

Graph-Based Anomaly Detection on Smart Grid Data

Lenin Mookiah, Chris Dean, William Eberle

lmookiah42@students.tntech.edu, chrisdean258@gmail.com and weberle@tntech.edu
Tennessee Tech University
Cookeville, TN 38505.

Abstract

With the rise in usage of smart meters, the issue of security has been an active area of research. One particular area of research is the potential tapping into and stealing of power that is connected to a single smart grid. In this research, we present current security issues related to homes whose power consumption is regulated, monitored, and ultimately billed to the consumers. While most approaches deal with detecting potential irregularities in electrical consumption based upon traditional statistical methods, few efforts have attempted to model a city or neighbourhood smart grid as a *network* in order to detect unusual patterns in a *structural* representation. In this work, we propose the use of a graph-based anomaly detection approach for detecting anomalies in power usage. We empirically evaluate a graph-based anomaly detection approach on actual smart home data, and demonstrate its potential effectiveness.

Introduction

Rapid technological advancement and integration into household appliances has produced a highly interconnected electrical network. However, with this network also comes the need for approaches to deal with new security issues. In general, any proposed security approach needs to be reliable, accurate, efficient, and privacy preserving. In addition, and the crux of this work, is that it also needs to detect not only malicious activity but also natural anomalies in the system that could be damaging to appliances or allowing security breaches. Groups of buildings and their associated appliances, connected to an electric grid, are called Smart Cities. As with any central repository of resources, they typically produce (and process) enormous amounts of data that would have to be carefully handled so as to avoid potential data leaks of sensitive information (Kitchin 2016) (Han and Xiao 2016). In addition to the smart grid infrastructure being an efficient power delivery mechanism, it is also susceptible to potential fraudulent usage of its resources.

A variety of approaches have been researched and analyzed for improving security and anomaly detection in smart cities. Machine learning techniques such as artificial neural

networks (ANNs) and rules based algorithms have been implemented with some success (Anwar and Mahmood 2014). Other approaches such as context awareness and graphical approaches such as petri nets and dependency graphs have also been implemented with varying levels of success, albeit primarily focused on securing the network (Peterson 1981) (Alcaraz, Cazorla, and Fernandez 2014) (Chen, Sanchez-Aarnoutse, and Buford 2011).

In this work, we examine some of the existing approaches, and then present a graph-based knowledge discovery approach for detecting anomalies. We will use this approach for analyzing actual data that contains the detailed electrical usage of households on a smart grid. These dwellings consist of many devices such as refrigerators, dishwashers, air conditioners, and even pool pumps, just to name a few. The objective of this work is to demonstrate that a graph-based approach is able to detect anomalies in power consumption on a network of homes connected to a smart grid. We propose that a graph-based knowledge discovery technique, where vertices (nodes) will represent smart appliances and edges (links) will represent usage between the parts of a home and their appliances, will allow us to uncover anomalies by taking advantage of the structure in such a representation. Such an approach could be a great complement to existing artificial neural networks and statistical approaches, providing a more comprehensive solution to detecting unusual behaviour in a smart grid network.

Related Work

Recently, several attempts have been made to define the problems a smart grid may face, along with suggesting requirements that must be met for the solution to be viable. Murillo suggested that any anomaly detection scheme must be reliable, efficient, and secure against various types of attacks that may be executed against a smart grid system (Murillo 2014). To this point, a two step mutual authentication protocol is suggested for securing the data transfer. Liu et al. suggest standards for smart grid communication based on the need for availability above integrity and confidentiality (Liu et al. 2012). They also address the need for an accurate and timely anomaly detection system, suggesting an artificial neural network and a multi-agent based fault location algorithm as a possible solution. Anwar et al. also outline the requirements for smart grid anomaly detection

(Anwar and Mahmood 2014). They discuss several types of attacks and how they may be detected using various hardware and software prevention methods.

Others have suggested schemes for anomaly detection. Alcaraz et al. suggest a system based on context awareness (Alcaraz, Cazorla, and Fernandez 2014). Their system implements different (separate) types of anomaly detection methods such as statistical, data mining, knowledge-based, and machine learning methods, based on the context of the targeted anomaly. They go on to outline the application and application area for where each anomaly detection method could be used. Kher et al. suggest a prediction model for anomaly detection in smart grids (Kher et al. 2013). This model is based on machine learning algorithms using clustering as a model. While the prediction algorithms were not able to detect all intrusions, the results were promising. With the introduction of a more refined and accurate prediction algorithm, this method could become a very effective method of intrusion detection.

Ten et al. propose an anomaly detection algorithm based on detection at the substation level (Ten, Hong, and Liu 2011). Although substation is only one of the levels at which a smart grid would have to detect anomalies, there is potential for this algorithm to be expandable to some of the larger structures within the smart grid system. A rules based intrusion detection scheme based on behavior was proposed by Mitchell et al. (Mitchell and Chen 2013). The model detects malicious activity with a low false-positive rate. In this work, the authors study attacks on three types of devices: securing head-ends (HEs), distribution access points/data aggregation points (DAPs) and subscriber energy meters (SEMs). Their proposed algorithm is able to detect the intrusion attack on all the three types of devices.

Researchers have also studied various graph based approaches. For instance, a dependency graph approach is presented by He et al. (He and Zhang 2011). The authors propose a Markov random field (MRF) in terms of physical parameter of networks. Then, using localization in the conditional correlation matrix of MRF, implements decentralized network inference of fault diagnosis. The authors mainly study two type of outages due to: (1) power line outage, and (2) change of physical parameters. The author shows that the proposed approach effectively catches the fault lines.

Calderaro et al. use another type of graph called a petri net to rapidly diagnose faults in the smart grid (Calderaro et al. 2011). The graph represents complex, large, distributed networks that involve protective relays and circuit breakers. The authors propose two case studies, each having different protection systems that require coordination. The case studies show that assessment of information in the network is much easier and effective in terms of not requiring complex data analysis.

In summary, there are two different research directions for smart grid in the literature. First, researchers are studying the smart grid in order to secure the data transfer using encryption. Second, researchers are studying the problem of anomaly detection in a smart grid. There are many devices attached to a smart grid, and such devices could either come under attack from hackers or run errant. Our research work

dataid	local_15min	furnace1	grid	lights_plugs2
5218	2014-04-03 00:30:00	0.0172	0.6242	0.08553
5218	2014-04-03 00:45:00	0.0180	0.7016	0.08600
5218	2014-04-03 01:00:00	0.0176	0.7092	0.08613

Table 1: Sample Data

falls under this second direction. Specifically, we are interested in discovering anomalies in a context graph, representing devices in the smart grid.

Data

We collected real-world data from *Pecanstreet*¹ during the period of January 2014 to December 2014. *Pecanstreet* works on two biggest problem in the world - water and energy. Specifically, Pecanstreet focuses on university research and accelerating innovation in water and energy. The total dataset consists of smart meter readings of power utilization collected from 67 electronic devices in 820 households. The devices include, but are not limited to, water heaters, air conditioners, light plugs, dish washers, jacuzzis, ice makers, clothes washers, furnaces, ovens, pools, refrigerators, sprinklers, microwaves, house fans, pool pumps, and winecoolers. The dataset is also further divided by city, of which 532 households are from Austin, Texas. Of those 532 households, 155 are of the housing type *Apartment*, and 377 are of the housing type *Single Family Home*. In this work, we decided to focus on only the 155 apartments, as we hypothesize that the demographic and usage would be similar. (We are planning on applying this same approach to the single family homes in our future work.) A data sample is shown in Table . In the table, *dataid* represents a unique id for each household, *local_15min* represents the datetime of the reading, *furnace1* represents power used by furnace 1, *grid* represents the power used in the grid, and *lights_plug2* represents the power used by light plug 2.

Graph Topology

In order to apply a graph-based approach for detecting *structural* anomalies in the smart grid data, we first must convert the data into a graph. Since there is not a standard way to represent this type of data as a graph, we designed our own topology, where nodes represent the household, grid, generator, rooms, and devices; and edges represent the existence of the corresponding room and the usage of a device. Each graph substructure consists of a node labelled as “Home” that is connected to two primary household nodes: one for the total solar power generated (“Gen”), and another for the total power supplied by the smart grid (“Grid”). In Figure. 1, the edge “with-usage-from-a” is to represent the devices of each house in the graph topology.

Per the schema defined by the provider of the smart grid data, household devices are categorized into one of six

¹<https://dataport.pecanstreet.org/>

Each targeted injection usage results in the creation of an anomalous edge and vertex for the device in the corresponding input graph file - a *structural anomaly*. The 24 injection anomalies are shown in Table . The anomalies are injected so as to reflect the corresponding malicious attack, and result in a structural change to the underlying graph structure. For example, the first anomaly in the table is from dataid 7512 with datetime 2014-01-10 04:45:00, where we inserted the usage of a “jacuzzi1” into a household that does not have one. In other words, from a graph topology perspective, an edge “with-usage-from-a” between node “poolpump” and its category node “outdoorSpace” (see Table) is created. Similarly, the other anomalies in the table are inserted into our daily graphs.

Regarding contextual structure, there are other types of scenarios which could create anomalies. For example, a malicious insider/employee of a utility company could try to inflate usage bills by corrupting the data with malicious code for some targeted households (Jiang et al. 2014) (Anwar and Mahmood 2014).

Graph Based Approach

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. The importance of this definition (which we more formally define below) lies in its relationship to any deceptive practices that are intended to illegally obtain or hide information (Hampton and Levi 1999).

Definition 1. A labeled graph $G = (V, E, F)$, where V is the set of vertices (or nodes), E is the set of edges (or links) between the vertices, and the function F assigns a label to each of the elements in V and E .

Definition 2. A subgraph SA is anomalous in graph G if $(0 < d(SA, S) < TD)$ and $(P(SA|S) < TP)$, where $P(SA|S)$ is the probability of an anomalous subgraph SA given the normative pattern S in G . TD bounds the maximum distance (d) an anomaly SA can be from the normative pattern S , and TP bounds the maximum probability of SA .

(It should be noted that for our implementation, SA will have a maximum value of 1, and value of d is optionally given by user, with a default value of 4.)

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

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. In order to test our approach, we will implement the publicly-available GBAD test suite , as defined by (Eberle and Holder 2007). Using a greedy beam search and a minimum description length

(MDL) heuristic, GBAD first discovers the “best” subgraph, or normative pattern, in an input graph. The MDL approach is used to determine the best subgraph(s) as the one that minimizes the following.

$$M(S, G) = DL(G|S) + DL(S) \quad (1)$$

where G is the entire graph, S is the subgraph, $DL(G|S)$ is the description length of G after compressing it using S , and $DL(S)$ is the description length of the subgraph. The complexity of finding the normative subgraph is constrained to be polynomial by employing a bounded search when comparing two graphs. Previous results have shown that a quadratic bound is sufficient to accurately compare graphs in a variety of domains (Kukluk, Holder, and Cook 2004).

For more details regarding the GBAD algorithms, the reader can refer to (Eberle and Holder 2007). In summary, the key to the GBAD approach is that anomalies are discovered based upon small *structural* deviations from the norm (e.g., insider threat, identity theft, etc.) – not outliers, which are based upon *statistical* deviations from the norm.

We will briefly discuss the input graph format for GBAD. Below is an example of how a single day usage of devices are represented in our graph. The example is from dataid 22 on date "2014-09-30 00:00:00". The first line indicates that there is a vertex with index id 1 and a label "furnace1". The first edge line, i.e., d 2 5 "with-usage-from-a", represents a directed edge from vertex 2 to 5, with an edge label "with-usage-from-a". An example of a portion of the graph-input file is shown here:

```
v 1 "furnace1"
v 2 "grid"
v 3 "kitchen1"
v 4 "livingroom1"
v 5 "home"
v 6 "smallAppliance"
v 7 "bigAppliance"
v 8 "UtilityRoom"
v 9 "BedRoom"
v 10 "LivingRoom"
v 11 "OutdoorSpace"
d 2 5 "with-usage-from-a"
d 5 6 "that-has-a"
d 5 7 "that-has-a"
d 5 10 "that-has-a"
d 5 11 "that-has-a"
d 5 9 "that-has-a"
d 5 8 "that-has-a"
d 9 1 "with-usage-from-a"
d 6 3 "with-usage-from-a"
d 10 4 "with-usage-from-a"
```

Experiments

With graph input files having the graph topology as shown in Figure 1, GBAD discover the normative (best) substructure and any anomalous substructures related to the corresponding home. We create daily-graph files for each day of the

year 2014, for a total of 365 days. We then run GBAD on each of the files. All of the experiments are run on an Intel(R) Xeon(R) CPU E5520, 2.27GHz, RAM 24 GB, 8 core machine.

Evaluation

We propose to use four different metrics - precision, recall, f1-score, and accuracy, which are represented in the Equation 2, Equation 3, Equation 4, and Equation 5 respectively, in order to evaluate this approach. The metrics are calculated using true positive (TP), true negative (TN), false positive (FP), and false negative (FN).

$$precision = \frac{TP}{TP + FP} \quad (2)$$

$$recall = \frac{TP}{TP + FN} \quad (3)$$

$$f1 - score = \frac{2 * precision * recall}{precision + recall} \quad (4)$$

$$accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (5)$$

Using GBAD, any anomalous node(s) and/or edge(s) are reported for each daily graph. It is also possible that an anomalous node (i.e., device) could be identified as an anomaly one or more times in the output of GBAD. If a node/edge is identified as an anomaly more than once (i.e., duplicate), it is removed from consideration as an anomaly. In other words, such nodes/edges are considered as true negatives. For example, on 2014-04-23, “dishwasher1” is marked as an anomaly 10 times. Hence, it would no longer be considered an anomaly (i.e., it would be considered as a true negative). In this dataset, there are 6,162 nodes marked as duplicates or true negatives. In contrast, on day 2014-06-12, node “dryg1” is marked as an anomaly as it occurs only once (i.e., it is considered a true positive).

Our approach has some limitations. GBAD works effectively in detecting single packet injections for a device, as shown in this paper. However, in the case of an attacker repeatedly injecting false packets targeting the same device, GBAD would need to be configured from its default state to recognize that the multiple instances of these particular substructures are not normative patterns.

Results and Analysis

As was discussed earlier, GBAD produces normative (best) substructure and any anomalous substructure related to the normative pattern. Figure 2 shows the normative substructure of dataids 7512 and 7585, and Figure 3 shows the anomalous substructure discovered marked in dashed lines for 7512. Similarly, Figure 4 shows the anomalous substructure discovered, visually represented with dashed lines, for dataid 7585. The GBAD approach is able to successfully discover all the anomalies from the Table, i.e., true positives. Table 5 shows the false positives reported by the GBAD approach (i.e., non-targeted anomalies). On this smart grid data, the GBAD approach gives a precision of 92.30, recall of 100.00, f1-score of 95.99, and accuracy of 99.41.

dataid	device	datetime
68	dryg1	2014-06-12 00:00:00
5252	disposal1	2014-09-10 00:00:00

Table 5: False positives using GBAD approach.

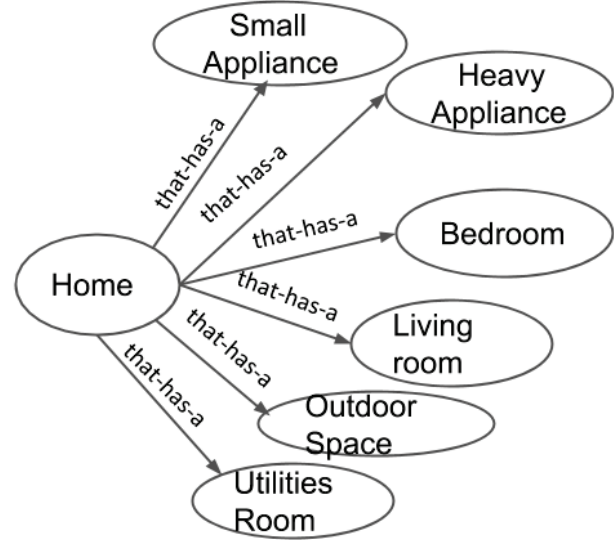


Figure 2: Normative Pattern of dataid 7512 and 7585 for the day 2014-01-10 and 2014-04-11 respectively.

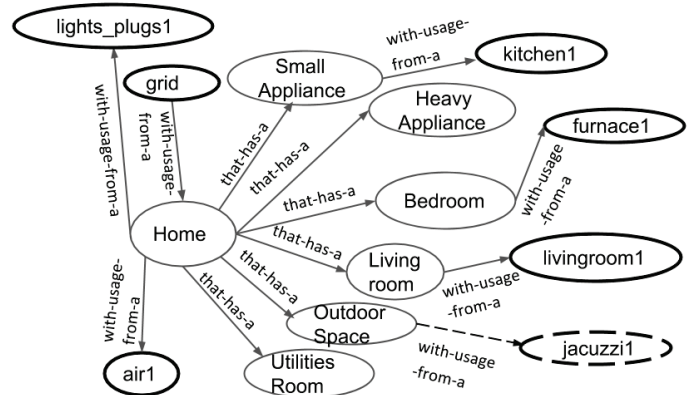


Figure 3: Anomalous substructure of dataid 7512 for the day 2014-01-10. Anomalous edge and node are marked in dashed lines.

Conclusions and Future Work

In this work, we study the problem of anomaly detection in smart grid data using a graph-based approach. We represent power usage of different devices from homes as context graphs. We show that an approach like GBAD effectively discovers anomalies with high precision, recall, and accuracy. The strength of our GBAD approach is that it can han-

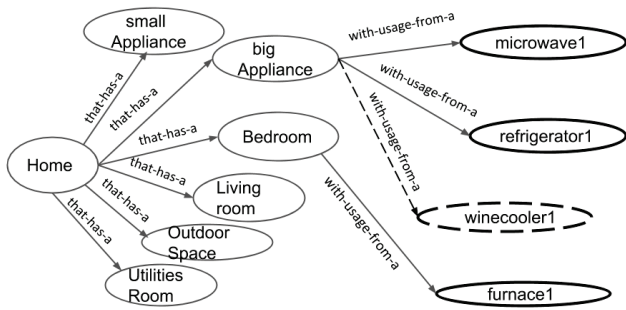


Figure 4: Anomalous substructure of dataid 7585 for the day 2014-04-11. Anomalous edge and node are marked in dashed lines.

dle contextual information, whereas methods such as neural networks are not designed to leverage *structure*. In addition, GBAD is an unsupervised approach compared to common classification algorithms such as neural networks and graphical models (Sutton and McCallum 2006). One of the weaknesses of our approach is that it has limitations on handling cases where attackers launch attacks on multiple packets of the same device in the smart grid.

In addition to analyzing other segments of the smart grid data, such as the single-family homes or households in other cities, there are several broader future directions for graph-based anomaly detection in the smart grid. First, we will study smart grid data as a *stream*, i.e., a continuous feed of information. This will not only allow us to potentially discover other types of anomalies, but also allow us to keep up in real-time. Second, we will study the application of graph-based anomaly detection on a complex power supply grid with sophisticated devices such as power relays, etc. Third, we will study a hybrid approach whereby we incorporate data visualization techniques with graph-based approaches.

Acknowledgement

We sincerely thank Vitaly Ford from the Cybersecurity Education, Research and Outreach Center (CEROC) for his expertise in smart grid data and security. This material is based upon work supported by the National Science Foundation under IIS Grant No. 1318657 and CNS Grant No. 1560434.

References

Alcaraz, C.; Cazorla, L.; and Fernandez, G. 2014. Context-awareness using anomaly-based detectors for smart grid domains. In *International Conference on Risks and Security of Internet and Systems*, 17–34. Springer.

Anwar, A., and Mahmood, A. N. 2014. Cyber security of smart grid infrastructure. *arXiv preprint arXiv:1401.3936*.

Calderaro, V.; Hadjicostis, C. N.; Piccolo, A.; and Siano, P. 2011. Failure identification in smart grids based on petri net modeling. *IEEE Transactions on Industrial Electronics* 58(10):4613–4623.

Chen, T. M.; Sanchez-Aarnoutse, J. C.; and Buford, J. 2011.

Petri net modeling of cyber-physical attacks on smart grid. *IEEE Transactions on Smart Grid* 2(4):741–749.

Eberle, W., and Holder, L. 2007. Anomaly detection in data represented as graphs. *Intelligent Data Analysis* 11(6):663–689.

Hampton, M. P., and Levi, M. 1999. Fast spinning into oblivion? recent developments in money-laundering policies and offshore finance centres. *Third World Quarterly* 20(3):645–656.

Han, W., and Xiao, Y. 2016. Non-technical loss fraud in advanced metering infrastructure in smart grid. In *International Conference on Cloud Computing and Security*, 163–172. Springer.

He, M., and Zhang, J. 2011. A dependency graph approach for fault detection and localization towards secure smart grid. *IEEE Transactions on Smart Grid* 2(2):342–351.

Jiang, R.; Lu, R.; Wang, Y.; Luo, J.; Shen, C.; and Shen, X. S. 2014. Energy-theft detection issues for advanced metering infrastructure in smart grid. *Tsinghua Science and Technology* 19(2):105–120.

Kayastha, N.; Niyato, D.; Hossain, E.; and Han, Z. 2014. Smart grid sensor data collection, communication, and networking: a tutorial. *Wireless communications and mobile computing* 14(11):1055–1087.

Kher, S.; Nutt, V.; Dasgupta, D.; Ali, H.; and Mixon, P. 2013. A prediction model for anomalies in smart grid with sensor network. In *Proceedings of the Eighth Annual Cyber Security and Information Intelligence Research Workshop*, 60. ACM.

Kitchin, R. 2016. Getting smarter about smart cities: Improving data privacy and data security.

Kukluk, J. P.; Holder, L. B.; and Cook, D. J. 2004. Algorithm and experiments in testing planar graphs for isomorphism. *Journal of Graph Algorithms and Applications* 8(2):313–356.

Liu, J.; Xiao, Y.; Li, S.; Liang, W.; and Chen, C. P. 2012. Cyber security and privacy issues in smart grids. *IEEE Communications Surveys & Tutorials* 14(4):981–997.

Mitchell, R., and Chen, R. 2013. Behavior-rule based intrusion detection systems for safety critical smart grid applications. *IEEE Transactions on Smart Grid* 4(3):1254–1263.

Murillo, A. F. 2014. Review of anomalies detection schemes in smart grids. Springer.

Peterson, J. L. 1981. Petri net theory and the modeling of systems.

Sutton, C., and McCallum, A. 2006. An introduction to conditional random fields for relational learning. *Introduction to statistical relational learning* 93–128.

Ten, C.-W.; Hong, J.; and Liu, C.-C. 2011. Anomaly detection for cybersecurity of the substations. *IEEE Transactions on Smart Grid* 2(4):865–873.

Xie, L.; Mo, Y.; and Sinopoli, B. 2010. False data injection attacks in electricity markets. In *Smart Grid Communications (SmartGridComm), 2010 First IEEE International Conference on*, 226–231. IEEE.