

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/321527157>

Deep Learning for Detection of BGP Anomalies

Conference Paper · September 2017

CITATION

1

READS

330

3 authors, including:



Obradovic Slobodan

University of East Sarajevo

40 PUBLICATIONS 48 CITATIONS

[SEE PROFILE](#)



Emina Junuz

University Dzemal Bijedic Mostar

15 PUBLICATIONS 10 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



CCMLL - Center for Curricula Modernization and Lifelong Learning - TEMPUS IV [View project](#)



TEMPUS UM_JEP-19015-2004 „Quality Assurance by Accreditation [View project](#)

Deep Learning for Detection of BGP Anomalies

Marijana Cosovic¹, Slobodan Obradovic¹, Emina Junuz²

¹Faculty of Electrical Engineering, University of East Sarajevo, Istocno Sarajevo, BiH

²Faculty of Information Technology, Dzemal Bijedic University, Mostar, BiH

marijana.cosovic@etf.unssa.rs.ba, slobodan.obradovic@gmail.com,
emina@edu.fit.ba

Abstract. The Internet uses Border Gateway Protocol (BGP) for exchange of routes and reachability information between Autonomous Systems (AS). Hence, BGP is subject to anomalous traffic that can cause problems with connectivity and traffic loss. Routing Table Leaks (RTL) are considered anomalous in the sense that they can disrupt Internet routing and cause slowdowns of varying severity, which leads to packet delivery reliability issues. Deep learning, a subfield of machine learning, could be applied in detection of BGP anomalies. Studying RTL events are of interest to network operators and researchers alike. In this paper we consider datasets of several RTL events, all of which caused large-scale Internet outages. We use artificial neural network (ANN) models based on a backpropagation algorithm for RTL event classification.

Keywords: Machine learning·Deep learning·Anomaly detection·BGP

1 Introduction

The Internet can be viewed as a graph of nodes in which autonomous systems, collections of routers with same routing policies, are represented by nodes, while the connection between the nodes are data paths used for exchanging reachability information between the ASs. BGP is a protocol that facilitates this exchange [1]. The Routing Information Service (RIS) project initiated by the Réseaux IP Européens Network Coordination Centre (RIPE NCC) is collecting routing data from Remote route collectors (RRC) positioned predominantly at Internet exchange points. RIS raw data come in two different type of files: all BGP packets created every five minutes and a complete BGP routing table that is created every eight hours [2].

By studying BGP packets files and, in particular, by extracting BGP update messages from them as they contain important reachability information, we can study connectivity disruption in the Internet during anomalous events. RTL events are in general initiated by router misconfigurations and, although not malicious in nature, can cause connectivity and traffic loss.

Machine learning techniques have been employed in anomaly classification tasks [3-6]. Deep learning, part of machine learning, has been used extensively in voice and

image recognition, language modelling, and information retrieval, amongst others, and has impacted the wide range of information processing tasks [7]. Routing data could be considered as time series data since data points are indexed in time order. Detection of anomalies in time series data has employed deep learning techniques in the past. ANNs are systems that can be trained to recognize patterns in data and classify anomalous from regular data instances [8], [9]. Routing data could be used to analyze past anomalous events and aid in classification of future anomalous events.

The paper is organized as follows. In Section 2, we describe ANNs. Introduction of RTL and particular RTL events are discussed in Section 3. In addition, extraction of BGP features from the datasets concludes Section 3. Classification methodology and used performance measures are discussed in Section 4. We conclude with Section 5.

2 ANN - Deep Learning

Application of Artificial Neural Networks (ANN) is present in the detection of anomalies [8], [9]. ANN are preferably developed to mimic basic biological systems and to learn based on examples in the way humans do. Neural networks learn gradually from the interdependence of data input properties, which can be linear or non-linear in nature. Neural network usage in supervised learning implies that input data are labeled; hence, it is known in advance which class they belong to. Based on a comparison between the output of the neural network and the target function, during the training process, ANN adjusts the weights as shown in Fig. 1.

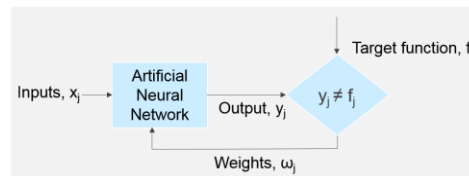


Fig. 1. Artificial Neural Network

Artificial neural networks can be classified as Feedforward or Feedbackward structures, depending on the direction of propagation of the information. The Feedbackward structure of neural networks refers to the spread of information backwards. When the input vector is applied to the input layer of the neural network it propagates through the network throughout all its layers, and it generates output values by using the output layer of the network. The output values are compared with a desired target function, and for each of the neurons in the output layer the difference is calculated. Further information about these differences propagate backwards until all the neurons in the neural network are affected by the difference of the original and the target output value. The value of the weighting factors are determined by the optimization technique (typically a minimizing of the loss function with respect to the weights in the network) which determines the weighting factors such that the loss function is minimized.

ANN are simple mathematical methods made up of basic processing elements called neurons. The structures of neural networks differ in the number of layers used. Between the first and the last layer of neural networks there are hidden layers: usually one hidden layer in simpler networks and more hidden layers in complex neural networks.

The architecture of the neural network is engaged in specific neuronal connectivity in a whole. Usually, the number of neurons in the input layer is equal to the number of features (number of columns in the feature matrix). Each neuron has one input, and all the outputs are connected to all neurons of the next layer, as shown in Fig. 2. When using a neural network for classification the output layer can have one or more neurons, depending on whether it is binary or multi-class classification. The most commonly used functions for the output neuron modeling are sigmoid or normalized exponential [10] functions.

Perceptron is a neuron model type developed in the original neural networks, in which each neuron has a number of inputs (x_j) associated with corresponding weight factors (ω_j), which show the effect of a particular input on the output. Thus, the output neuron classifies information by comparing the value of the sum (1) and the threshold value, which is a parameter of the neuron.

$$\sum_j \omega_j x_j \quad (1)$$

Modeling of neurons with perceptron has the following disadvantage: a small change in the weight factor of any perceptron can lead to a sudden change in its output. This in turn can lead to a complicated change in the rest of the network, which may be difficult to control. The most commonly used artificial neuron model, which solves the aforementioned problem, is the sigmoid neuron, shown by the following expression:

$$\frac{1}{1 + \exp(-\sum_j \omega_j x_j - b)} \quad (2)$$

where ω_j are weighting factors, x_j are input neurons and b is bias. It turns out that a change in the output of sigmoid neurons linear function of changes in weighting factors and bias. In this way it is easier to determine how changes in weighting factors and bias may influence the change of the output neuron; hence, the neural network could be considered more resilient to changes of data and the ability to learn.

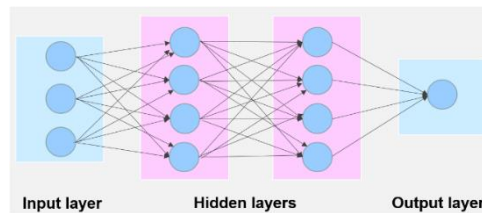


Fig. 2. Architecture of ANN with four layers: one input layer, two hidden layers and one output layer

3 Routing Leaks

The BGP routing system is subject to frequent incidents that result in significant interruptions of Internet connectivity. Many of the events that cause connectivity issues are classified as routing leaks. It is often unclear what is meant by that term. Based on research of actual events on the Internet, which can be of use to network operators and Internet users, the authors in [11] define routing leaks as a propagation of announced paths beyond the intended scope. This means that the BGP path announcement from one AS to another in some way violates the routing agreements between a sending AS, a receiving AS or any transit AS. The consequence of routing leaks is traffic redirection through a path not originally planned, and thus, various malicious attacks from analyzing data to eavesdropping could be performed. The most common reasons why routing leaks occur are errors in the router's configuration [12].

3.1 Routing Table Leak Events

In this paper we consider the following routing table leak events: Routing Leak AS9121 [13], AWS Route Leak [14], Telekom Malaysia AS4788 Route Leak [15], and Indosat Routing Table Leak [16], all of which showed an increased number of announced IP prefixes throughout the duration of the events.

AS9121 Routing Table Leak.

The AS9121 Routing Table Leak took place on December, 24 2004. AS9121 announced to other AS's through BGP sessions that were used to reach almost 70% of all prefixes, which at that time amounted to more than 106k prefixes. As a result, the data of tens of thousands of networks were either lost or diverted. AS9121 started to announce prefixes to its neighbors around 9:20 GMT, and the event lasted until just after 10:00 GMT. AS9121 continued announcing the prefixes for the rest of the day. The prefix announcement rate reached a second peak at 19:47 GMT. The number of announced IP prefixes during the routing leak event is shown in Figure 4.10 (a). Picture 4.10 (b) shows that during the routing leak event, the maximum edit distance (the measure of similarity between two ASPATH attributes) increases within one minute. This could indicate that the choice of the paths differed from the common ones and it was sign of disruption between commonly connected ASs.

AWS Route Leak.

The AWS Route Leak started at 17:10 UTC on April, 22 2016 and affected a large number of ASs and prefixes. Loss of traffic and connectivity were present since networks with high traffic prefixes, such as Google, Amazon, and Twitter, were affected, amongst others. The event occurred due to maintenance issues on Innofield AG (AS 200759) that is connected to Swiss Internet eXchange (SwissIX). Innofield AG normally announces one IPv4 and IPv6 prefix to SwissIX. During maintenance reactiva-

tion of BGP sessions, AS 200759 distributed prefixes belonging to Amazon as belonging to private AS 65021. Prefix announcements were propagated through AS 6939 Hurricane Electric (HE) that peers at SwissIX. This resulted in a redirection of traffic passing through HE to a private AS, and hence, it compromised the reachability of Amazon AS. Since the event was widespread and likely caused by a misconfigured route optimizer, we observed an increase in announced IP prefixes at CIPX, as shown in Fig. 2.

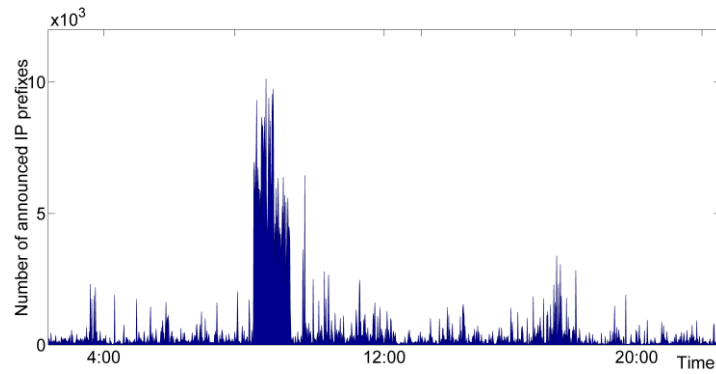


Fig. 3. Number of announced Network Layer Reachability Information (NLRI) prefixes during AS9121 Routing Leak Event as observed on RIPE Route Collector rrc04, CIPX

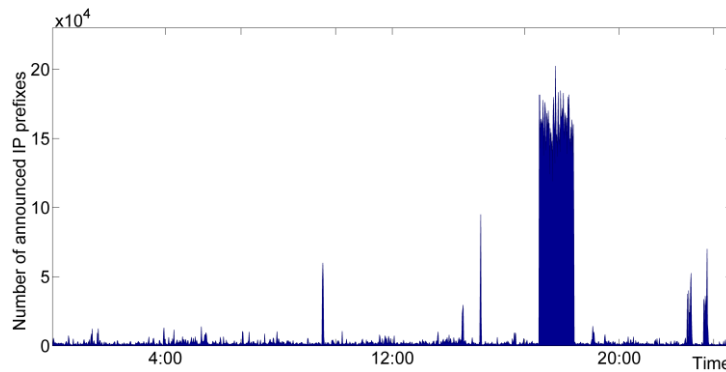


Fig. 4. Number of announced NLRI prefixes during AWS Routing Leak Event as observed on RIPE Route Collector rrc04, CIPX

Telecom Malaysia Route Leak.

The Malaysian Telecom (AS 4788) leaked one third of all IP prefixes in the global routing table to the backbone provider Level3 (AS 3549). The event, triggered by routers misconfiguration at Telecom Malaysia, started on June, 12 2015 at 8:43 UTC and lasted until 11:45 UTC. Level3 (AS 3549) propagated traffic from its peers and customers via Telecom Malaysia, which was not capable of handling traffic volume, re-

sulting in major packet loss and performance degradation. The performance degradation was especially pronounced between the Asia Pacific region and the rest of the Level3 network. Fig. 5 shows an increased number of announced IP prefixes (left) and also an increase of maximum AS-PATH length (right) for the duration of the route leak event.

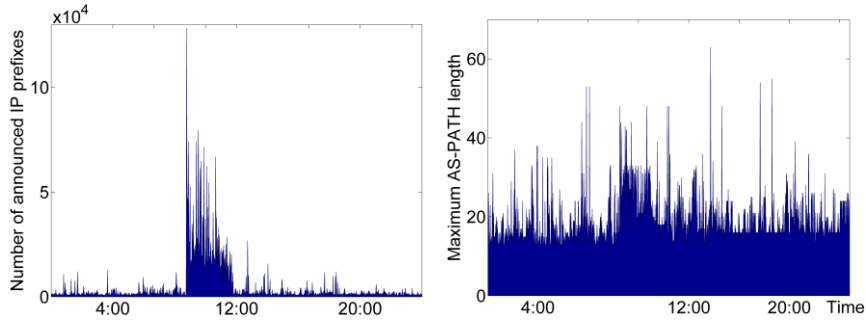


Fig. 5. Number of announced NLRI prefixes (left) and maximum AS-path length (right) during Telecom Malaysia Routing Leak Event as observed on RIPE Route Collector rrc04, CIPX

Indosat Routing Table Leak.

The Indosat routing table leak occurred on April, 2 2014. At the time of the event the global routing table consisted of nearly half a million routes. AS 4761 (Indosat) leaked around 320,000 routes, which happened during scheduled maintenance starting at 18:25 UTC. The reason behind Indosat originating prefixes that were not assigned to it is assumed to be that BGP was redistributed with bad upstream filtering. This inadvertent error had an impact that was observed on various route collectors through an increase of announced IP prefixes, as shown in Fig.6 (left). Several hundreds of those prefixes were widely accepted, and services of some networks such as Akamai, a leading content delivery network (CND) and cloud service provider, were disrupted.

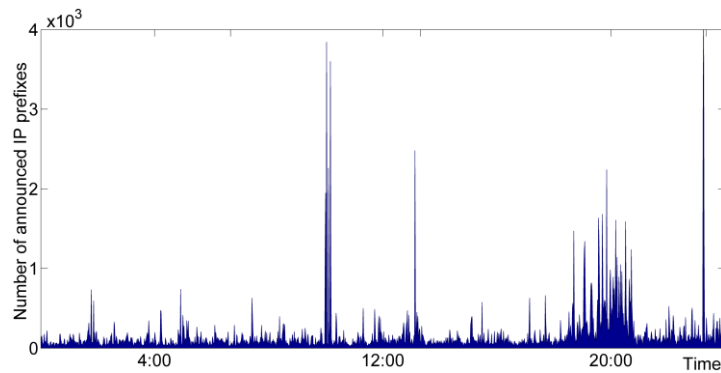


Fig. 6. Number of announced NLRI prefixes during Indosat Routing Leak Event as observed on RIPE Route Collector rrc04, CIPX

3.2 Routing Leak Datasets

We obtain datasets from the RIPE NCC that collects Internet routing data by using Routing Information Service (RIS) Remote Route Collectors (RRC) positioned in various location throughout the world. Since the effects of all events considered in this paper and presented in Tab. 1 caused globally visible connectivity issues, we have used routing updates collected at RRC located in CIPX, Geneva.

We have used BGP update messages during the occurrence of the routing leak events stored in MRT format described in [17]. We have observed BGP update messages during a five day period, two days before and two days after the actual event. After MRT to ASCII conversion, python code was written in order to extract information from the datasets. We have observed fifteen volume and AS-PATH features on a minute level during five days period, hence producing a feature matrix of 7200x15 size.

The volume features that we have observed are: the number of BGP messages announcing new routes, the number of BGP messages withdrawing already existing routes, the number of announced IP prefixes (Fig. 3, 4, 5, and 6), the number of withdrawn IP prefixes, the number of duplicate announced messages, the number of duplicate withdrawn messages, the number of implicitly withdrawn messages, the number of BGP messages which NLRI originates from Exterior Gateway Protocol (EGP), the number of BGP messages which NLRI originates from Interior Gateway Protocol (IGP), and the number of BGP messages which NLRI originates from unknown sources.

Duplicate announcements and withdrawal messages are defined as BGP update messages that announce the same combination of IP prefix and AS-PATH attribute that has previously been announced. Implicit withdrawal implies that the same IP prefix has been announced with a different AS-PATH attribute, hence it is an implicit withdrawal of a previous announcement (same IP prefix but different AS-PATH).

The features we computed based on AS-PATH attribute are: the average length of AS-PATH attribute, the maximum length of AS-PATH attribute, the average length of unique AS-PATH attribute, the average edit distance, and the maximum edit distance. While extracting information from AS-PATH attribute, we considered regular and unique AS-PATH's. We also considered AS-PATHs as a string of ASNs (autonomous system number) and computed the similarity of two adjacent AS-PATHs by finding their edit distance [18].

Features belong to three types, namely, continuous, categorical and binary. All of the volume features belong to the continuous type since features may have an infinite number of values. On the other hand, features derived from the AS-PATH attribute may have a finite number of values and hence, are categorical. The class feature is of the binary type: given volume and AS-PATH features, we either have anomalous instances or not.

We have labeled all 7200 time instances (described by 15 features) as either belonging to anomalous or regular class in accordance to the information regarding beginning, duration and end of each of the events. We have referred to several sources [19] in order to label our data as correctly as possible.

Considering that global routing tables increase in size from the time of the first event, we needed to normalize feature values to account for Internet size growth. Normalization is done such that each feature vector has a zero mean and a standard deviation of one [18]. We also performed feature discretization for the features of the continuous type prior to training the neural network. We did not encounter any missing data during the four events observed, although we did have an increased number of outliers in the case of Indosat RTL dataset, which can be observed in Fig. 6.

Table 1. RTL events

Dataset	Regular Class	Anomaly Class	Number of features
AS9121 RTL	7121	79	15
AWS RTL	7085	115	15
Malaysian Telecom RTL	7018	182	15
Indosat RTL	7050	150	15

4 Classification of Routing Leaks

4.1 Methodology

We used the Keras Python library with the Theano backend for development and evaluation of deep learning models. Models, based on a backpropagation algorithm for training of fully connected multilayer perceptron (MLP) neural networks, are defined as sequences of layers: an input layer, hidden layers and an output layer. Only for the first layer in the sequence the shape of the input data needs to be specified. In Keras, using Dense class is one of the ways to define fully connected layers. Network weights can be initialized to random numbers using either uniform or Gaussian distribution. Use of appropriate activation function allows for better training of the network [19]. Traditionally, sigmoid and tanh activation functions are used, but the authors in [19] have shown that better performance can be achieved using a rectifier activation function. In the output layer we use a sigmoid function as we are dealing with binary classification. We use 10-fold cross validation for determining accuracy on the test dataset, and as we increase the number of hidden layers beyond two, classification accuracy decreases. We found that a neural network with two hidden layers is the optimal model for routing table leak datasets. Using either too few or too many neurons in the hidden layers may result in problems of underfitting and overfitting, respectively. General guidelines for determining number of neurons within each hidden layer are used. We selected neural network architecture based on trial and error, but in accordance with the following general guidelines: the number of neurons in hidden layers should be between the sizes of input and output layers, and they should be the sum of $2/3$ of the input layer neurons and output layer neurons. Hence, we trained the neural network with two dense hidden layers with 15 and 10 neurons, respectively.

4.2 Performance Measures

The performance measures employed in this paper, needed for comprehensive comparison of different deep learning models, are accuracy, f-measure, the Matthews Correlation Coefficient (MCC), the area under Precision-Recall (PR), and the area under Receiver Operating Characteristics (ROC). Accuracy, considering our datasets are highly imbalanced (Tab. 1), might not be the most accurate performance measure. This is due to the fact that misclassification should have different costs associated with points belonging to either the regular or anomalous class. Accuracy is defined as the ratio of points belonging to the regular/anomalous class that are classified as regular/anomalous and the total number of points in the dataset. In order to define f-measure we first define recall (R) as the ratio of detected anomalous points and all points labeled as anomalous. On the other hand, precision (P) is a ratio of detected anomalous points and all anomalous points. F-measure is given as a double ratio of product of P and R and the sum of P and R. MCC is given by (3) where N is the number of all points and TP is the number of data points classified as anomalous.

$$MCC = \frac{TP / N - PR}{\sqrt{PR(1 - P)(1 - R)}} \quad (3)$$

4.3 Classification Results

We have used a neural network with two hidden layers and obtained the performance measure values shown in Tab. 2. Accuracy is not the best approach to compare classification of different events, as the datasets are highly imbalanced. Tab. 1 shows that Malaysian Telecom RTL has the largest set of data labeled as anomalous – 182 compared to the AS 9121 RTL event in which only 79 instances are labeled as anomalous. The Indosat RTL event shows the worst performance of all datasets, and we can contribute that to noise in the dataset (Fig. 6). We have used under- and oversampling techniques as in [6] to balance regular and anomalous instances in the RTL datasets. In the case of oversampled and undersampled datasets, their imbalance ratio is around 1, meaning the classes are balanced; hence, accuracy and f-measure are approximately the same values.

Oversampling techniques are algorithms that create additional instances of the class that is represented with a smaller number of instances in the dataset. We used Synthetic Minority Oversampling Technique (SMOTE), Support Vector Machine (SVM)-SMOTE, Borderline1-SMOTE, Borderline2-SMOTE, Adaptive Synthetic Sampling (ADASYN), and Random Oversampling (ROS) algorithms as discussed in [6]. By using balancing techniques of the datasets, we have achieved better performance measures as shown in Tab. 3. The best results were achieved using the SVM-SMOTE oversampling technique for AS9121 RTL, AWS RTL, and Indosat RTL, while the Malaysian Telecom RTL dataset, when oversampled by ROS algorithm, had the best performance measure that was better by a small margin than when oversampled by the SVM-SMOTE algorithm.

Undersampling techniques are algorithms that remove instances from the dataset that belong to the more represented class. We used ten undersampling algorithms, namely, Near Miss-1, Near Miss-2, Near Miss-3, Tomek Links, Cluster Centroids, One-sided selection, Random undersampling, Edited Nearest Neighbours, Neighbourhood Cleaning Rule, and Condensed Nearest Neighbours. By using undersampling balancing techniques of the datasets, we have achieved better performance measures as shown in Tab. 4. When comparing Tab. 3 and Tab. 4, the values of performance measures are greater in the case of oversampling techniques, and this is due to possible overfitting. The best results were achieved using the RUS undersampling technique for AS9121 RTL and AWS RTL, while Indosat RTL and Malaysian Telecom RTL datasets, when under-sampled by the Near-Miss 1 algorithm, had the best performance measure that was only better by a small margin than when under-sampled by the RUS algorithm.

Table 2. Performance measures of the original RTL events

Dataset	Acc	F-measure	MCC	ROC	PR
AS9121 RTL	0.99375	0.945	0.942	0.998	0.946
AWS RTL	0.99431	0.808	0.807	0.961	0.848
Malaysian Telecom RTL	0.9925	0.852	0.848	0.979	0.883
Indosat RTL	0.93056	0.753	0.707	0.897	0.802

Table 3. Performance measures of RTL events using oversampling techniques

Dataset	Acc	F-measure	MCC	ROC	PR
AS9121 RTL	0.99816	0.998	0.996	0.999	0.999
AWS RTL	0.99167	0.992	0.983	0.994	0.984
Malaysian Telecom RTL	0.98953	0.990	0.979	0.995	0.994
Indosat RTL	0.92087	0.923	0.844	0.958	0.940

Table 4. Performance measures of RTL events using undersampling techniques

Dataset	Acc	F-measure	MCC	ROC	PR
92AS9121 RTL	0.98734	0.987	0.975	0.999	0.999
AWS RTL	0.96087	0.960	0.923	0.979	0.986
Malaysian Telecom RTL	0.95055	0.950	0.901	0.975	0.981
Indosat RTL	0.88333	0.878	0.770	0.927	0.948

5 Conclusion

We have developed a model for anomaly detection based on artificial neural networks with two hidden layers, which are optimal since with choosing additional hidden layers performance indices deteriorated. We used a cross-validation technique to determine

the number of neurons in each of the layers. Balancing techniques (dataset over-sampling and undersampling) were employed as the original datasets are highly imbalanced. Classification of the Indosat RTL dataset achieved the worst performance measures due to noise in the dataset. We concluded in routing table leak datasets presented in the study that employing volume and AS-PATH features from BGP update messages could lead to reliable classification of RTL events.

References

1. Y. Rekhter, T. Li, S. Hares, "A Border Gateway Protocol 4 (BGP-4)," RFC 4271, IETF, 2006 [Online]. Available: <http://ietf.org/rfc/rfc4271>
2. RIPE RIS raw data [Online]. Available: <http://www.ripe.net/data-tools/stats/ris/ris-raw-data>.
3. M. Ćosović, S. Obradović, Lj. Trajković, "Performance evaluation of BGP anomaly classifiers," Proceedings of the International Conference on Digital Information, Networking and Wireless Communication, Moscow, Russia, Feb. 2015, pp. 115–120.
4. M. Cosovic, S. Obradovic, Lj. Trajkovic, "Classifying anomalous events in BGP datasets," in Proceedings of the 29th Annual IEEE Canadian Conference on Electrical and Computer Engineering (CCECE 2016), Vancouver, Canada, May 2016, pp. 697-700.
5. M. Cosovic, S. Obradovic, "Ensemble methods for classifying BGP anomalies," Industrial Technologies, ISSN: 13149911, vol. 4, no. 1, pp. 12-20, June 2017.
6. M. Cosovic, S. Obradovic, "BGP anomaly detection with balanced datasets," Tehnički vjesnik/Technical Gazette, vol. 25, no. 3, in press, June 2018.
7. Li Deng and Dong Yu. 2014. Deep Learning: Methods and Applications. Found. Trends Signal Process. 7, 3–4 (June 2014), 197-387. DOI: <http://dx.doi.org/10.1561/20000000039>
8. H. A. Dau, V. Ciesielski, A. Song, "Anomaly Detection Using Replicator Neural Networks Trained on Examples of One Class," In Proceedings of the 10th International Conference on Simulated Evolution and Learning (SEAL 2014), vol. 8886, Springer-Verlag New York, Inc., New York, NY, USA, pp. 311-322, 2014.
9. Z. Jadidi, V. Muthukumarasamy, E. Sithirasanen, M. Sheikhan, "Flow-Based Anomaly Detection Using Neural Network Optimized with GSA Algorithm," In Proceedings of the 33rd IEEE International Conference on Distributed Computing Systems Workshops (ICDCSW '13), Washington, DC, USA, pp. 76-81, 2013.
10. Christopher M. Bishop. 2006. Pattern Recognition and Machine Learning (Information Science and Statistics). Springer-Verlag New York, Inc., Secaucus, NJ, USA.
11. <https://tools.ietf.org/html/draft-sriram-route-leak-problem-definition-00> <https://www.rfc-editor.org/rfc/rfc7908.txt>
12. R. Mahajan, D. Wetherall, T. Anderson, "Understanding BGP misconfiguration," In Proceedings of the Conference on Applications, technologies, architectures, and protocols for computer communications (SIGCOMM '02), New York, NY, USA, pp. 3-16, 2002.
13. A. C. Popescu, B. J. Premore, T. Underwood. (May 19, 2005). Anatomy of a Leak: AS9121. Renesys Corporation. Manchester, NH, USA. [Online]. Available: <http://research.dyn.com/content/uploads/2013/05/renesys-nanog34.pdf>.
14. (May. 22, 2017) North American Network Operators Group Mailing List [Online]. Available: <https://mailman.nanog.org/pipermail/nanog/2016-April/085410.html>
15. (May. 22, 2017) North American Network Operators Group Mailing List [Online]. Available: <https://mailman.nanog.org/pipermail/nanog/2015-June/076187.html>
16. (May. 22, 2017) North American Network Operators Group Mailing List [Online]. Available:

<https://mailman.nanog.org/pipermail/nanog/2014-April/065920.html>

17. T. Manderson, “Multi-threaded routing toolkit (MRT) Border Gateway Protocol (BGP) routing information export format with geo-location extensions,” RFC 6397, IETF, [Online]. Available: <http://www.ietf.org/rfc/rfc6397.txt>.
18. Yann LeCun, Léon Bottou, Genevieve B. Orr, and Klaus-Robert Müller. 1998. Efficient BackProp. In *Neural Networks: Tricks of the Trade*, this book is an outgrowth of a 1996 NIPS workshop, Genevieve B. Orr and Klaus-Robert Müller (Eds.). Springer-Verlag, London, UK, 9-50.
19. Vinod Nair and Geoffrey E. Hinton. 2010. Rectified linear units improve restricted Boltzmann machines. In *Proceedings of ICML*. 807-814.