel and interpretable uniform framework for dense block detection, in large tensor. CATCHCORE accurately detects dense blocks and hierarchical dense subtensors for the synthetic and real data, and outperforms state-of-the-art baseline methods, it can identify anomaly behaviors and interesting patterns, like periodic attack and dynamic researcher group. In addition, CATCHCORE also scales up linearly in term of all aspects of tensor.

Acknowledgments

This material is based upon work supported by the Strategic Priority Research Program of CAS (XDA19020400), NSF of China (61772498, 61425016, 61872206), and the Beijing NSF (4172059).

References

1. Akoglu, L., Tong, H., Koutra, D.: Graph based anomaly detection and description: a survey. *Data mining and knowledge discovery* (2015)
2. Andersen, R., Chellapilla, K.: Finding dense subgraphs with size bounds. *WAW*
3. Balalau, O.D., Bonchi, F., Chan, T.H.H., Gullo, F., Sozio, M.: Finding subgraphs with maximum total density and limited overlap. In: *WSDM'15* (2015)
4. Chen, J., Saad, Y.: Dense subgraph extraction with application to community detection. *IEEE Trans. on knowledge and Data Engineering* (2010)
5. Coleman, T.F., Li, Y.: An interior trust region approach for nonlinear minimization subject to bounds. *SIAM Journal on optimization* (1996)
6. Edler, D., Bohlin, L., et al.: Mapping higher-order network flows in memory and multilayer networks with infomap. *Algorithms* **10**(4), 112 (2017)
7. Gibson, D., Kumar, R., Tomkins, A.: Discovering large dense subgraphs in massive graphs. In: *VLDB'05. VLDB Endowment* (2005)
8. Gorovits, A., Gujral, E., Papalexakis, E.E., Bogdanov, P.: Larc: Learning activity-regularized overlapping communities across time. In: *SIGKDD'18. ACM* (2018)
9. Grippo, L., Sciandrone, M.: On the convergence of the block nonlinear gauss-seidel method under convex constraints. *Oper. Res. Lett.* (2000)
10. Hooi, B., Song, H.A., Beutel, A., Shah, N., Shin, K., Faloutsos, C.: Fraudar: Bounding graph fraud in the face of camouflage. In: *SIGKDD'16*. pp. 895–904 (2016)
11. Jiang, M., Beutel, A., Cui, P., Hooi, B., Yang, S., Faloutsos, C.: A general suspiciousness metric for dense blocks in multimodal data. In: *ICDM'15* (2015)
12. Kolda, T.G., Bader, B.W.: *Tensor decompositions and applications*. SIAM (2009)
13. Kumar, R., Novak, J., Tomkins, A.: Structure and evolution of online social networks. In: *Link mining: models, algorithms, and applications* (2010)
14. Leskovec, J., Lang, K.J., Dasgupta, A., Mahoney, M.W.: Statistical properties of community structure in large social and information networks. In: *WWW* (2008)
15. Lin, C.J.: Projected gradient methods for nonnegative matrix factorization (2007)
16. Lin, C.J., Moré, J.J.: Newton's method for large bound-constrained optimization problems. *SIAM J. on Optimization* (1999)
17. Papadimitriou, S., Sun, J., Faloutsos, C., Philip, S.Y.: Hierarchical, parameter-free community discovery. In: *ECML-PKDD'08. Springer* (2008)
18. Rissanen, J.: A universal prior for integers and estimation by minimum description length. *The Annals of Statistics* (1983)
19. Sariyüce, A.E., Pinar, A.: Fast hierarchy construction for dense subgraphs. *VLDB*
20. Shin, K., Hooi, B., Faloutsos, C.: M-zoom: Fast dense-block detection in tensors with quality guarantees. In: *ECML-PKDD'16* (2016)
21. Shin, K., Hooi, B., Kim, J., Faloutsos, C.: D-cube: Dense-block detection in terabyte-scale tensors. In: *WSDM'17. ACM* (2017)
22. Shin, K., Hooi, B., Kim, J., Faloutsos, C.: Densealert: Incremental dense-subtensor detection in tensor streams (2017)

23. Siddique, B., Akhtar, N.: Temporal hierarchical event detection of timestamped data. ICCCA'17 (2017)
24. Tsourakakis, C., Bonchi, F., Gionis, A., Gullo, F., Tsiarli, M.: Denser than the densest subgraph: extracting optimal quasi-cliques with quality guarantees. In: SIGKDD'13. ACM
25. Yang, B., Di, J., Liu, J., Liu, D.: Hierarchical community detection with applications to real-world network analysis. DKE (2013)
26. Zhang, S., Zhou, D., Yildirim, M.Y., Alcorn, S., He, J., Davulcu, H., Tong, H.: Hidden: hierarchical dense subgraph detection with application to financial fraud detection. In: SDM'17. SIAM (2017)