

Scalable anomaly detection in graphs

William Eberle^{a,*} and Lawrence Holder^b

^a*Department of Computer Science, Tennessee Technological University, Cookeville, TN, USA*

^b*School of Electrical Engineering and Computer Science, Washington State University, Pullman, WA, USA*

Abstract. The advantage of graph-based anomaly detection is that the relationships between elements can be analyzed for structural oddities that could represent activities such as fraud, network intrusions, or suspicious associations in a social network. Traditionally, methods for discovering anomalies have ignored information about the relationships between people, e.g., who they know, or who they call. One approach to handling such data is to use a graph representation and detect normative patterns and anomalies in the graph. However, current approaches to detecting anomalies in graphs are computationally expensive and do not scale to large graphs. In this work, we describe methods for scalable graph-based anomaly detection via graph partitioning and windowing, and demonstrate its ability to efficiently detect anomalies in data represented as a graph.

Keywords: Anomaly detection, graph mining, dynamic graphs

1. Introduction

Institutions build data warehouses that store terabytes of information ranging from simple personal information to complex attributes that represent relationships between entities. Hellerstein calls it “the revolution of data” – so much data that even the term “big data” is being used to describe the phenomena [20]. Users then try to “mine” these databases for patterns to help them understand not only what the data means, but also discover those patterns that are unexpected or anomalous. Traditionally, methods for discovering patterns consist of “machine learning” approaches that use techniques such as classification, clustering, nearest neighbors, and statistics [17]. These methods have also focused mostly on the attributes of entities and ignored the *relationships* between entities.

Recent research efforts have involved the representation of this type of complex data as a *graph*, in order to analyze the relational *structure* in the data. This research has touched on a wide range of graph-theoretic approaches that have been applied to a wide variety of domains. While some successes have been demonstrated, they have either been specific to a particular data set, a particular type of graph, or a particular graph algorithm. In addition, they have not dealt with the scalability issues associated with “big data” when attempting to learn patterns and anomalies in data represented as a graph. For instance, in the case of computer network traffic, a graph representation of the traffic might consist of nodes representing computers and edges representing communications between the corresponding computers. In addition, other potential data sources for aiding in the analysis of the network traffic could include details about the individual users, location of the computer nodes, or even switch information. Adding

*Corresponding author: William Eberle, Department of Computer Science, Tennessee Technological University, Box 5101, Cookeville, TN 38505, USA. Tel.: +1 931 372 3278; Fax: +1 931 372 3686; E-mail: weberle@tntech.edu.

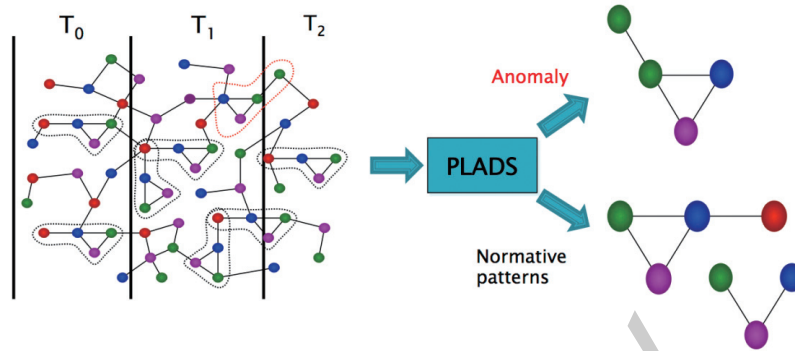


Fig. 1. Network monitoring scenario for PLADS. Information about entities and relationships streams in over time, and PLADS maintains a current set of normative patterns and anomalies. (Colours are visible in the online version of the article; <http://dx.doi.org/10.3233/IDA-140696>)

these heterogeneous data sets to the network traffic, represented as a graph, could provide the basis for discovering interesting structural patterns and anomalies, which may alert a security analyst to the potential threat in the form of a network intrusion attempt, denial-of-service attack, or worms. However, computer network traffic is typically voluminous, or acquired in real-time as a stream of information. For example, CAIDA (www.caida.org) provides a data repository to the research community for the analysis of internet traffic [5]. In one example of network traffic collected by CAIDA, representing a dynamic denial-of-service (DDOS) attack at a single location, every second produced an average of 3,992 transactions, for a total of 2,395,234 transactions over a 10 minute span.

In order to lay the foundation for this effort, we hypothesize that a real-world, meaningful definition of a *graph-based anomaly* is an *unexpected deviation to a normative pattern*. Assuming that an anomaly is not random, for instance in the case of committing fraud, we believe that this type of anomaly should only be a minor deviation from the normal pattern. Since anyone who is attempting to commit fraud or hide devious activities would not want to be caught, it only makes sense that they would want their activities to look as real as possible. For example, the United Nations Office on Drugs and Crime states the first fundamental law of money laundering as “The more successful money-laundering apparatus is in imitating the patterns and behavior of legitimate transactions, the less the likelihood of it being exposed.” [14]. In other words, such anomalies are associated with illicit activity that tries to mimic normal behavior. Thus, if some set of data is represented as a graph, any nefarious activities should be identifiable by small modifications, insertions or deletions to the normative patterns within that graph.

Our initial approach to graph-based anomaly detection, called GBAD [8], used a compression-based measure to find normative patterns, and then analyzed the close matches to the normative patterns to determine if they meet the above definition of an anomaly. However, while this approach has demonstrated its effectiveness in a variety of domains [9], the issue of *scalability* has limited this approach when dealing with domains containing millions of nodes and edges. Furthermore, many graphs of interest are dynamic, i.e., changes to the graph are streaming in over time. This further complicates the analysis, because we cannot just analyze a static graph, but would need to analyze snapshots of the graph over time. However, this streaming graph scenario also offers an opportunity for methods that can update the current set of patterns and anomalies based on only the changes to the graph, rather than repeated analyses on the large graph snapshots.

We have developed such a method for pattern learning and anomaly detection in streams (PLADS) depicted in Fig. 1. In this paper we describe the PLADS approach and demonstrate its effectiveness and

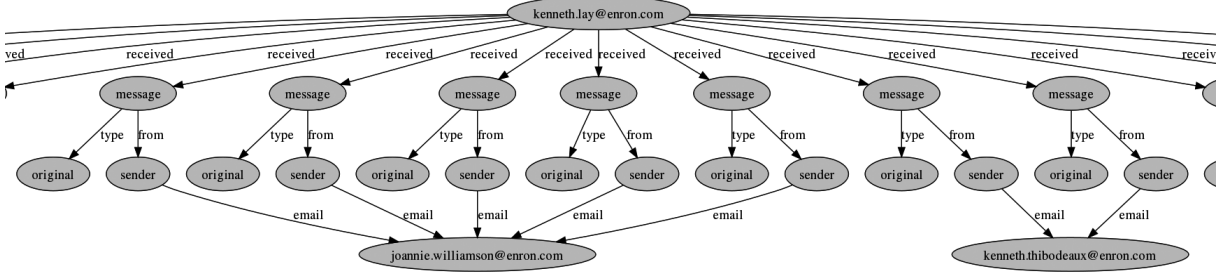


Fig. 2. Example partial graph of Enron e-mail correspondences.

scalability for large datasets. While the intent of PLADS is to operate on a streaming graph, at this point we are evaluating the approach on *partitions* of static graphs which we view as the output of a stream partitioner.

2. Graph-based anomaly detection

A graph is a set of nodes and a set of links, where each link connects either two nodes or a node to itself. More formally, we use the following definition.

Definition 1. A labeled graph $G = (V, E, L)$ consists of three sets V , E and L . V is the set of vertices (or nodes), E is the set of edges (or links) between the vertices, and L is the set of string labels assigned to each of the elements of V and E . [13]

Much work has been done using *graph*-based representations of data. Using *vertices* to represent entities such as people, places and things, and *edges* to represent the relationships between the entities, such as friend, lives-in and owns, allows for a much richer expression of data than is present in the standard textual or tabular representation of information. Representing various data sets like telecommunications call records, financial information and social networks in a graph form allow us to discover *structural* properties in data that are not evident using traditional data mining methods.

The idea behind our approach to graph-based anomaly detection is to find anomalies in graph-based data where the anomalous substructure in a graph is part of (or attached to or missing from) a *normative pattern*.

Definition 2. A substructure S_A is anomalous in graph G if $(0 < d(S_A, S) < T_D)$ and $(P(S_A|S) < T_P)$, where S is a normative pattern in G , T_D bounds the maximum distance an anomaly S_A can be from the normative pattern S , and T_P bounds the maximum probability of S_A .

Definition 3. The anomalous score of an anomalous substructure S_A based on the normative substructure S in graph G as $d(S_A, S) * P(S_A|S)$, where the smaller the score, the more anomalous the substructure.

The distance between two graphs can be due to the addition, removal or modification of structure from one graph to the other. The probability of S_A given S is based on the frequency of S_A among all graphs within distance T_D of S . Therefore, the more anomalous substructure is that which is closer to the normative pattern and appears with lower probability. The importance of this definition lies in its relationship to any deceptive practices that are intended to illegally obtain or hide information.

The advantage of *graph-based anomaly detection* is that the *relationships* between entities can be analyzed for structural oddities in what could be a rich set of information, as opposed to just the entities' attributes. For instance, take the example situation that occurred at the Enron Corporation [16]. Using anomaly detection on graphs that represent e-mail correspondences (Fig. 2), such as those between executives at Enron, might help in the prevention of lost revenues, pensions, and jobs. However, graph-based approaches have been prohibitive due to computational constraints. Because graph-based approaches typically perform subgraph isomorphisms, a known NP-complete problem, in order to address this issue, most approaches use some type of heuristic to arrive at an approximate solution. However, this is still problematic, and in order to use graph-based anomaly detection techniques in a real-world environment, we need to take advantage of the structural/relational aspects found in dynamic, streaming data sets.

3. GBAD

The PLADS approach is based on our previous work on static graph-based anomaly detection (GBAD) [8]. Here we briefly review the GBAD approach. There are three general *categories of anomalies* in a graph: insertions, modifications and deletions. Insertions would constitute the presence of an unexpected vertex or edge. Modifications would consist of an unexpected label on a vertex or edge. Deletions would constitute the unexpected absence of a vertex or edge. GBAD discovers each of these types of anomalies. Using a greedy beam search and a minimum description length (MDL) heuristic, GBAD first discovers the best substructure, or normative pattern, in an input graph. The MDL approach is used to determine the *best substructure(s)* as the one that minimizes the following:

$$M(S, G) = DL(G|S) + DL(S)$$

where G is the entire graph, S is the substructure, $DL(G|S)$ is the description length of G after compressing it using S , and $DL(S)$ is the description length of the substructure. In short, using the Minimum Description Length principle, the graph is represented as the number of bits needed to encode an adjacency matrix representation of the graph. Thus, the substructure that reduces the number of bits the most (by replacing all instances of the substructure with a single vertex) is considered the best substructure.

The GBAD approach is based on the exploitation of *structure* in data represented as a graph. We have found that a structural representation of such data can improve our ability to detect anomalies in the behaviors of entities being tracked [10]. GBAD discovers anomalous instances of structural patterns in data that represent entities, relationships and actions. GBAD uncovers the relational nature of the problem, rather than solely the traditional statistical deviation of individual data attributes. Attribute deviations are evaluated in the context of the relationships between structurally similar entities. In addition, most anomaly detection methods use a supervised approach, requiring labeled data in advance (e.g., illicit versus legitimate) in order to train their system. GBAD is an *unsupervised* approach, which does not require any baseline information about relevant or known anomalies. In summary, GBAD looks for those activities that appear to match normal/legitimate/expected transactions, but in fact are *structurally* different.

Once GBAD finds a normative pattern and anomalies in a graph, it can then iterate to find additional anomalies. First, GBAD compresses the graph using the normative pattern, i.e., replacing each instance of the normative pattern with a single node with a new label. Then, GBAD is executed on this compressed graph to again find normative patterns and anomalies. This process can continue for multiple iterations to find more and more normative patterns, and anomalies to them, throughout the graph, and at different levels of abstraction as the graph is further compressed.

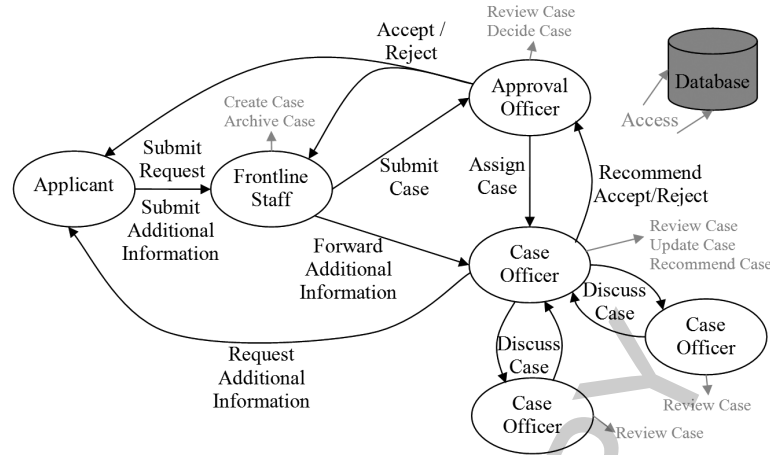


Fig. 3. Process flow of a passport application.

4. Empirical evaluations

First we show the performance of GBAD on single, large graphs. Then, to simulate a streaming setting, we evaluate GBAD on the same graphs where the approach is applied *independently* on roughly equal-sized *partitions*. In this case, a partitioning is composed of n subgraphs from a graph G , and the assumption is that these partitions represent portions of the graph streaming in over time. The value for n is chosen based on several considerations: size and diversity of the input graph, stream rate of the input graph, available processing resources. Ideally, we chose n so as to maintain partitions that are appropriately sized to allow GBAD to find the best normative patterns and anomalies, while keeping pace with the data stream.

We perform our experiments using a Debian 6.0, 64-bit server with 24 GB of RAM and 8 processors.

4.1. GBAD results

4.1.1. Business process

We simulate a passport processing scenario that is motivated by the CERT Insider Threat documents (http://www.cert.org/insider_threat/) that involve privacy violations in a government identification card processing organization and fraud in an insurance claim processing organization. The general process flow associated with a passport application is depicted in Fig. 3 [6].

As part of any business process, there is always the potential for someone to not adhere to the defined process. In some instances, perhaps the person did not understand the process, and inadvertently made a mistake. However, in some cases, people might be trying to illegally access information or manipulate data – what is typically referred to as an *insider threat*.

For instance, in the case of processing a passport (again, see Fig. 3), perhaps a Frontline Staff does not submit the case directly to the Approval Officer, or the Approval Officer does not review a case in the Database after it has been ruled upon by the Case Officer, or the Approval Officer does not assign a case to the Case Officer. Since these activities are required, in these instances the proper procedures for processing a passport would be violated.

Representing the processing of 1,000 passport applications, we generated a graph of approximately 5,000 vertices and 13,000 edges, and proceeded to replicate the above scenarios. (There are obviously

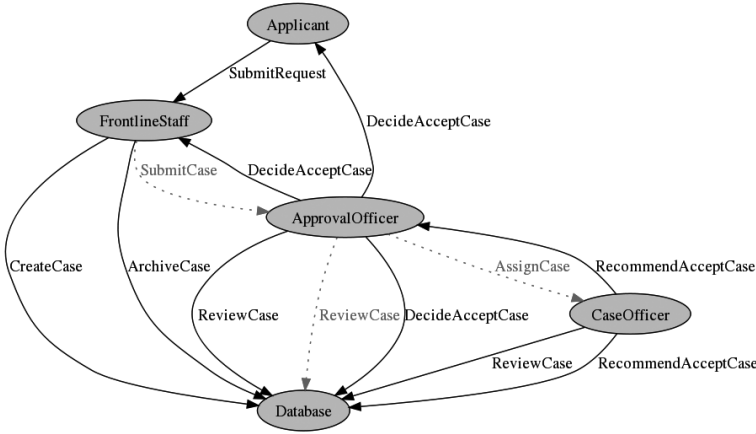


Fig. 4. Graph substructure with known anomalies.

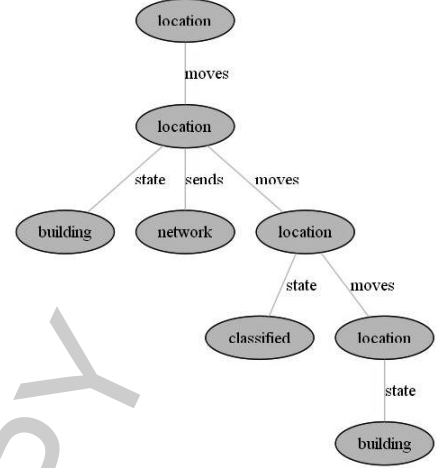


Fig. 5. Normative pattern in embassy data of movements and transactions.

many other potential insider threat scenarios, but for the sake of example, we will just show results on these three. This also allows us to control the anomalies present for verification of the accuracy of our approach.) We randomly inserted the three examples above into the graph. Running GBAD on this graph, all three instances of potential insider threat behavior are discovered. Figure 4 shows the normative pattern and the missing “SubmitCase” edge between “FrontlineStaff” and “ApprovalOfficer”, the missing second “ReviewCase” edge between “ApprovalOfficer” and “Database”, and the missing “AssignCase” edge between “ApprovalOfficer” and “CaseOfficer”. The actual anomalies in Fig. 4 are shown as dashed, red edges, indicating their absence from the graph.

The one drawback of this approach is that because of the size of the graph, it took GBAD 6,530 seconds to discover the anomalous structures. However, if we divide the graph into 10 roughly-equal-sized partitions (i.e., similar number of subgraph instances across each partition), and run GBAD on each of the partitions independently, the longest running partition only takes 42 seconds, and the anomalous substructures are discovered with no false positives.

It should be noted that the value 10 is chosen arbitrarily, and experiments with different values for n produce similar results. Obviously, the smaller the partition, the running times for individual partitions will decrease while comparisons *across* partitions will increase.

4.1.2. Cyber-threat

An example of a cyber-threat is the leaking of information by employees with access to confidential and sensitive information. One of the Visual Analytics Science and Technology (VAST) 2009 mini-challenges involved various aspects of a fictional insider threat scenario where someone is leaking information [11]. The goal of these challenges is to allow contestants to apply various visual analysis techniques to discover the spy and their associated actions. The VAST data set consists of the activities (card swipes and network traffic) of 60 employees at an embassy over the month of January in 2008.

On a graph of embassy employee activity data on January 29, 2008, consisting of 180 transactions (5,058 vertices and 4,878 edges), we randomly inserted an extra edge and vertex into the graph, representing a potentially anomalous *insertion*. We then run GBAD on the entire graph, targeting anomalous *insertions*. This results in the targeted anomaly being discovered in 2,364 seconds with no false positives. The normative pattern (shown in Fig. 5), consists of 8 vertices and 7 edges.

We then divide the graph into 10 graph partitions, where each partition consists of 18 transactions, and ran GBAD on each *partition* individually. This results in a much shorter running time, with the longest running partition taking only 215 seconds. However, while the targeted anomaly is discovered, it also results in 189 false positives being reported. (Again, choosing 10 partitions was arbitrary.)

The GBAD anomalous score (by our definition) is based upon the number of instances of a substructure that matches in every way except for a single extension. In these experiments, the partition that contains the targeted anomalous substructure contains 125 instances of a best substructure, whereas the partition that reports the most anomalous substructures has 138 instances of a best substructure. In other words, while all of the reported anomalies are single, unique instances (within the partition) of single edge/vertex insertions, the greater the number of normative instances, the lower the anomalous score (and greater the anomaly).

In these experiments, the normative pattern is the same across all partitions (except for partition 3, which contains the ninth ranked anomalous substructure). This leads us to make three observations: First, if we know the total number of instances of the best substructure, the targeted anomalous substructure would have a similar anomalous score as the reported most anomalous substructure. Again, based upon our definition of anomaly, because the anomalousness of a substructure is based upon its comparison to the best substructure, if we know the true best substructure we can determine the true anomaly. Second, and even more important, if we know the numbers of instances of each anomalous substructure across all of the partitions, our targeted anomalous substructure would come out on top by itself. This is due to the fact that there are fewer instances of the targeted anomalous substructure as compared to the reported anomalous substructures – hence a lower anomalous score. And third, if we keep track of the best/most anomalous score across the partitions, we would be able to remove some false positives. In other words, if we examine anomalousness globally (i.e., not just limited to the local partition), we gain a better picture of what substructures are truly anomalous.

4.2. Evaluating across partitions

One potential way to improve the accuracy of performing anomaly detection on a graph that has been partitioned (i.e., increase the number of true positives and decrease the number of false positives), is to examine the anomalies across all partitions. In other words, after all partitions have been processed, we can *evaluate all reported anomalies collectively* (i.e., across the partitions) in an effort to determine which one(s) is/are the true anomalous substructures.

While a simple partitioning of the passport application graph yields only the targeted anomalous substructures and no false positives, GBAD is not as effective on the graph of employee activity in the insider threat scenario. While there are 190 reported anomalous substructures, and only one targeted anomalous substructure, the anomalous scores vary across the different partitions. However, when analyzing the anomalous scores across the partitions, the targeted anomalous substructure ranks only fourth most anomalous, has 103 anomalous substructures ranked lower (i.e., more anomalous), and is equal in score to 63 other substructures. So, if we use just the anomalous score across all of the partitions, the targeted anomalous substructure would not be reported, and 14 false positives (i.e., substructures that are not anomalous across the entire graph) would be reported (i.e., the 14 substructures with the same lowest anomalous score, all from partition 7). However, if we factor in the *frequency* of the 190 reported anomalous substructures *across all of the partitions*, we discover that only two substructures are the most anomalous: the targeted anomalous substructure (from partition 1) and a substructure from partition 8. The GBAD algorithm factors in frequency along with the amount of change needed to make

two subgraphs match in order to get a final score of anomalousness. But, in this case, many of the substructures are not frequent within their individual partition. So, just evaluating the anomalies across the partitions will not eliminate the false positive reported in partition 8. While this particular substructure is prevalent in several other partitions (like 2, 3 and 5), and is thus not reported as anomalous in those instances, within its partition it only occurs once.

So, the question becomes, how do we deal with the false positive reported in partition 8? We know (by visually inspecting the data) that this particular substructure is prevalent in several other partitions where it is *not* reported as anomalous (so we do not have it among the original list of 190 potentially anomalous substructures). One way we can programmatically handle this is by allowing GBAD to make one more pass through the data, searching for the total number of instances of each extension. To demonstrate, if we provide each partition with the best normative pattern (from examining all of the partitions) as our “starting point”, and then total up all of the instances of each unique extension off of that normative pattern, the result is that our targeted anomalous substructure only occurs once (and is the only one that occurs once), and there are 15 instances of the false positive substructure. While this may seem to be time-consuming and not advantageous for scalability, specifying a normative pattern and counting the instances of the extensions took less than a second for each partition. So, this is a potentially valid approach for dealing with false positives that are too few to be normative and too numerous to be anomalous across the partitions.

If we can effectively detect anomalies across multiple graphs, we can attempt to handle not only very large graphs that are static (by partitioning them into multiple smaller graphs), but also graphs that represent a continuous *stream* of information.

5. Related work on mining large or streaming graphs

One potential solution to handling very large graphs is to view the graph as a “stream” and process the graph one, or a few edges, at a time. Previous work in this area has provided a few different approaches to the handling of *streaming graphs*.

One approach is to use what is called a *semi-streaming model* as a way of studying massive graphs whose edge sets cannot be stored in memory. For example, Feigenbaum et al.’s work presents semi-streaming constant approximation algorithms for un-weighted as well as weighted matching problems, as well as an improvement for handling bipartite graphs [12]. By considering a set of classical graph problems in their semi-streaming model, they were able to demonstrate that certain approximations to the problems can be achieved. Other work has generalized this approach to different graph problems, such as the shortest paths in directed graphs, and used intermediate temporary streams as a means of resolving the space issues [3,7,19]. Basically, these approaches propose a tradeoff between the available internal memory and the number of passes it requires.

Another approach is to examine the problem of *clustering* massive graph streams and use a technique for creating *hash-compressed micro-clusters* from graph streams [1]. Addressing the issues with large disk-resident graphs, the compressed micro-clusters are designed using a hash-based compression of the edges onto a smaller domain space. Again, the approach attempts to address the issue of graphs that are not available at one time on disk, but are continuously received in the form of a stream.

Recently, others have attempted to mine frequent closed subgraphs in non-stationary data streams. Once such approach called AdaGraphMiner, maintains only the current frequent closed graph, utilizing estimation techniques with theoretical guarantees [4]. Empirical experiments have demonstrated the

effectiveness of this approach on graph streams representing chemical molecules and structural representations of cancer data.

In addition, there have been recent attempts to discover outliers in massive network streams. Using what is called a structural connectivity model, some researchers have attempted to handle the issue of sparseness in massive networks by dynamically partitioning the network [2]. Using techniques such as reservoir sampling methods that compress a graph stream, one can search for structural summaries of the underlying network. The goal of this type of outlier detection is to identify graph objects which contain unusual bridging edges, or edges between regions of a graph that rarely occur together.

However, all of the approaches so far have not addressed the issue of *scalability* associated with performing *graph-based anomaly detection*. While some approaches have detected outliers in graph streams, their objective is to identify unusual clusters of subgraphs in the graph by analyzing the *statistical* nature of the existence of edges, as opposed to discovering anomalies in the *structure* of a graph, or graph stream. In addition, while some work has attempted to discover anomalous subgraphs using an ensemble-based approach [18] based on the GBAD approach [8], that type of approach does not address the issue of scalability.

6. A streaming-partitioned approach to graph-based anomaly detection

The advantages associated with graph-based anomaly detection are well-documented, providing a myriad of approaches for discovering structural and relational anomalies. However, they have been limited to static domains, or data sets that are relatively small in size – certainly nothing on the order of what we would call “big data”. What our preliminary experiments have shown us is that we *can* devise an approach whereby if we take into account smaller, individual partitions (i.e., a segment of the data that is processed individually, in parallel with other partitions) *in terms of what we know* about other partitions, we can not only provide similar accuracy but do it in a fraction of the time.

In order to formalize our approach, we propose the following PLADS algorithm. PLADS accepts as input a set of N graph partitions either by partitioning a static graph, or fed in over time.

PLADS (input graph partitions)

1. Process N partitions in parallel
 - a. Each partition discovers top M normative patterns.
 - b. Each partition waits for all partitions to discover their normative patterns.
2. Determine best normative pattern, P , among NM possibilities.
3. Each partition discovers anomalous substructures based upon P .
4. Evaluate anomalous substructures across partitions and report most anomalous substructure(s).
5. Process new partition
 - a. If oldest partition(s) has exceeded a threshold T (based upon criteria such as the number of available partitions or the time-stamped-age of the partition), remove partition(s) from further processing.
 - b. Determine top M normative patterns from new partition.
 - c. Determine best normative pattern, P' , among all active partitions.
 - d. If $(P' \neq P)$, each partition discovers new anomalous substructures based upon P' .
 - e. Else, only new partition discovers anomalous substructure(s).
 - f. Evaluate anomalous substructures across partitions and report most anomalous substructure(s).
 - g. Repeat.

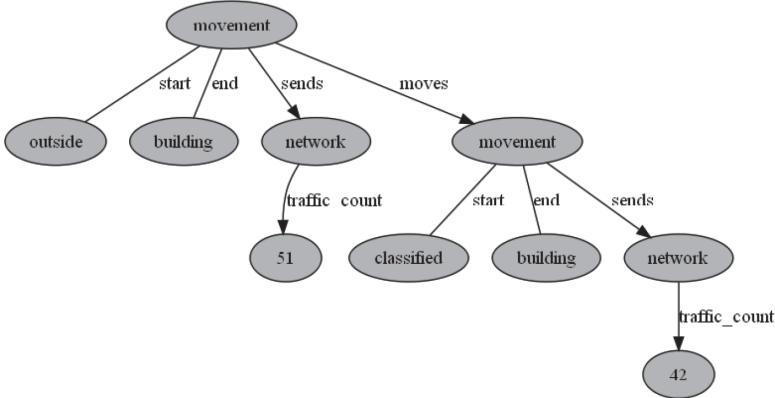


Fig. 6. Anomalous substructure in embassy data of movements and transactions.

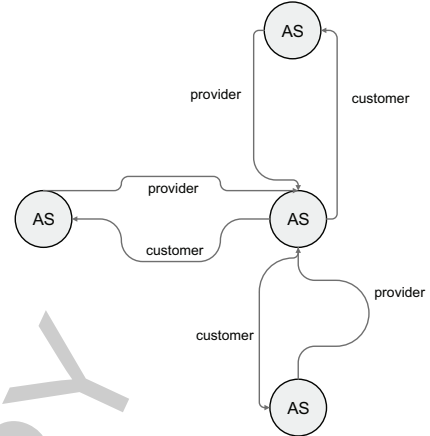


Fig. 7. Normative pattern discovered in the CAIDA dataset.

This is a very generic algorithm for applying graph-based anomaly detection methods to streaming data. The user can apply any normative pattern discovery techniques and any graph-based anomaly detection algorithms with this approach.

We now show how an approach like GBAD performs using such an algorithm.

6.1. Experiments using cyber-threat data

First, we will show the PLADS algorithm applied to a subset of the cyber-threat data presented earlier. In this experiment, we analyze just the movements of the employees throughout the embassy over the specified month of January in 2008. This set consists of card swipes as employees enter various rooms in the embassy.

6.1.1. GBAD

As input to GBAD, the data is represented as a graph, composed of 39,331 vertices and 38,052 edges, where movement, building, and type of room are depicted as vertices and edges indicating direction and movement between rooms. The normative pattern for this graph is depicted in Fig. 5. After running GBAD on the entire graph, two anomalous substructures are discovered (one of the substructures is shown in Fig. 6). However, it took 14,347 seconds to discover the anomalous substructure when analyzing the entire graph.

6.1.2. Streaming GBAD

In order to demonstrate the potential effectiveness of the PLADS approach, we applied the algorithm to the dataset that represents cyber-threat activity as follows.

1. *Process N partitions in parallel.*

We arbitrarily chose to initially process the first 5 (N) partitions of the graph. Running them in parallel, all of the partitions finish processing in 293 seconds, each producing 3 (M) normative patterns.

2. *Determine best normative pattern, P , among NM possibilities.*

We then examine all of the partitions' normative patterns, searching for the best normative substructure among them (i.e., the substructure that maximizes the value of size * frequency). The result is a normative pattern that is identical to the normative pattern shown in Fig. 5.

3. *Each partition discovers anomalous substructures based upon P .*

Based upon the best substructure from among all of the partitions (previous step), we then search for all anomalous substructures related to that normative pattern. The result is that only 1 substructure is reported as anomalous across all of the partitions, with the longest running partition taking 328 seconds.

4. *Evaluate anomalous substructures across partitions and report most anomalous substructure.*

For this example, since there is only one anomalous substructure reported, evaluation is trivial. (It should also be noted that this is one of the targeted anomalous substructures discovered when the graph was processed in its entirety.)

5. *Process new partition.*

Processing data as streams can be handled in two ways. Either we can always remove the oldest partition, or we can remove any partitions that are older than some time threshold T (i.e., a sliding window). For this example, we will do the former, removing the oldest partition and processing a new partition (e.g., removing partition 1 and processing partition 6). We then discover the best substructure on the new partition, so that we can determine the best normative pattern among all of the remaining partitions. However, while the reported normative pattern in partition 6 is different, it is not better than the best substructure reported by the other four partitions. So, we use the best substructure on partition 6, and no anomalous substructures are discovered (in 106 seconds). Also, since we are using the best substructure from a previous iteration, we do not have to re-discover any anomalous substructures in the older partitions.

At the next iteration (e.g., partition 2 is removed and partition 7 is added), we discover that the normative pattern has not changed (i.e., it is still the best substructure across all of the active partitions). Again, only the new partition needs to be analyzed for any anomalous substructures, as the anomalies would not change for the already processed partitions. Analysis of the results from the new partition (partition 7) yields (in a total of 257 seconds) no anomalous substructures.

This same behavior continues over partitions 8 and 9, using 207 and 301 seconds respectively. However, on partition 10, the same best substructure is reported, but a new anomalous substructure is reported of equal “anomalousness” (in 501 seconds) to the substructure discovered in partition 3. This happens to be the second anomalous substructure discovered when the entire, non-partitioned graph was processed.

So, we are able to implement a graph-based anomaly detection approach on data that represents movements of people, and successfully discover the same two anomalous substructures (with no false positives) within a streaming approach in a fraction of the time (1,993 seconds) it took to process the entire graph (14,347 seconds). However, this graph is rather sparse (i.e., few edges compared to the number of vertices). So, now we will examine results on a denser graph that also represents data that can be streamed.

6.2. Experiments using actual network data

The Cooperative Association for Internet Data Analysis (CAIDA) is a publicly available resource for the analysis of IP traffic. Through a variety of workshops, publications, tools, and projects, CAIDA provides a forum for the dissemination of information regarding the interconnections on the internet. One of the core missions of CAIDA is to provide a data repository to the research community that will allow for the analysis of internet traffic and its performance (<http://www.caida.org/data/>). Using GBAD, we analyzed the CAIDA AS (Autonomous Systems) data set for normative patterns and possible anomalies. The AS data set represents the topology of the internet as the composition of various Autonomous Systems. Each of the AS units represents routing points through the internet.

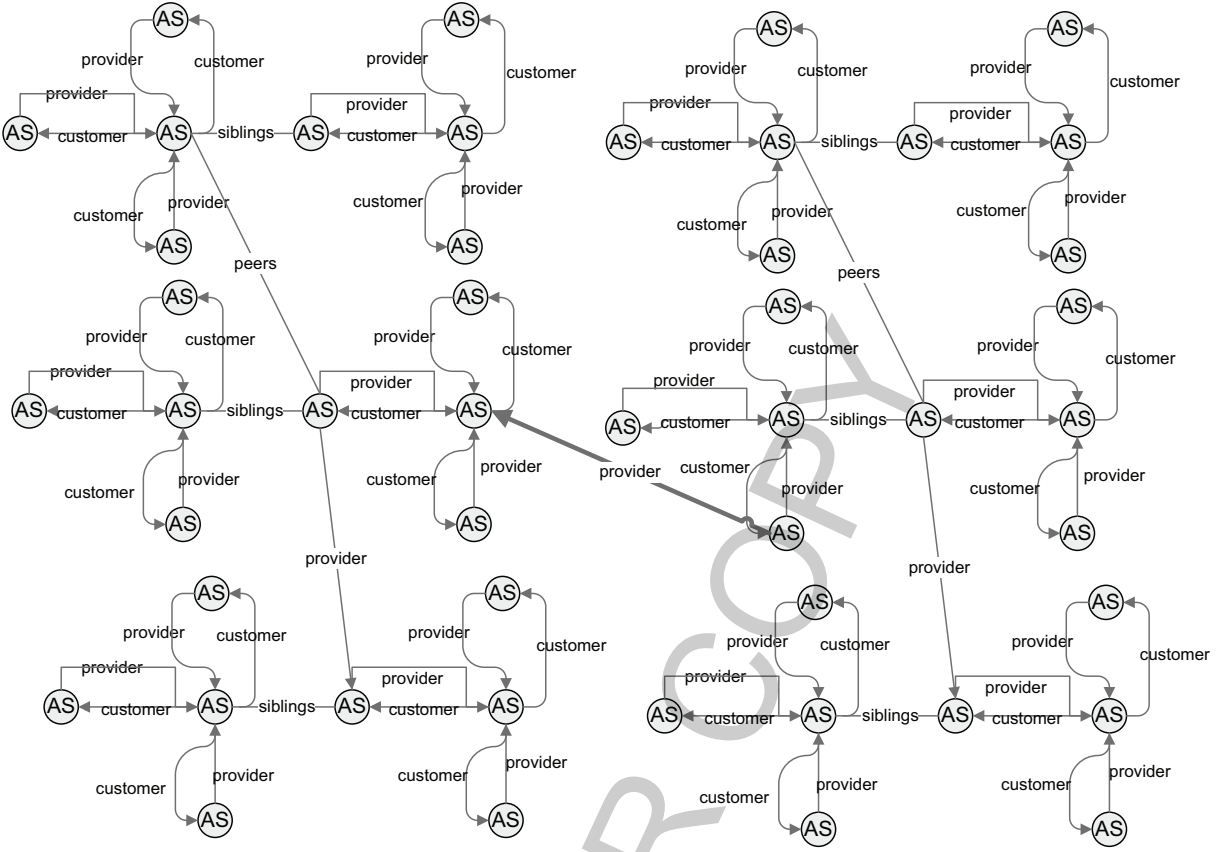


Fig. 8. Anomalous pattern discovered in the CAIDA dataset, where anomalous component is the central “provider” edge.

6.2.1. GBAD

For the purposes of analysis, the data is represented as a graph, composed of 24,013 vertices and 98,664 edges, with each AS depicted as a vertex and an edge indicating a peering relationship between the AS nodes. The normative pattern for this graph is depicted in Fig. 7. After running GBAD on the AS graph, the anomalous substructure discovered is shown in Fig. 8.

This example shows the advantage of using a graph-based approach on a complex structure. While the data indicates many provider/customer relationships, of which the norm is a particular AS being the provider to three different customers, this single substructure indicates an unusual connection between two ASes. Such an inconspicuous structure would probably be missed by a human analyst, and shows the potential of an approach like GBAD to find these anomalies in network traffic data.

However, just like the examples shown earlier, it took 59,743 seconds to discover the anomalous substructure using a tool like GBAD.

6.2.2. Streaming GBAD

In order to demonstrate the potential effectiveness of a *streaming* approach to graph-based anomaly detection, we apply the PLADS algorithm to this same CAIDA data set as follows.

1. *Process N partitions in parallel.*

We again arbitrarily chose to initially process the first 5 (N) partitions of the graph. Running

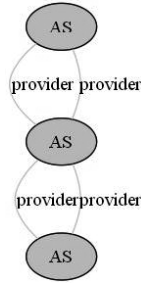


Fig. 9. Normative pattern early in the “stream”.

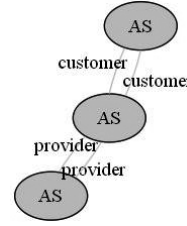


Fig. 10. Normative pattern later in the “stream”.

them in parallel, all of the partitions finish processing in 210 seconds, each with 3 (M) normative patterns.

2. *Determine best normative pattern, P , among NM possibilities.*

We then examine all of the partitions' normative patterns, searching for the best normative substructure among them. The result is the normative pattern shown in Fig. 9, which is smaller than the normative pattern found when running on the entire graph (see Fig. 7).

3. *Each partition discovers anomalous substructures based upon P .*

Based upon the best substructure from among all of the partitions (previous step), we then search for all anomalous substructures related to that normative pattern. The result is that 166 substructures are reported as anomalous across all of the partitions, with the longest running partition taking 112 seconds.

4. *Evaluate anomalous substructures across partitions and report most anomalous substructure.*

We then examine all of the reported anomalous substructures across the partitions, and the result is that 2 substructures are reported as equally anomalous. However, at this point, neither of the substructures are the targeted anomalous substructures. (The timing for this step takes less than a second.)

5. *Process new partition.*

Similar to the previous example, we handle the data as a stream by removing the oldest partition and processing a new partition. We then discover the best substructure on the new partition, so that we can determine the best normative pattern among all of the remaining partitions. The result is the discovery of the normative pattern shown in Fig. 7 (i.e., the same normative pattern from processing the entire graph) in 92 seconds.

Since the normative pattern has changed since the last iteration, we have each partition *re-discover* any anomalous substructures *based upon* the new normative pattern. Examining all of the reported anomalous substructures across the partitions, we discover that the most anomalous substructure (found in partition 5) is the one that was identified when we ran GBAD on the entire graph!

At the next iteration (e.g., partition 2 is removed and partition 7 is added), we discover in 45 seconds that the normative pattern has not changed (i.e., it is still the best substructure across all of the active partitions). In this case, only the new partition needs to be analyzed for any anomalous substructures, as the anomalies would not change for the already processed partitions. Analysis of the results from the new partition (partition 7) yields (in 58 seconds) no substructures more anomalous than what were already discovered.

Taking this scenario one more iteration (e.g., partition 3 is removed and partition 8 is added), we discover that the best normative pattern across all of the partitions is different from the previous

iteration (see Fig. 10). So, similar to two iterations back, all of the active partitions need to be re-evaluated based upon this new best substructure. The result is two new anomalous substructures. However, if you compare their “anomalousness” to the one reported earlier (shown in Fig. 8), one can see that there are more instances of this newly reported anomalous substructure, so the anomalous substructure discovered earlier would still be the most anomalous.

After two more iterations of adding and removing partitions (i.e., processing all of the partitions that represented the single graph), the new normative pattern stays the same, and the anomalousness of reported substructures lessens (i.e., becomes more common), still leaving us with the targeted anomalous instance.

So, we are able to implement a graph-based anomaly detection approach on network data that is able to successfully discover the same anomalous substructure within a streaming approach in a fraction of the time (642 seconds) it took to process the entire graph (59,743 seconds). Even the overhead associated with comparing normative patterns and anomalous substructures across partitions is negligible, as the number of substructures to evaluate from each partition is minimal.

6.3. Experiments using larger synthetic graphs

While the previous real-world data set experiments analyze some interesting scenarios, the data sets are relatively small and the ability to control the anomalies is limited. So, in order to validate our approach on larger graphs and vary the substructures, we used a synthetic graph generator to generate a sparse graph of ~ 2 M vertices and edges. While not on the order of what most would define as “big data” [15], where graphs consist of billions of nodes, processing partitions of this size using the PLADS approach could quickly add up to this scale.

6.3.1. GBAD

The graph consists of a specified normative subgraph of 10 vertices and 9 edges, with random substructures of varying levels of anomalousness (i.e., frequency of their existence) injected into the graph. The normative pattern for this graph is depicted in Fig. 11. After running GBAD on the complete graph, the anomalous substructure, consisting of an unexpected edge and vertex, is discovered as shown in Fig. 12. However, the normative pattern and anomalous substructure are discovered in 276,873 seconds (i.e., over 3 days) – hardly useful in a real-world environment.

6.3.2. Streaming GBAD

Again, to demonstrate the potential effectiveness of a *streaming* approach to graph-based anomaly detection, we apply the PLADS algorithm to this same synthetic graph. For this example, we have divided the original graph into 100 partitions, where each partition consists of approximately 19,000 vertices and edges.

1. *Process N partitions in parallel.*

We initially process the first 20 (N) partitions of the graph. The choice of an initial 20 partitions is somewhat arbitrary, as our goal is to just get a representative sample with which to start analyzing. Running them in parallel, all of the partitions finish processing in 9 seconds, each with 3 (M) normative patterns. It should be noted that even if the partitions are not processed in parallel, it only takes 136 seconds to process the 20 partitions serially. However, we will also use $T = 20$ as the size of our processing window (i.e., how many partitions will be retained in memory).

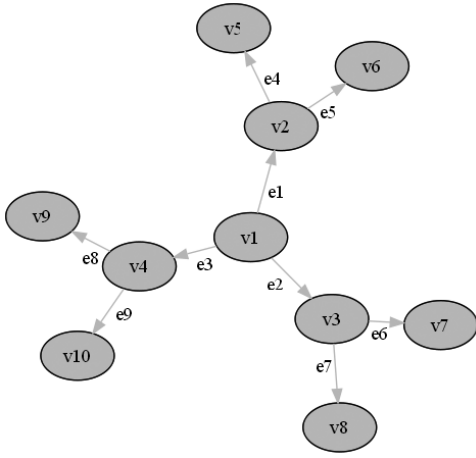


Fig. 11. Normative pattern in synthetic graph.

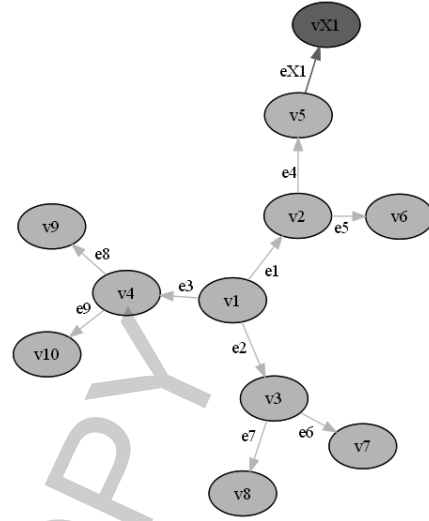


Fig. 12. Anomalous substructure in synthetic graph.

2. *Determine best normative pattern P among NM possibilities.*

We then examine all of the partitions' normative patterns, searching for the best normative substructure among them. The result is the normative pattern shown in Fig. 11, which is the same (targeted) normative pattern found when running on the entire graph.

3. *Each partition discovers anomalous substructures based upon P .*

Based upon the best substructure from among all of the partitions (previous step), we then search for all anomalous substructures related to that normative pattern. The result is that only 2 anomalous instances (from partition 17) are reported as anomalous across all of the partitions.

4. *Evaluate anomalous substructures across partitions and report most anomalous substructure.*

We then examine all of the reported anomalous substructures across the partitions (in this case, they are both in the same partition), and the result is that the substructure is the same, so it is reported as the lone anomalous substructure. At this point, the anomalous substructure is not the targeted anomalous substructure. (Again, the timing for this step takes less than a second.)

5. *Process new partition.*

As before, this is a repetitive step where we process one partition at a time, handling the data as a stream by removing the oldest partition and processing a new partition. We then discover the best substructure on the new partition, so that we can determine the best normative pattern among all of the remaining partitions. If the normative pattern changes, we then evaluate all partitions against the new normative pattern, searching for new anomalous substructures. If the best normative pattern does not change, we only need to analyze the new partition for potential anomalous substructures. Processing partitions 21–33 results in no new normative patterns or anomalies, at a total processing time of 92 seconds. Processing partition 34 does not report a new normative pattern, but 2 instances of an anomalous substructure are discovered. The new anomalous substructure is evaluated against the current best anomalous substructure, and it is discovered to be the same anomalous substructure (albeit, still not the targeted anomaly at this point). Thus, we now have 4 instances of the current best anomalous substructure. This step takes 8 seconds.

Processing partitions 35–57 results in no new normative patterns or anomalies, at a total processing time of 165 seconds. Processing partition 58 does not report a new normative pattern, but 3

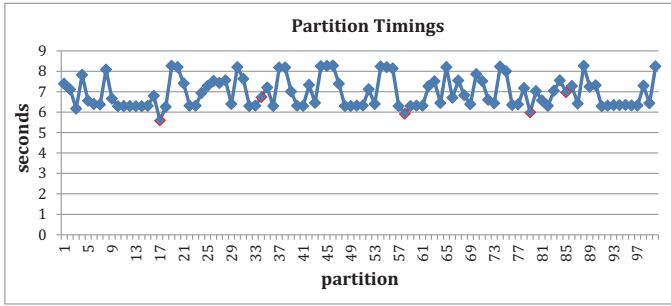


Fig. 13. Partition timings. (Colours are visible in the online version of the article; <http://dx.doi.org/10.3233/IDA-140696>)

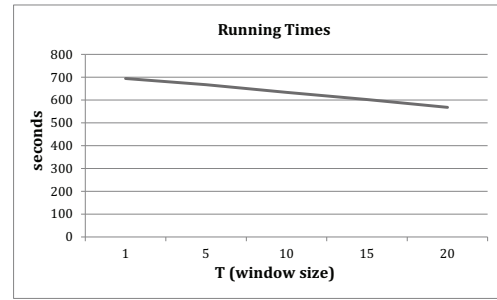


Fig. 14. Running times on synthetic graph based upon initial window size (T).

instances of a new anomalous substructure are discovered. The anomalous substructure is evaluated against the current best anomalous substructure, and is found to be different and more anomalous. Thus, this new substructure becomes the current best anomalous substructure. This step takes 7 seconds, and again, it is still not the targeted anomaly.

Processing partitions 59–78 results in no new normative patterns or anomalies, at a total processing time of 141 seconds. Processing partition 79 does not report a new normative pattern, but 3 instances of an anomalous substructure are discovered. The anomalous substructure is evaluated against the current best anomalous substructure, and it is discovered to be the same anomalous substructure. Thus, we now have 6 instances of the current best anomalous substructure. This step takes 7 seconds. It is also interesting to note that in terms of “global anomalousness”, our visibility is limited to the partitions that are retained in memory. If we could compare the current best anomalous substructure to the one discovered in partition 17 (as well as partition 34), the older substructures would be more anomalous.

Processing partitions 80–84 results in no new normative patterns or anomalies, at a total processing time of 35 seconds. Processing partition 85 does not report a new normative pattern, but 1 instance of an anomalous substructure is discovered. The anomalous substructure is evaluated against the current best anomalous substructure, and is found to be different and more anomalous. Thus, this new substructure, which is the targeted anomaly that is discovered when analyzing the entire graph (shown in Fig. 12), becomes the current best anomalous substructure. This step takes 7 seconds.

Processing partitions 86–100 results in no new normative patterns or anomalies, at a total processing time of 103 seconds.

Again, we are able to implement a graph-based anomaly detection approach on a larger graph that is able to successfully discover the same targeted anomalous substructure within a streaming approach in a fraction of the time (574 seconds) it took to process the entire graph (276,873 seconds). Timings for each partition are shown in Fig. 13, with partitions containing anomalous substructures shown in red.

As noted earlier, the size of the window (i.e., the number of partitions retained in memory) does affect what anomalous substructures are discovered. At the end of processing all 100 partitions, the targeted anomalous substructure is reported as the best, but, for instance, if the targeted anomalous substructure was in an early partition, a later partition may report a less anomalous substructure (as opposed to when we process the non-partitioned graph), as the targeted anomalous substructure will have fallen out of the window when $T = 20$. However, this fits into the concept of handling data as a stream, whereby subgraphs that are reported as anomalous may lessen in their “anomalousness” as new information is received.

Running this same experiment using different values for T (the initial number of partitions, N , is set to the same value), we discover that:

- $T = 1$: the newest partition always has the current best anomalous substructure (if any).
- $T = 5$: due to the scarcity of injected anomalous substructures, only one partition in the window contains an anomalous substructure, leaving the newest partition to always have the current best anomalous substructure.
- $T = 10$: when the targeted anomalous substructure is discovered at partition 85, within the window is the previous best anomalous substructure, which would be replaced by the targeted anomalous substructure because it is more anomalous (score wise). (Same results for $T = 15$.)

In addition to accuracy, the running times are also slightly affected by the change in T , and are shown in Fig. 14. Here we see that total running time decreases as the window size T increases due to the increased ability to exploit parallelism across partitions in the window. Thus, T (and N) is bound by the number of processors available. While there is some overhead associated with having to compare substructures (normative and anomalous) across partitions, the overhead is minimal both from a run-time perspective as well as system resources. It should also be noted that the value chosen for M (number of normative patterns), while defaulted to 3, can be increased with minimal impact to performance. For a user perspective, one must determine how truly “normative” is a pattern that is not in the top M .

7. Conclusions and future work

Handling large or streaming graphs provides the opportunity to handle complex data sets that are well-suited for graph-based approaches. We have proposed a method for analyzing graphs using a parallel approach that can discover anomalous substructures – particularly ones that indicate potential criminal activity. We have also demonstrated the scalability of our approach with an order-of-magnitude improvement in the running-times of a graph-based anomaly detection approach. Using real-world data from a cyber-threat scenario, and actual network traffic between autonomous systems, we have been able to discover the anomalies with minimal false-positives using a parallel approach to processing segments of the entire graph. In addition, we have demonstrated the ability to discover targeted anomalous substructures in synthetic graphs that have been scaled-up to represent potentially larger domains.

While there are several different approaches to graph-based knowledge discovery and anomaly detection, the algorithm presented is not dependent on a single approach. In future work, we hope to examine the scalability of such an approach to “big data” sizes as well as high-speed streams. In addition, we would like to further examine this approach in a variety of other real-world domains, such as social networks, process controls, biological networks, and credit-card fraud.

Acknowledgements

This material is based upon work supported by the National Science Foundation under Grant Nos. 1318913 and 1318957. The authors would also like to thank Mr. Abhik Ray, a PhD student at Washington State University, for his aide in partitioning the graphs. Support for CAIDA’s Internet Traces is provided by the National Science Foundation, the US Department of Homeland Security, and CAIDA Members.

References

- [1] C. Aggarwal, Y. Zhao and P. Yu, On clustering graph streams, *SIAM* (2010).

- [2] C. Aggarwal, Y. Zhao and P. Yu, Outlier detection in graph streams, *ICDE* (2011).
- [3] G. Aggarwal, M. Datar, S. Rajagopalan and M. Ruhl, On the streaming model augmented with a sorting primitive, in: *Proceedings of the 45th Annual IEEE Symposium on Foundations of Computer Science (FOCS)* (2004).
- [4] A. Bifet, G. Holmes, B. Pfahringer and R. Gavaldà, Mining frequent closed graphs on evolving data streams, in: *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge discovery and Data Mining (KDD'11)* (2011).
- [5] The CAIDA as relationships dataset, <http://www.caida.org/data/active/as-relationships>.
- [6] A. Chun, An AI framework for the automatic assessment of e-government forms, *AI Magazine, Spring* **29** (2008).
- [7] C. Demetrescu, I. Finocchi and A. Ribichini, Trading off space for passes in graph streaming problems, in: *Proceedings of the 17th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)* (2006).
- [8] W. Eberle and L. Holder, Anomaly detection in data represented as graphs, *Intelligent Data Analysis* **11**(6) (2007).
- [9] W. Eberle, L. Holder and J. Graves, Insider threat detection using a graph-based approach, *Journal of Applied Security Research* **6**(1) (January 2011), 32–81.
- [10] W. Eberle, L. Holder and D. Cook, Identifying threats using graph-based anomaly detection, in: *Machine Learning in Cyber-Trust*, Springer, J. Tsai and P. Yu, eds, (May 2009).
- [11] W. Eberle, L. Holder and J. Graves, Detecting employee leaks using badge and network IP traffic, *IEEE Symposium on Visual Analytics Science and Technology* (October 2009).
- [12] J. Feigenbaum, S. Kannan, A. McGregor, S. Suri and J. Zhang, On graph problems in a semi-streaming model, *Theoretical Computer Science* (2005).
- [13] J. Gross and J. Yellen, Graph theory and its applications, CRC Press, 1999.
- [14] M. Hampton and M. Levi, Fast spinning into oblivion? Recent developments in money-laundering policies and offshore finance centres, *Third World Quarterly* **20**(3) (1999), 645–656.
- [15] U. Kang and C. Faloutsos, Big graph mining: Algorithms and discoveries, *ACM SIGKDD Explorations* **14**(2) (2012), 29–36.
- [16] B. McLean and P. Elkind, The smartest guys in the room: The amazing rise and scandalous fall of enron, Portfolio Trade Publishing, 2004.
- [17] T. Mitchell, Machine learning, McGraw Hill, 1997.
- [18] P. Parveen, J. Evans, B. Thuraisingham, K. Hamlen and L. Khan, Insider threat detection using stream mining and graph mining, *Privacy, Security, Risk and Trust (PASSAT), 2011 IEEE Third International Conference on and 2011 IEEE Third International Conference on Social Computing (SocialCom)* (2011), 1102–1110.
- [19] A. Sarma, S. Gollapudi and R. Panigrahy, Estimating PageRank on graph streams, *AMC PODS* (2008).
- [20] D. Weinberger, Data, data everywhere, *The Economist* (February 27, 2010).