# Application of Machine Learning Techniques for Detecting Anomalies in Communication Networks

Qingye Ding
qingyed@sfu.ca

Communication Networks Laboratory
http://www.ensc.sfu.ca/~ljilja/
School of Engineering Science
Simon Fraser University

SFU

SIMON FRASER
UNIVERSITY

# Roadmap

- Introduction

- Border Gateway Protocol datasets

- Extraction of features from BGP update messages

- Performance metrics

- Long Short-Term Memory

- Comparison of classification algorithms

- Discussion

- Future work and conclusion

- References

# Roadmap

- **Introduction**

- Border Gateway Protocol datasets

- Extraction of features from BGP update messages

- Performance metrics

- Long Short-Term Memory

- Comparison of classification algorithms

- Discussion

- Future work and conclusion

- References

Application of Machine Learning Techniques for Detecting Anomalies in Communication Networks

SFU

SIMON FRASER UNIVERSITY

# Motivation

- The Internet is a critical asset of information and it contains multiple Autonomous Systems (ASes)

- An AS is a collection of Internet Protocol (IP) routing prefixes administrated by a single domain

- Border Gateway Protocol (BGP) plays an essential role in routing data between ASes

- Cyber attacks and threats significantly impact the Internet performance

# Border Gateway Protocol

- Forwards IP traffic between Autonomous Systems (ASes)

- BGP 4: a standard for exchanging information among the Internet Service Providers (ISPs)

- Relies on the Transport Control Protocol (TCP) to establish a connection between routers

- Exchange the update message to advertise routing information:
  - an available route
  - withdraw multiple unavailable routes

# Sample of BGP update message

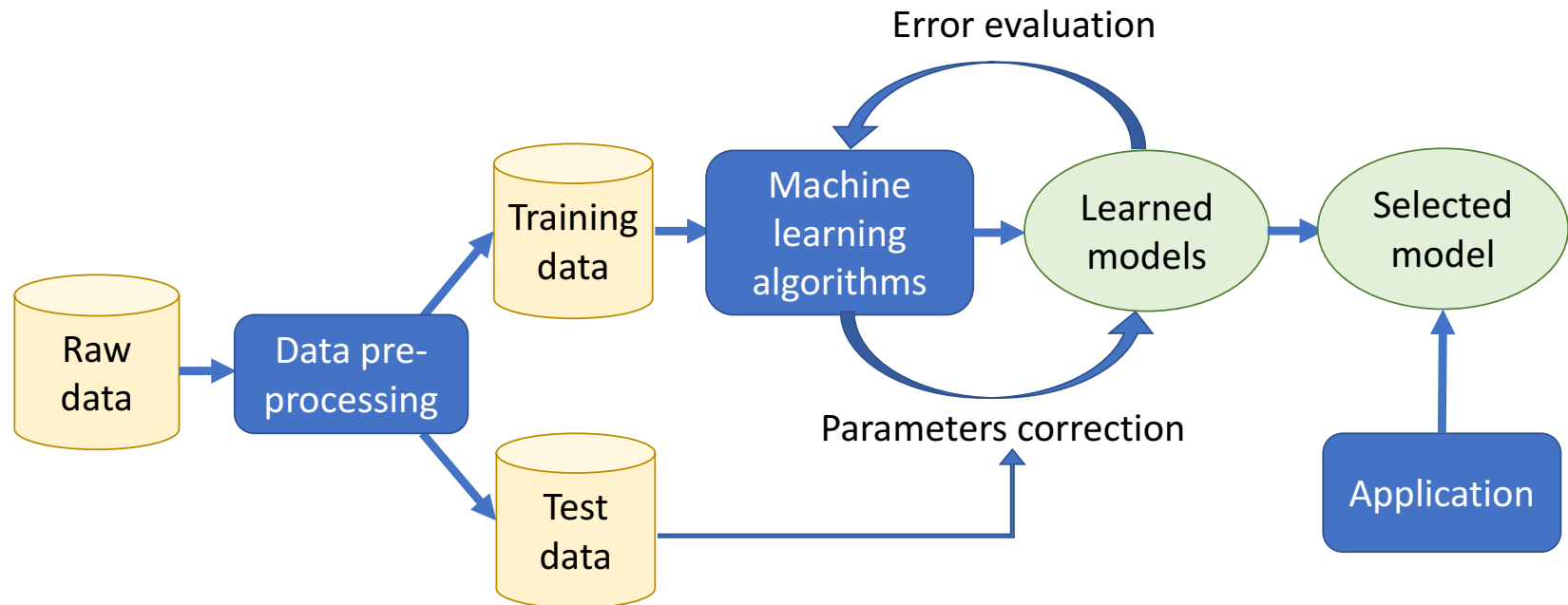| Field | Value |
|---|---|
| TIME | 2003 1 24 00:39:53 |
| TYPE | BGP4MP/BGP4MP_MESSAGE AFI_IP |
| FROM | 192.65.184.3 |
| TO | 193.0.4.28 |
| BGP PACKET TYPE | UPDATE |
| ORIGIN | IGP |
| AS-PATH | 513 3320 7176 15570 7246 7246 |
| NEXT-HOP | 192.65.184.3 |
| ANNOUNCED NLRI PREFIX | 198.155.189.0/24 |
| ANNOUNCED NLRI PREFIX | 198.155..0/24 |

IGP: Interior Gateway Protocol

NLRI: Network Layer Reachability Information

Application of Machine Learning Techniques for Detecting Anomalies in Communication Networks
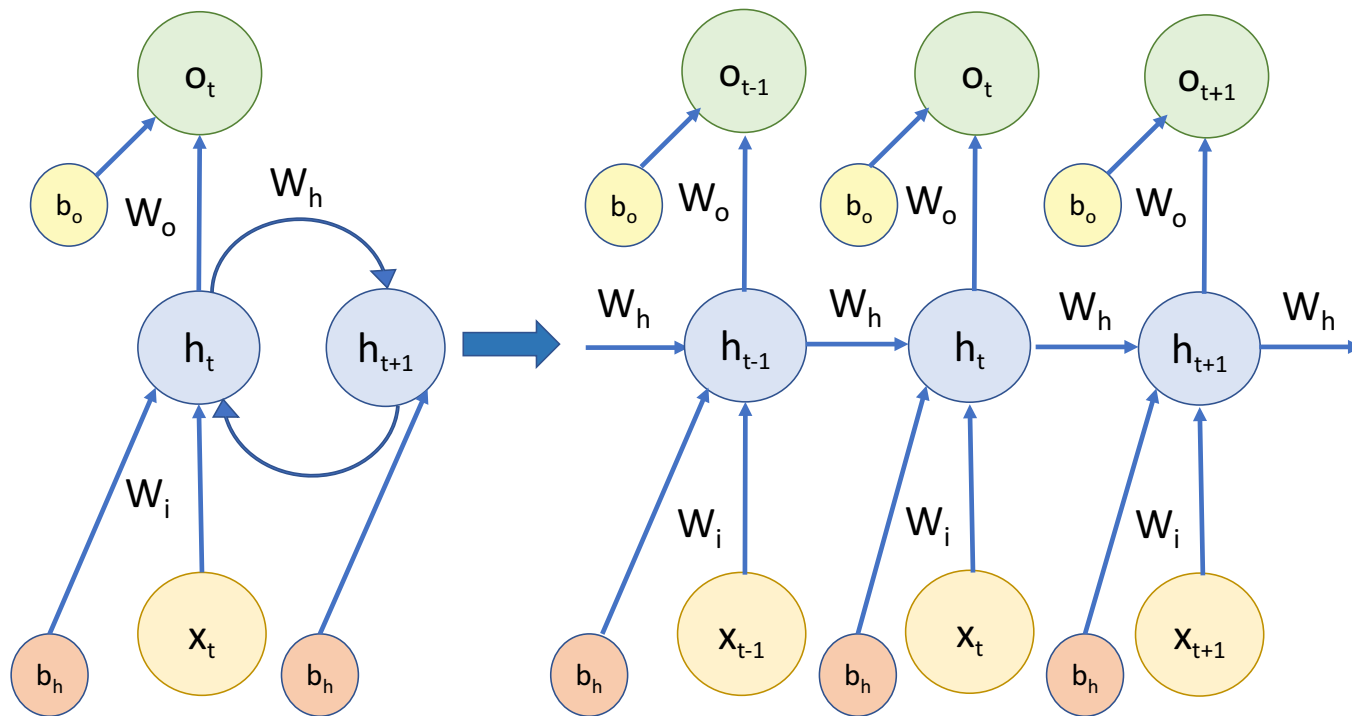
# Machine learning techniques

- Unsupervised learning:
  - aims to learn a function that represents the underlying structure from "unlabeled" data
  - motivation: labeled data is difficult to obtain
  - data clustering

- Supervised learning:
  - trains data based on the observation to predict labels for new events
  - Long Short-Term Memory, Support Vector Machine, Naïve Bayes, Decision Tree, and Extreme Learning Machine

# Typical procedure for machine learning

Application of Machine Learning Techniques for Detecting Anomalies in Communication Networks

SFU

SIMON FRASER UNIVERSITY

# Recurrent Neural Network: RNN

- Used for sequence recognition, pattern classification, and temporal prediction tasks

# Research contributions

- View detection of BGP anomalies as a classification problem

- Apply Long Short-Term Memory algorithm to develop classification models

- Extract BGP features based on the attributes of BGP update messages

- Create balanced datasets by randomly reducing a subset of regular data points

- Improve classification results emanating from previous studies

- Show feasibility of LSTM for detecting BGP anomalies

# Roadmap

- Introduction

- Border Gateway Protocol datasets

- Extraction of features from BGP update messages

- Performance metrics

- Long Short-Term Memory

- Comparison of classification algorithms

- Discussion

- Future work and conclusion

- References

# BGP anomaly: Slammer

- Attacked Microsoft SQL servers on January 25, 2003

- Generated random IP addresses and replicated itself

- The number of infected machines doubled approximately every 9 seconds

- The update messages consumed most of the routers' bandwidth causing routers to:
  - slow down
  - crash

SFU

SIMON FRASER
UNIVERSITY

# BGP anomaly: Nimda

- Released on September 18, 2001

- Exploited vulnerabilities in the Microsoft Internet Information Services web servers for the Internet Explorer 5

- Three methods of propagation:
  - email messages
  - web browsers
  - file systems
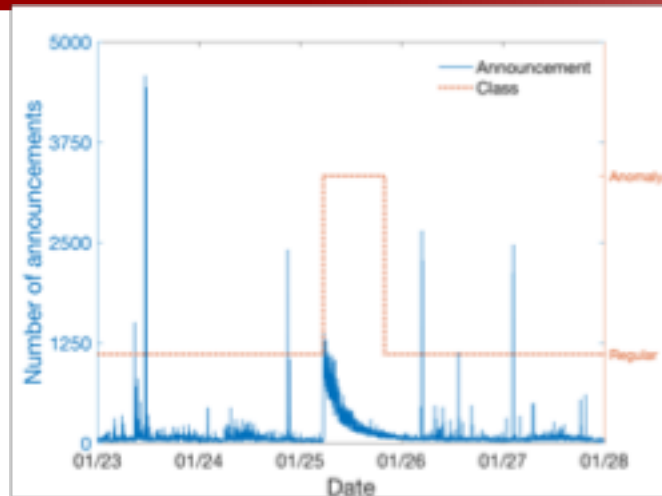
# BGP anomaly: Code Red I

- Attacked web servers on July 19, 2001
- Affected approximately 500,000 IP addresses a day
- Searched for vulnerable servers and replicated itself
- Rate of infection was doubling every 37 minutes

June 19, 2018
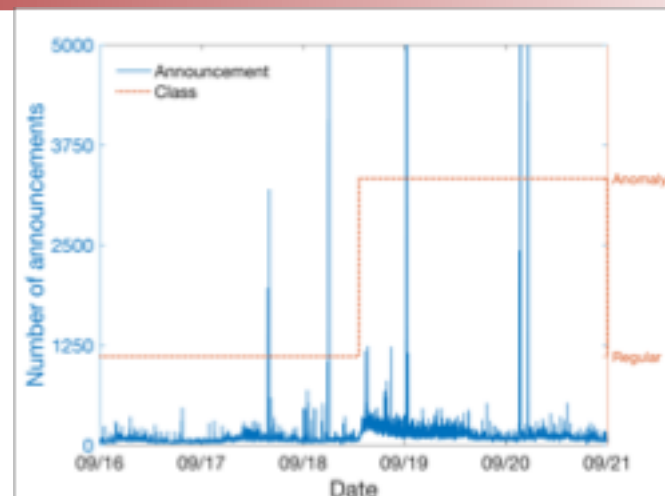
Application of Machine Learning Techniques for Detecting Anomalies
in Communication Networks

14/58

SFU
SIMON FRASER
UNIVERSITY

# BGP anomalies

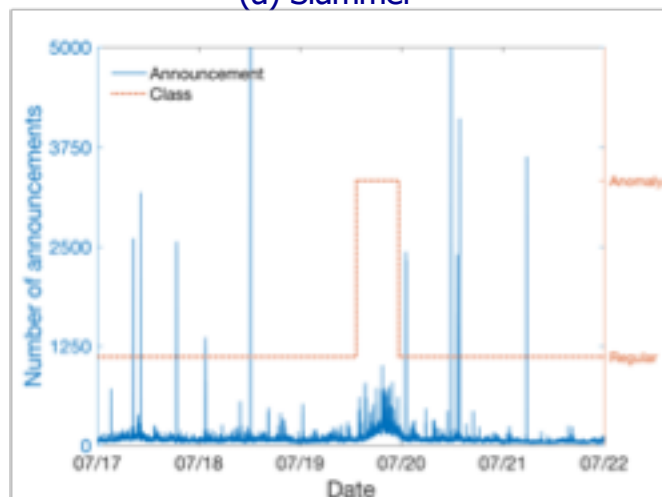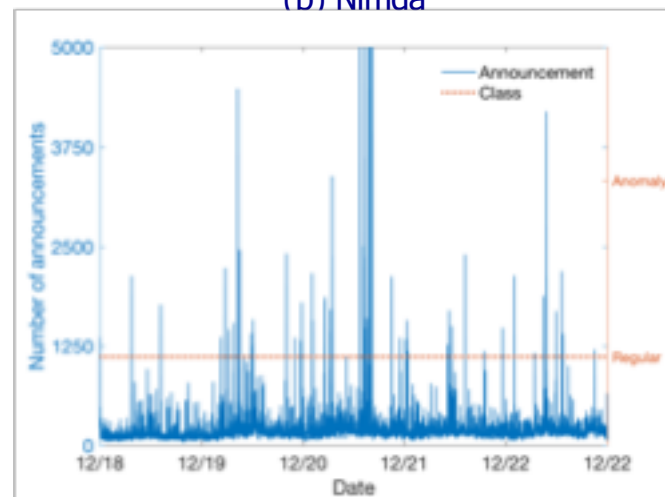| Dataset | Class | Date | | Duration |
|---|---|---|---|---|
| | | Beginning of the event | End of the event | (min) |
| Slammer | Anomaly | 25.01.2003 at 5:31 GMT | 25.01.2003 at 19:59 GMT | 869 |
| Nimda | Anomaly | 18.09.2001 at 13:19 GMT | 20.09.2001 at 23:59 GMT | 3,521 |
| Code Red I | Anomaly | 19.07.2001 at 13:20 GMT | 19.07.2001 at 23:19 GMT | 600 |

GMT: Greenwich Mean Time

# Number of announcements



(a) Slammer

(b) Nimda

(c) Code Red I

(d) Regular data

SFU

SIMON FRASER UNIVERSITY

# BGP datasets

- Réseaux IP Européens (RIPE) Network Coordination Centre:
  - Regional Internet Registry for Europe, Middle East, and parts of Central Asia
  - collects BGP update messages by the remote route collectors (rrc)
  - multi-threaded routing toolkit (MRT) binary format
  - AS 513 (rrc04, CIXP, Geneva, Switzerland)

- BCNET
  - Regular BCNET dataset
  - BCNET location in Vancouver, British Columbia, Canada

CIXP: CERN Internet eXchange Point

# Collected data: unbalanced datasets

5 days

| Regular traffic (3,165 min) | Slammer (869 min) | Regular traffic (3,166 min) |

| Regular traffic (3,679 min) | Nimda (3,521 min) |

| Regular traffic (3,300 min) | Code Red I (600 min) | Regular traffic (3,300 min) |

# Collected data: balanced datasets

regular:anomaly = 1:1

Regular traffic (869 min) | Slammer (869 min)

Regular traffic (3,521 min) | Nimda (3,521 min)

Regular traffic (600 min) | Code Red I (600 min)

# Pre-processing of the collected data

| | Training dataset | Test dataset |
|---|---|---|
| 1 | Slammer and Nimda | Code Red I |
| 2 | Slammer and Code Red I | Nimda |
| 3 | Nimda and Code Red I | Slammer |

- Datasets are concatenated to increase the size of training datasets

SFU

SIMON FRASER
UNIVERSITY

# Roadmap

- Introduction

- Border Gateway Protocol datasets

- Extraction of features from BGP update messages

- Performance metrics

- Long Short-Term Memory

- Comparison of classification algorithms

- Discussion

- Future work and conclusion
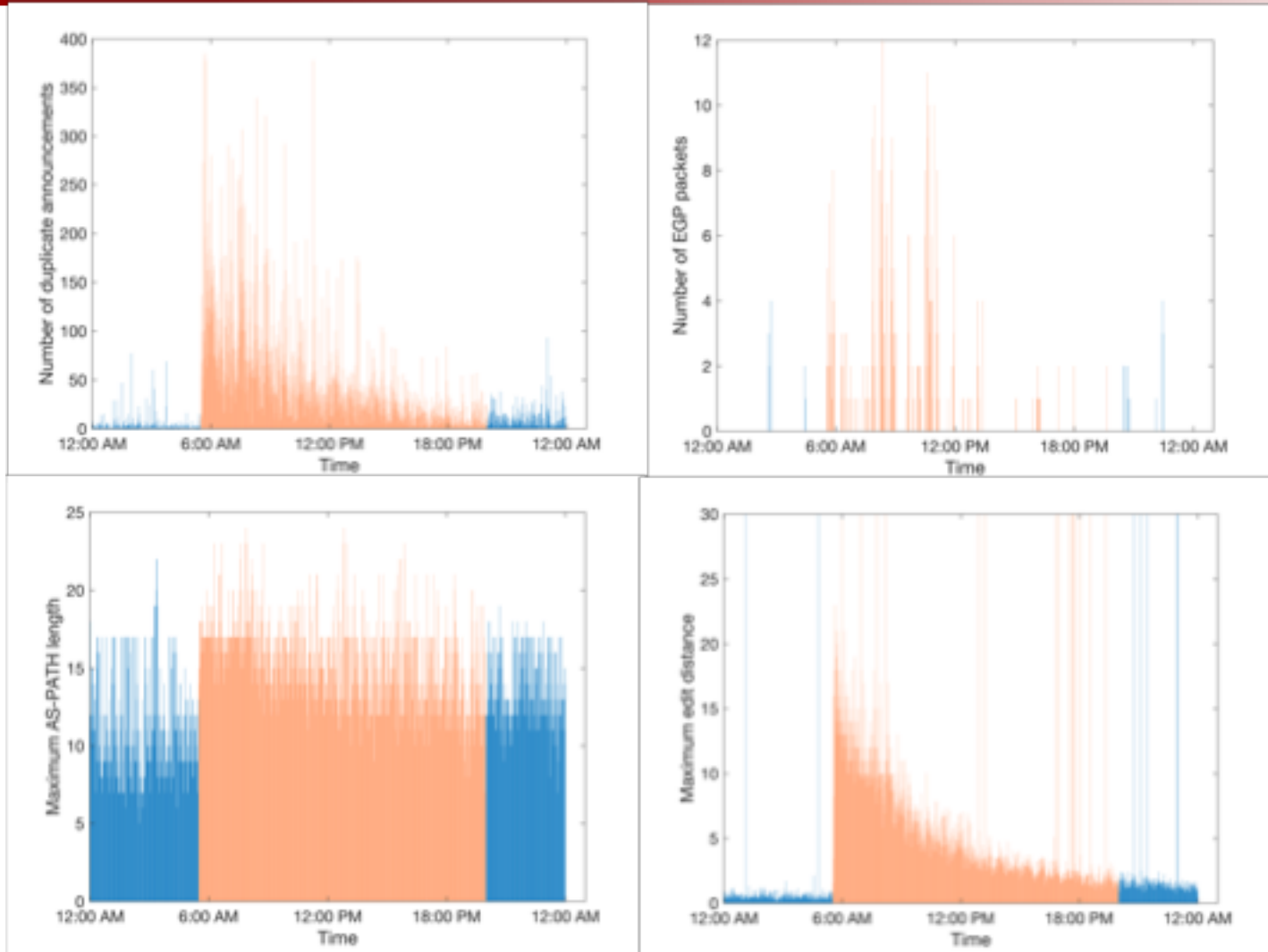
- References

# Feature extraction method

- Converted BGP update messages from MRT into American Standard Code for Information Interchange (ASCII) format

- Used LibBGPdump library on a Linux platform

- C# tool was used to extract features:
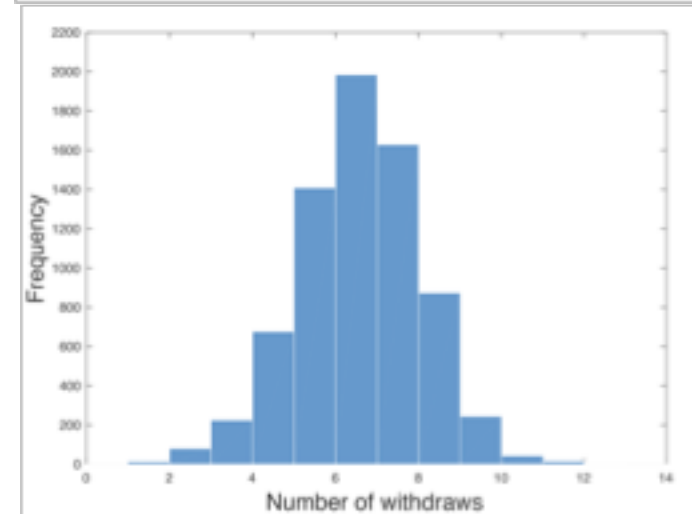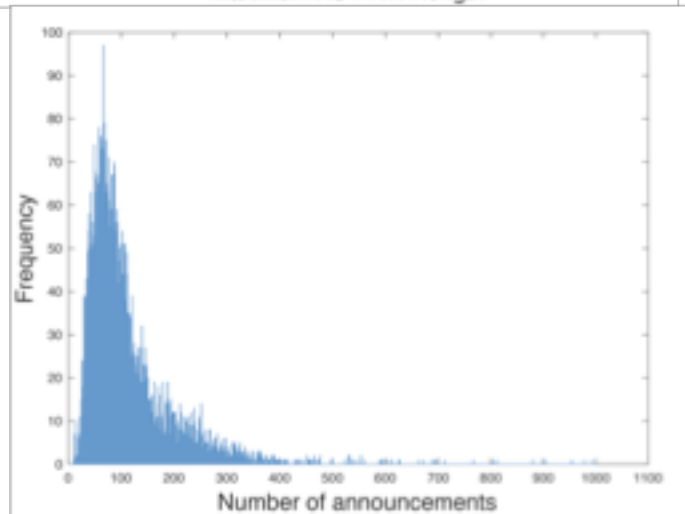  - volume
  - AS-path

# Extracted features

| Feature | Name | Category |
|---------|------|----------|
| 1 | Number of announcements | *volume* |
| 2 | Number of withdrawals | *volume* |
| 3 | Number of announced NLRI prefixes | *volume* |
| 4 | Number of withdrawn NLRI prefixes | *volume* |
| 5 | Average *AS-path* length | *AS-path* |
| 6 | Maximum *AS-path* length | *AS-path* |
| 7 | Average unique *AS-path* length | *AS-path* |
| 8 | Number of duplicate announcements | *volume* |
| 9 | Number of duplicate withdrawals | *volume* |
| 10 | Number of implicit withdrawals | *volume* |
| 11 | Average edit distance | *AS-path* |
| 12 | Maximum edit distance | *AS-path* |
| 13 | Inter-arrival time | *volume* |
| 14-24 | Maximum edit distance $= n$, where $n = (7, ..., 17)$ | *AS-path* |
| 25-33 | Maximum *AS-path* length $= n$, where $n = (7, ..., 15)$ | *AS-path* |
| 34 | Number of Interior Gateway Protocol (IGP) packets | *volume* |
| 35 | Number of Exterior Gateway Protocol (EGP) packets | *volume* |
| 36 | Number of incomplete packets | *volume* |
| 37 | Packet size $(B)$ | *volume* |

# Volume and AS-path features: Slammer worm

# Distribution of features: Slammer worm

June 19, 2018

Application of Machine Learning Techniques for Detecting Anomalies
in Communication Networks

25/58

SFU
SIMON FRASER
UNIVERSITY

# Roadmap

- Introduction
- Border Gateway Protocol datasets
- Extraction of features from BGP update messages
- Performance metrics
- Long Short-Term Memory
- Comparison of classification algorithms
- Discussion
- Future work and conclusion
- References

SFU

SIMON FRASER
UNIVERSITY

June 19, 2018

Application of Machine Learning Techniques for Detecting Anomalies
in Communication Networks

26/58

# Performance metrics

Confusion matrix:

| | Predicted class | |
|---|---|---|
| Actual class | Anomaly (positive) | Regular (negative) |
| Anomaly (positive) | TP | FN |
| Regular (negative) | FP | TN |

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$\text{F-Score} = 2 \times \frac{precision \times sensitivity}{precision + sensitivity}$$

$$\text{precision} = \frac{TP}{TP + FP}$$

$$\text{sensitivity (recall)} = \frac{TP}{TP + FN}$$

# Roadmap

- Introduction

- Border Gateway Protocol datasets

- Extraction of features from BGP update messages

- Performance metrics

- Long Short-Term Memory

- Comparison of classification algorithms

- Discussion

- Future work and conclusion

- References

Application of Machine Learning Techniques for Detecting Anomalies
in Communication Networks

SFU
SIMON FRASER
UNIVERSITY

# Long Short-term Memory: LSTM

- A special form of the recurrent neural networks (RNNs):
  - LSTM cell (memory block)

- Connects time intervals (short-term memories) to form a continuous memory

- Overcomes long-term dependency

- Prevents vanishing gradient problems

SFU
SIMON FRASER
UNIVERSITY

June 19, 2018

Application of Machine Learning Techniques for Detecting Anomalies in Communication Networks

29/58

# LSTM cell: components

- Input node $g_{nt}$:
  - contains input information

- Input gate $i_{nt}$:
  - controls the information to be updated in the LSTM cell

- Internal state $c_t$:
  - stores the cell's memory

- Forget gate $f_{nt}$:
  - determines whether to remember or discard the memories

- Output gate $o_{nt}$:
  - filters and clears irrelevant memories

$n$: The $n$ th LSTM cell

# LSTM repeating module

# LSTM components

Input node: $\quad g_{nt} = tanh(U_{gn}h_{t-1} + W_g x_t + b_{gn})$

Input gate: $\quad i_{nt} = \sigma(U_{in}h_{t-1} + W_i x_t + b_{in})$

Forget gate: $\quad f_{nt} = \sigma(U_{fn}h_{t-1} + W_f x_t + b_{fn})$

Output gate: $\quad o_{nt} = \sigma(U_{on}h_{t-1} + W_o x_t + b_{on})$

Internal state: $\quad c_t = f_{nt} * c_{t-1} + i_{nt} * tanh(U_c h_{t-1} + W_c x_t + b_c)$

Output node: $\quad h_t = o_{nt} * ReLU(c_t)$

ReLU: Rectified Linear Unit

SFU

SIMON FRASER
UNIVERSITY

# LSTM components

- $tanh$: tangent activation function

- $U_*$ and $W_*$: weight parameters

- $h_{t-1}$: hidden layer at the previous time step

- $x_t$: input at the current time step

- $b_{*n}$: bias of the $n$th LSTM cell
- $\sigma$: sigmoid activation function

# Internal state: actions

| Input gate | Forget gate | Actions |
|---|---|---|
| 0 | 1 | Keep the memory from the previous time step |
| 1 | 1 | Add the current information to the memory |
| 0 | 0 | Discard both current and past information |
| 1 | 0 | Overwrite the memory by current information |

# LSTM classification procedure

- Keras:
  - open source neural networks Application Program Interface (API) written in Python
  - enables fast experimentation with deep neural networks
  - runs on top of TensorFlow

- Import sequential model from Keras

- Normalize data points and scale their values within the range [0, 1]

- Replace anomaly labels by 0

- Length of time sequence: 20

- Keras: Deep Learning library for Theano and TensorFlow. [Online]. Available: https: //keras.io/ [Mar. 2018].
- TensorFlow. [Online]. Available: https://www.tensorflow.org [Mar. 2018].

# LSTM classification procedure

- Adam optimizer

- Learning rate: 0.001

- Random seed: 77

- Batch size: 32

- Validation dataset: 20% of the original training dataset
- Epochs: 30

SFU
SIMON FRASER
UNIVERSITY

June 19, 2018

Application of Machine Learning Techniques for Detecting Anomalies
in Communication Networks

36/58

# LSTM model: implementation

# Roadmap

- Introduction

- Border Gateway Protocol datasets

- Extraction of features from BGP update messages

- Performance metrics

- Long Short-Term Memory

- **Comparison of classification algorithms**

- Discussion

- Future work and conclusion

- References

SFU

SIMON FRASER
UNIVERSITY

# Support Vector Machine: SVM

- Supervised learning algorithm used for classification and regression tasks

- Used as a binary classifier for detecting BGP anomalies

- Two types of SVM models:
  - hard-margin SVMs require each data point to be correctly classified
  - soft-margin SVMs allow some data points to be misclassified

# Soft-Margin SVM



- Aims to find the maximum margin between both classes

- Support vectors determine the position of the decision boundary

SFU
SIMON FRASER
UNIVERSITY
June 19, 2018
Application of Machine Learning Techniques for Detecting Anomalies
in Communication Networks
40/58

# Soft-Margin SVM: kernel function



Kernel function: $k(x_n, x_m) = \Phi(x_n)^\top \Phi(x_m)$

SFU SIMON FRASER UNIVERSITY

# Naïve Bayes

- Used as supervised classifiers

- One of the most efficient machine learning classification techniques

- Assumes that features are conditionally independent for a given class

- Low complexity

- Trained effectively with smaller datasets

- Suitable for online real time detection of anomalies

SFU
SIMON FRASER
UNIVERSITY

June 19, 2018

Application of Machine Learning Techniques for Detecting Anomalies in Communication Networks

42/58

# Decision Tree

- Used in data mining to predict the class labels

- A tree is "learned" by splitting the input dataset into subsets based on appropriate features:
  - root: source dataset
  - internal (non-leaf) node: input feature
  - tree branches: prediction outcomes
  - leaf node: class or class probability distribution

- Advantages:
  - does not require feature selection
  - does not require linear datasets

- Software tool: C5.0

June 19, 2018

Application of Machine Learning Techniques for Detecting Anomalies in Communication Networks

43/58

SFU
SIMON FRASER
UNIVERSITY

# Extreme learning machine: ELM

- Feed-forward neural network with single hidden layer

- Avoids the iterative tuning of the weights used in traditional neural networks

- Suitable for applications that require fast response and real-time predictions

# ELM: architecture



Randomly generated weight and bias

Optimized weight

W, b

$\beta$

$x_1$  $x_2$  $\vdots$  $x_d$

$f(\cdot)$  $f(\cdot)$  $\vdots$  $f(\cdot)$

$y_1$  $y_2$  $\vdots$  $y_m$

Input layer    Hidden layer    Output layer

Output:    $y_m = \sum_{i=1}^{k} \beta_i f(w_i x_d + b_i)$

SFU
SIMON FRASER
UNIVERSITY

| Training model | Test datasets | | | |
|---|---|---|---|---|
| | Accuracy (%) | | | F-Score (%) |
| | Code Red I | RIPE regular | BCNET | Code Red I |
| LSTM$_u$1 | 95.22 | 65.49 | 57.30 | 83.17 |
| SVM$_u$1 | 78.65 | 69.17 | 57.22 | 39.51 |
| Naïve Bayes$_u$1 | 82.03 | 82.99 | 79.03 | 29.52 |
| Decision Tree$_u$1 | 85.36 | 89.00 | 77.22 | 47.82 |
| ELM$_u$1 | 80.92 | 75.81 | 69.03 | 36.27 |
| | Nimda | RIPE regular | BCNET | Nimda |
| LSTM$_u$2 | 53.94 | 51.53 | 50.80 | 11.81 |
| SVM$_u$2 | 55.50 | 89.89 | 82.08 | 24.29 |
| Naïve Bayes$_u$2 | 62.56 | 82.85 | 86.25 | 48.78 |
| Decision Tree$_u$2 | 58.13 | 94.19 | 81.18 | 26.16 |
| ELM$_u$2 | 54.42 | 96.15 | 91.88 | 13.72 |
| | Slammer | RIPE regular | BCNET | Slammer |
| LSTM$_u$3 | 95.87 | 56.74 | 58.55 | 84.62 |
| SVM$_u$3 | 93.04 | 73.92 | 59.24 | 75.93 |
| Naïve Bayes$_u$3 | 83.58 | 84.79 | 81.18 | 51.12 |
| Decision Tree$_u$3 | 95.89 | 89.42 | 77.78 | 84.34 |
| ELM$_u$3 | 86.96 | 78.57 | 73.47 | 55.31 |

SVM, Naïve Bayes, Decision Tree, and ELM results have been reported in:
- Z. Li, Q. Ding, S. Haeri, and Lj. Trajković, "Application of machine learning techniques to detecting anomalies in communication networks: classification algorithms," in *Cyber Threat Intelligence*, M. Conti, A. Dehghantanha, and T. Dargahi, Eds., Berlin: Springer, pp. 71-92, 2018.

SFU
SIMON FRASER
UNIVERSITY
June 19, 2018
Application of Machine Learning Techniques for Detecting Anomalies
in Communication Networks
46/58

# Performance comparison: balanced datasets

| Training model | Test datasets | | | |
|---|---|---|---|---|
| | Accuracy (%) | | | F-Score (%) |
| | Code Red I | RIPE regular | BCNET | Code Red I |
| LSTMb1 | 56.43 | 60.48 | 62.78 | 26.59 |
| | Nimda | RIPE regular | BCNET | Nimda |
| LSTMb2 | 56.32 | 44.27 | 53.58 | 65.96 |
| SVMb2 | 69.26 | 51.81 | 44.86 | 72.32 |
| | Slammer | RIPE regular | BCNET | Code Red I |
| LSTMb3 | 82.98 | 55.00 | 48.20 | 58.54 |
| SVMb3 | 87.19 | 63.31 | 51.11 | 64.76 |

SVM results have been reported in:
- Z. Li, Q. Ding, S. Haeri, and Lj. Trajković, "Application of machine learning techniques to detecting anomalies in communication networks: classification algorithms," in *Cyber Threat Intelligence*, M. Conti, A. Dehghantanha, and T. Dargahi, Eds., Berlin: Springer, to appear.

SIMON FRASER UNIVERSITY

# Performance comparison: unbalanced vs. balanced LSTM models

| Training model | Test datasets | | | |
|---|---|---|---|---|
| | Accuracy (%) | | | F-Score (%) |
| | Code Red I | RIPE regular | BCNET | Code Red I |
| LSTMu1 | 95.22 | 65.49 | 57.30 | 83.17 |
| LSTMb1 | 56.43 | 60.48 | 62.78 | 26.59 |
| | Nimda | RIPE regular | BCNET | Nimda |
| LSTMu2 | 53.94 | 51.53 | 50.80 | 11.81 |
| LSTMb2 | 56.32 | 44.27 | 53.58 | 65.96 |
| | Slammer | RIPE regular | BCNET | Code Red I |
| LSTMu3 | 95.87 | 56.74 | 58.55 | 84.62 |
| LSTMb3 | 82.98 | 55.00 | 48.20 | 58.54 |

SFU

SIMON FRASER UNIVERSITY

# Roadmap

- Introduction

- Border Gateway Protocol datasets

- Extraction of features from BGP update messages

- Performance metrics

- Long Short-Term Memory

- Comparison of classification algorithms

- Discussion

- Future work and conclusion

- References

June 19, 2018

Application of Machine Learning Techniques for Detecting Anomalies
in Communication Networks

49/58

SFU

SIMON FRASER
UNIVERSITY

# Discussion

- Sources of labeled anomalous datasets:
  - artificial datasets may not contain properties of the real-world data
  - RIPE and BCNET data were collected from deployed networks

- Selection of performance metrics:
  - Accuracy: ratio of correct predictions for the entire dataset
  - F-Score: more suitable because it emphasizes importance of the anomaly class

- Selection of appropriate machine learning approach:
  - application dependent
  - based on algorithm advantages and limitations

June 19, 2018

Application of Machine Learning Techniques for Detecting Anomalies in Communication Networks

50/58

SFU
SIMON FRASER
UNIVERSITY

# Roadmap

- Introduction

- Border Gateway Protocol datasets

- Extraction of features from BGP update messages

- Performance metrics

- Long Short-Term Memory

- Comparison of classification algorithms

- Discussion

- **Future work and conclusion**

- References

SFU

SIMON FRASER
UNIVERSITY

June 19, 2018

Application of Machine Learning Techniques for Detecting Anomalies
in Communication Networks

51/58

# Future work

- Optimize the LSTM performance:
  - use dropout technique in the input layer to learn independent representations of the dataset
- Tune hyperparameters to improve LSTM convergence:
  - number of LSTM cells
  - number of epochs
- Consider other LSTM architectures:
  - Gated Recurrent Unit (GRU): simplified LSTM
    - more efficient
    - requires smaller training dataset

Application of Machine Learning Technique to Detecting Anomalies in Communication Networks

SFU
SIMON FRASER
UNIVERSITY

# Conclusions

- Classified anomalies in BGP traffic traces using a number of classification models

- Extracted features and created unbalanced and balanced datasets

- Compared the performance of LSTM models to SVM, Naïve Bayes, Decision Tree, and ELM classifiers

- Performance of classifiers is influenced by the employed datasets

- No single classifier performs the best across all used datasets

- Machine learning is a feasible approach to successfully classify BGP anomalies

SFU

SIMON FRASER UNIVERSITY

# Roadmap

- Introduction
- Border Gateway Protocol datasets
- Extraction of features from BGP update messages
- Performance metrics
- Long Short-Term Memory
- Comparison of classification algorithms
- Discussion
- Future work and conclusion
- **References**

SFU

SIMON FRASER
UNIVERSITY

# References: BGP

- Y. Rekhter and T. Li, "A Border Gateway Protocol 4 (BGP-4)," RFC 1771, *IETF*, Mar. 1995. [Online]. Available: http://tools.ietf.org/rfc/rfc1771.txt [Mar. 2018].

- Y. Rekhter and T. Li, "A Border Gateway Protocol 4 (BGP-4)," RFC 4271, *IETF*, Jan. 2016. [Online]. Available: http://tools.ietf.org/rfc/rfc5271.txt [Mar. 2018].

- RIPE NCC: RIPE Network Coordination Center. [Online]. Available: http://www.ripe.net/data-tools/stats/ris/ris-raw-data [Mar. 2018].

- BCNET. [Online]. Available: http://www.bc.net [Mar. 2018].

- Bgpdump [Online]. Available: https://bitbucket.org/ripencc/bgpdump/wiki/Home [Mar. 2018].

# References: Machine learning algorithms

- C. M. Bishop, *Pattern Recognition and Machine Learning*. Secaucus, NJ, USA: Springer-Verlag, 2006, pp. 325–358.

- G. E. Hinton, S. Osindero, and Y-W. Teh, "A fast learning algorithm for deep belief nets," *Neural Compt.*, vol. 18, no. 7, pp. 1527–1554, July 2006.

- S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.,* vol. 9, no. 8, pp. 1735–1780, Oct. 1997.

- F. A. Gers, J. Schimidhuber, and F. Cummins, "Learning to forget: Continual prediction with LSTM," *Neural Computation*, vol, 12, no. 10, pp. 2451–2471, Oct. 2000.

- L. Rokach and O. Maimon, "Top-down induction of decision trees classifiers–a survey," *IEEE Trans. Syst., Man, Cybern., Appl. and Rev.*, vol. 35, no. 4, pp. 476–487, Nov. 2005.

- G. B. Huang, Q. Y. Zhu, and C. K. Siew, "Extreme learning machine: theory and applications," *Neurocomputing*, vol. 70, pp. 489–501, Dec. 2006.

SFU
SIMON FRASER
UNIVERSITY

June 19, 2018

Application of Machine Learning Technique to Detecting Anomalies in Communication Networks

56/58

# Publications:

- Q. Ding, Z. Li, S. Haeri, and Lj. Trajković, "Application of machine learning techniques to detecting anomalies in communication networks: datasets and feature selection algorithms," in *Cyber Threat Intelligence*, M. Conti, A. Dehghantanha, and T. Dargahi, Eds., Berlin: Springer, pp. 47-70, 2018.

- Z. Li, Q. Ding, S. Haeri, and Lj. Trajković, "Application of machine learning techniques to detecting anomalies in communication networks: classification algorithms," in *Cyber Threat Intelligence*, M. Conti, A. Dehghantanha, and T. Dargahi, Eds., Berlin: Springer, pp. 71-92, 2018.

- Q. Ding, Z. Li, P. Batta, and Lj. Trajković, "Detecting BGP anomalies using machine learning techniques," in *Proc. IEEE Trans. Syst., Man, Cybern.*, Budapest, Hungary, Oct. 2016, pp. 3352–3355.

- P. Batta, M. Singh, Z. Li, Q. Ding, and Lj. Trajković, "Evaluation of support vector machine kernels for detecting network anomalies," *IEEE Int. Symp. Circuits and Systems*, Florence, Italy, May 2018, pp. 1–4.

- H. Ben Yedder, Q. Ding, U. Zakia, Z. Li, S. Haeri, and Lj. Trajković, "Comparison of virtualization algorithms and topologies for data center networks," *The 26th Int. Conf. Comput. Commun., Netw., 2nd Workshop on Netw. Security Anal. Automat.*, Vancouver, Canada, Aug. 2017.

- S. Haeri, Q. Ding, Z. Li, and Lj. Trajković, "Global resource capacity algorithm with path splitting for virtual network embedding," in *Proc. IEEE Int. Symp. Circuits and Systems*, Montreal, Canada, May 2016, pp. 666–669.

# Acknowledgements

- Chair:
  - Prof. Ivan V. Bajić

- Senior supervisor:
  - Prof. Ljiljana Trajković

- Supervisor:
  - Prof. Parvaneh Saeedi

- SFU examiner:
  - Prof. Qianping Gu

June 19, 2018

Application of Machine Learning Technique to Detecting Anomalies in Communication Networks

58/58

SFU

SIMON FRASER
UNIVERSITY

# Thank you!

## Questions?

SFU
SIMON FRASER
UNIVERSITY