

Improving the Generalizability of BGP Anomaly Detection Models

JIAOJIAO CHENG (862187837)*, COLIN LEE (862185718)*, ZHUOCHENG SHANG (862188698)*, and WILLIAM SHIAO (861280411)*,

Group #2, CS 240 (Fall 2020), University of California Riverside, USA

Final presentation slides: [Google Slides](#)

1 EFFORT

Jiaojiao used Pyod[16] to test several classic anomaly detection algorithms. The test based on three dataset: code red, Nimda and slammer. Models are trained on two datasets and tested on the other one in turns. The methods include KNN, isolation forest, PCA, feature bagging and ABOD. Also, OCSVM and basic Auto Encoder in pyod have been test to provide result comparisons for our reimplemented edition. She compared these methods by their F1 scores and accuracy, and gave a general feature importance evaluation and analysis of those algorithms.

Colin rebuilt the Deep Support Vector Data Description (Deep SVDD) method [11]. The original architecture of the network was built to work with images, so this involved reworking the architecture of the network as well and investigating the effects of the hypersphere parameters to achieve convergence. He additionally worked on the abstract and problem statement of the report and wrote the interpretation of the feature ablation results.

Zhuocheng is responsible for implementing Entropy-One Class SVM, and training three different datasets through basic OCSVM implemented with Scikit-learn API. For general OCSVM¹, she trained model with all normal data, and tested on the mixed data (both normal and outliers). The next challenge solved was calculating normalized entropy [10] for each distinct features (total is 41) partitioned based on sequential time interval (24 hours period in our dataset). Following that, the normalized entropy should be applied on training the Entropy-OCSVM model, the modified model was expected to train and predict directly on test data.

William worked on implementing the Fence-GAN (FGAN) [9] and autoencoder (AE) models. He reimplemented FGAN in PyTorch based off of the reference TensorFlow implementation², and adapted it from a convolutional architecture for images to a feed-forward model for our dataset (although it did not work well, as described below). He modified the PyTorch Lightning AE bolt³ to work on our dataset instead of images. He also wrote the code to load the dataset into the correct combinations to evaluate the explainability of our models.

2 ABSTRACT

The practicality of supervised BGP anomaly detection is limited to detecting events that have been previously observed and trained on. Much of the previous work has focused primarily on improving classification scores on datasets drawn from the same events as the training data. This begs the

* All authors equally contributed to this project.

¹https://scikit-learn.org/stable/auto_examples/svm/

²https://github.com/phuccuongngo99/Fence_GAN/

³<https://pytorch-lightning-bolts.readthedocs.io/en/latest/autoencoders.html#basic-ae>

question: *Can current BGP anomaly detection methods generalize well across different anomalies? And if so, which features of the data are most important to this generalizability?*

Our project includes a multitude of different machine learning anomaly detection methods applied, rebuilt, or re-implemented in Python for BGP anomaly detection. We use these methods in an ablation experiment to determine the networking features key to generalizability and show the performance of these methods without those features. In all, this totals about 2900 lines of code. The application of the more complex machine learning anomaly detection methods has not yet seen application to the specific area of BGP anomaly detection, and the simple ones have not seen an analysis of their generalizability. Our project offers the first analysis that compares the generalizability of all of these methods and determines which networking features are vital to good generalization.

3 PROBLEM

Given a model M ,
Train M on a BGP anomaly dataset D
Detect BGP anomalies in other dataset(s)

As with any machine learning literature, the current literature for BGP anomaly detection does test for generalization, but only between training and testing sets *drawn from the same anomaly*. While these sorts of results are indicative of good performance for the tested anomalies, they cannot tell us much about the ability of models to generalize across multiple anomalies.

In this project, we present a novel analysis of BGP anomaly detection methods *between* different anomalies. We compare the generalization performance of several anomaly detection approaches, some of which have never been applied to this domain before. We measure the generalizability of these models using accuracy and F1 scores on test sets consisting of multiple anomalies. Furthermore, we explore the features responsible for generalization performance and attempt to interpret their importance.

3.1 Scope

An "anomaly" in networking can be difficult to define, as it is difficult to precisely characterize the appearance of "normal" traffic. A 2017 survey of BGP anomaly detection approaches by Al-Musawi et al. [2] constructs a taxonomy based on the cause of the anomalous behavior. E.g. "Direct" anomalies are caused by problems with the network itself, such as prefix hijacks ("direct intended") or origin misconfigurations("direct unintended"), while indirect anomalies occur as a result of events such as a worm spreading across the Web.

Finding data for many different anomalies was difficult itself, but compounding that problem was the fact that datasets generally share very few of the same features, which makes comparisons between them difficult. As a result, we limit our dataset to three indirect anomalies caused by computer worms in the early 2000s: Code Red I, Slammer, and Nimda. It would be ideal to include different types of anomalies from different time periods, but we cannot both gather the necessary data and perform the analysis with the time we have. The relative similarity of these anomalies can be beneficial though, as any deficiencies in generalization will represent a sort of upper bound on generalization performance.

3.2 Previous and Related Work

BGP anomaly detection technique is a hot topic aiming at detecting and alerting anomalous events so as to minimize the damage it causes. According to the survey[1], typical types of anomalies

can be concluded as point anomaly, contextual anomaly and collective anomaly. Worms, power outages, and BGP router configuration errors are all considered as anomalous events. These attacks are sharp, resulting in sustained increases in the number of announcement or withdrawal messages exchanged by BGP routers[6].

Machine learning based(classification, clustering, statistical analysis) network anomaly detection is an important category in all the detection methods, thus making both datasets and feature selection algorithms significant. Well-known datasets like Route Views, Réseaux IP Européens (RIPE), and BCNET[6] are widely chosen by researchers to apply various feature selection algorithms. One outstanding approach is the automatic feature extraction by neural network[15] for it solving the problem of manually selecting and deciding the statistical features used for anomaly detection. Generally, the output of anomaly detection techniques are scores and/or binary labels, and evaluation is generated based on both criteria.

Neural networks have been commonly applied to anomaly detection. For supervised anomaly detection, many recent papers use LSTMs (Long Short-Term Memory) [8] and GRUs (Gated Recurrent Unit) [5], which work on a timeseries, rather than an individual sample. Generative Adversarial Networks (GANs) [7] were original developed for image generation, but have since been applied to a wider variety of tasks, including anomaly detection. Fence-GAN [9]

For one-class learning, most models use the idea of fitting a hyperplane or hypersphere to the data and then using that to classify anything outside of it as anomalous. One common model for this purpose is the Support Vector Data Description (SVDD) [14] model, which is inspired by the Support Vector Machine (SVM) [13]. Another relatively common family of models for this task are the one-class neural networks (OC-NN) [4]. One example of this is the Deep SVDD [12], which has been successfully applied to image classification.

However, most of these one-class techniques have yet to be applied specifically to BGP anomaly detection. The closest work we were able to find was by Allahdadi *et al.*[3], which uses a one-class SVM. Thus far, machine learning for BGP anomaly detection has not grown very complex, and still uses relatively simple methods such as SVMs [13], or simple RNNs such as LSTMs (Long Short-Term Memory) [8] or GRUs (Gated Recurrent Unit) [5].

Other anomaly detection tasks have seen the use of more complex models involving deep learning such as, GANs (Generative Adversarial Networks) [7] [9] and Deep SVDD (Support Vector Data Description) [12]. These have yet to be applied to the specific domain of BGP anomaly detection.

As stated above, there currently is no literature which compares the inter-anomaly generalizability of these models, so we will have to design our experimental framework from scratch.

4 SOLUTION

Since previous work in BGP anomaly detection has not thus far examined the generalization performance between different anomalies, we begin by surveying the generalizability of different anomaly detection methods on our three separate worm datasets. Using the information gleaned from this survey, we select the method with the highest generalizability (which we define as the highest mean F1 score), and perform a feature ablation experiment with it in order to ascertain which networking features are most important for generalizing between different the three different worms.

4.1 Generalization Performance

In order to evaluate the inter-anomaly generalization performance of each method, we use three combinations of training and testing data where each method is trained on one of the anomalies and tested on the other two.

Model	Nimda		Code Red		Slammer		Mean	
	Acc.	F1	Acc.	F1	Acc.	F1	Acc.	F1
One-class Methods								
OC-SVM	0.102	0.185	0.253	0.404	0.237	0.384	0.198	0.324
Entropy OC-SVM	0.102	0.185	0.253	0.404	0.237	0.384	0.198	0.324
Autoencoder	0.903	0.501	0.608	0.282	0.745	0.530	0.752	0.438
Deep SVDD	0.921	0.479	0.804	0.354	0.726	0.388	0.817	0.407
Unsupervised Methods								
KNN	0.895	0.488	0.786	0.453	0.778	0.449	0.820	0.463
Isolation Forest	0.856	0.395	0.769	0.355	0.769	0.324	0.798	0.358
PCA-based	0.827	0.295	0.711	0.245	0.747	0.177	0.762	0.239
LOF w/ Feature Bagging	0.779	0.249	0.679	0.178	0.695	0.176	0.717	0.201
Angle-based Outlier Detector	0.892	0.488	0.768	0.372	0.772	0.382	0.811	0.414

Table 1. All of the scores shown are the scores when the model is trained on the listed dataset and evaluated on the remaining two datasets. LOF stands for Local Outlier Factor.

We note that autoencoders appear to be the best one-class method and KNN is the best unsupervised method for generalizing between these datasets. In general, the unsupervised methods outperform the one-class methods, which is an expected result, given that the unsupervised methods learn from both anomalous and unanomalous data, whereas the one-class methods are limited to learning their decision boundaries from the unanomalous data.

4.2 Feature Ablation Analysis

In order to explore *what* makes these models generalizable, we ablatively remove features from our data and note which features are most important for the best methods. We then note the most and least important features based on how much they decrease the F1 score of the methods.

Method	Most important	Least important
KNN	Number of withdrawn NRLI prefixes Number of announcements	Packet size Number of duplicate withdrawals
ABOD	Number of withdrawn NRLI prefixes Avg unique AS-path	Packet size Number of duplicate withdrawals
AE	Avg AS-path length Max AS-path length	Max AS-path length = 8 Number of duplicate withdrawals

Table 2. Top two most important and least important features for the three best methods: K-nearest neighbor, Angle-based Outlier Detector, Autoencoder

In general, it seems that the number of withdrawn NLRI prefixes and one of the AS-path length features are important to the generalization performance of the each of the models. This suggests that generalizable models must key in on features that are informative about the reachability of other ASs. Worms such as the ones investigated in these datasets typically affect networks most when attempt to rapidly replicate themselves, overloading the capacity of the network and leading

to Denial of Service events. As such, it would follow that some nodes become unreachable and different AS-paths must be found.

Least important features included packet size, and the number of duplicate withdrawals. Packet size is negligible as packet size can vary normally for any number of reasons. That the number of duplicate withdrawals is unimportant to all three suggests that normal traffic sees a similar number of duplicate withdrawals as anomalous traffic.

5 CHALLENGES

Since nobody else has performed this type of analysis on BGP data before, we also had difficulty deciding how exactly to evaluate our methods. We had a hard time deciding on how to evaluate the generalizability of our methods.

Ideally, the anomalies in our dataset would have represented a diverse array of anomalies from different time periods, but almost every dataset we encountered used a different set of features, which would have made comparison and analysis impossible with our current experimental framework.

We had to write code to load the datasets properly because of how we chose to evaluate the data. We also to re-implement several of methods because there were no implementations online that worked with our data. Examples of these were the Fence-GAN [9], which had old code online, but it only worked with images and no longer ran on newer hardware. We reimplemented this from scratch only to find that it would take too long to train to a reasonable level of accuracy (which is why it is not included in the table). The autoencoder method was another example of this, where the papers describing it for anomaly detection lacked detail and a sample implementation, so we had to guess about how to fill in some of the blanks. Another challenge is the low accuracy and F1-score performance by the OC-SVM and entropy OCSVM, it lead our team to consider which entropy or one class method would fit better on such type of BGP dataset.

REFERENCES

- [1] Hu J. Ahmed H., Mahmood A.N. 2016. A survey of network anomaly detection techniques. *Journal of Network and Computer Applications* 60 (Jan. 2016), 19–31. <https://doi.org/10.1016/j.jnca.2015.11.016>
- [2] Bahaa Al-Musawi, Philip Branch, and Grenville Armitage. 2016. BGP anomaly detection techniques: A survey. *IEEE Communications Surveys & Tutorials* 19, 1 (2016), 377–396.
- [3] Anisa Allahdadi, Ricardo Morla, and Rui Prior. 2017. A Framework for BGP Abnormal Events Detection. (8 2017). <http://arxiv.org/abs/1708.03453>
- [4] Raghavendra Chalapathy, Aditya Krishna Menon, and Sanjay Chawla. 2018. Anomaly Detection using One-Class Neural Networks. *CoRR* abs/1802.06360 (2018). arXiv:1802.06360 <http://arxiv.org/abs/1802.06360>
- [5] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. 2014. Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling. (12 2014). <http://arxiv.org/abs/1412.3555>
- [6] Haeri S. Trajković L. Ding Q., Li Z. 2018. Application of Machine Learning Techniques to Detecting Anomalies in Communication Networks: Datasets and Feature Selection Algorithms. *Advances in Information Security* 70 (April 2018). https://doi.org/10.1007/978-3-319-73951-9_3
- [7] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative Adversarial Networks. arXiv:1406.2661 [stat.ML]
- [8] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long Short-Term Memory. *Neural Computation* 9, 8 (11 1997), 1735–1780. <https://doi.org/10.1162/neco.1997.9.8.1735>
- [9] Cuong Phuc Ngo, Amadeus Aristo Winarto, Connie Kou Khor Li, Sojeong Park, Farhan Akram, and Hwee Kuan Lee. 2019. Fence GAN: Towards Better Anomaly Detection. *Proceedings - International Conference on Tools with Artificial Intelligence, ICTAI 2019-November* (4 2019), 141–148. <http://arxiv.org/abs/1904.01209>
- [10] Darsh Patel, Kathiravan Srinivasan, Chuan-Yu Chang, Takshi Gupta, and Aman Kataria. 2020. Network Anomaly Detection inside Consumer Networks—A Hybrid Approach. *Electronics* 9, 6 (Jun 2020), 923. <https://doi.org/10.3390/electronics9060923>
- [11] Lukas Ruff, Robert Vandermeulen, Nico Goernitz, Lucas Deecke, Shoaib Ahmed Siddiqui, Alexander Binder, Emmanuel Müller, and Marius Kloft. 2018. Deep one-class classification. In *International conference on machine learning*. 4393–4402.

- [12] Lukas Ruff, Robert Vandermeulen, Nico Goernitz, Lucas Deecke, Shoaib Ahmed Siddiqui, Alexander Binder, Emmanuel Müller, and Marius Kloft. 2018. Deep One-Class Classification (*Proceedings of Machine Learning Research*, Vol. 80), Jennifer Dy and Andreas Krause (Eds.). PMLR, Stockholmsmässan, Stockholm Sweden, 4393–4402. <http://proceedings.mlr.press/v80/ruff18a.html>
- [13] Bernhard Schölkopf. 1998. SVMs - A practical consequence of learning theory. *IEEE Intelligent Systems and Their Applications* 13, 4 (7 1998), 18–21. <https://doi.org/10.1109/5254.708428>
- [14] David M.J. Tax and Robert P.W. Duin. 2004. Support Vector Data Description. *Machine Learning* 54, 1 (01 Jan 2004), 45–66. <https://doi.org/10.1023/B:MACH.0000008084.60811.49>
- [15] Li X. Xu M. 2020. BGP Anomaly Detection Based on Automatic Feature Extraction by Neural Network. *2020 IEEE 5th Information Technology and Mechatronics Engineering Conference (ITOEC)* (2020). https://doi.org/10.1007/978-3-319-73951-9_3
- [16] Yue Zhao, Zain Nasrullah, and Zheng Li. 2019. PyOD: A Python Toolbox for Scalable Outlier Detection. *Journal of Machine Learning Research* 20, 96 (2019), 1–7. <http://jmlr.org/papers/v20/19-011.html>