

San Jose State University

SJSU ScholarWorks

Mineta Transportation Institute Publications

1-2022

Taming the Data in the Internet of Vehicles

Shahab Tayeb

California State University, Fresno

Follow this and additional works at: https://scholarworks.sjsu.edu/mti_publications



Part of the [Graphics and Human Computer Interfaces Commons](#), [OS and Networks Commons](#), [Programming Languages and Compilers Commons](#), [Software Engineering Commons](#), and the [Transportation Commons](#)

Recommended Citation

Shahab Tayeb. "Taming the Data in the Internet of Vehicles" *Mineta Transportation Institute Publications* (2022). <https://doi.org/10.31979/mti.2022.2014>

This Report is brought to you for free and open access by SJSU ScholarWorks. It has been accepted for inclusion in Mineta Transportation Institute Publications by an authorized administrator of SJSU ScholarWorks. For more information, please contact scholarworks@sjsu.edu.

Taming the Data in the Internet of Vehicles

Shahab Tayeb



Mineta Transportation Institute

Founded in 1991, the Mineta Transportation Institute (MTI), an organized research and training unit in partnership with the Lucas College and Graduate School of Business at San José State University (SJSU), increases mobility for all by improving the safety, efficiency, accessibility, and convenience of our nation's transportation system. Through research, education, workforce development, and technology transfer, we help create a connected world. MTI leads the [Mineta Consortium for Transportation Mobility](#) (MCTM) funded by the U.S. Department of Transportation and the [California State University Transportation Consortium](#) (CSUTC) funded by the State of California through Senate Bill 1. MTI focuses on three primary responsibilities:

Research

MTI conducts multi-disciplinary research focused on surface transportation that contributes to effective decision making. Research areas include: active transportation; planning and policy; security and counterterrorism; sustainable transportation and land use; transit and passenger rail; transportation engineering; transportation finance; transportation technology; and workforce and labor. MTI research publications undergo expert peer review to ensure the quality of the research.

Education and Workforce

To ensure the efficient movement of people and products, we must prepare a new cohort of transportation professionals who are ready to lead a more diverse, inclusive, and equitable transportation industry. To help achieve this, MTI sponsors a suite of workforce development and education opportunities. The Institute supports educational programs offered by the Lucas Graduate School of Business: a

Master of Science in Transportation Management, plus graduate certificates that include High-Speed and Intercity Rail Management and Transportation Security Management. These flexible programs offer live online classes so that working transportation professionals can pursue an advanced degree regardless of their location.

Information and Technology Transfer

MTI utilizes a diverse array of dissemination methods and media to ensure research results reach those responsible for managing change. These methods include publication, seminars, workshops, websites, social media, webinars, and other technology transfer mechanisms. Additionally, MTI promotes the availability of completed research to professional organizations and works to integrate the research findings into the graduate education program. MTI's extensive collection of transportation-related publications is integrated into San José State University's world-class Martin Luther King, Jr. Library.

Disclaimer

The contents of this report reflect the views of the authors, who are responsible for the facts and accuracy of the information presented herein. This document is disseminated in the interest of information exchange. MTI's research is funded, partially or entirely, by grants from the California Department of Transportation, the California State University Office of the Chancellor, the U.S. Department of Homeland Security, and the U.S. Department of Transportation, who assume no liability for the contents or use thereof. This report does not constitute a standard specification, design standard, or regulation.

Report 22-02

Taming the Data in the Internet of Vehicles

Shahab Tayeb

January 2022

A publication of the
Mineta Transportation Institute
Created by Congress in 1991

College of Business
San José State University
San José, CA 95192-0219

TECHNICAL REPORT DOCUMENTATION PAGE

1. Report No. 22-02	2. Government Accession No.	3. Recipient's Catalog No.	
4. Title and Subtitle Taming the Data in the Internet of Vehicles		5. Report Date January 2022	
		6. Performing Organization Code	
7. Authors Shahab Tayeb, PhD: 0000-0002-7466-1042		8. Performing Organization Report CA-MTI-2014	
9. Performing Organization Name and Address Mineta Transportation Institute College of Business San José State University San José, CA 95192-0219		10. Work Unit No.	
		11. Contract or Grant No. ZSB12017-SJAUX	
12. Sponsoring Agency Name and Address State of California SB1 2017/2018 Trustees of the California State University Sponsored Programs Administration 401 Golden Shore, 5 th Floor Long Beach, CA 90802		13. Type of Report and Period Covered	
		14. Sponsoring Agency Code	
15. Supplemental Notes			
16. Abstract As an emerging field, the Internet of Vehicles (IoV) has a myriad of security vulnerabilities that must be addressed to protect system integrity. To stay ahead of novel attacks, cybersecurity professionals are developing new software and systems using machine learning techniques. Neural network architectures improve such systems, including Intrusion Detection System (IDSs), by implementing anomaly detection, which differentiates benign data packets from malicious ones. For an IDS to best predict anomalies, the model is trained on data that is typically pre-processed through normalization and feature selection/reduction. These pre-processing techniques play an important role in training a neural network to optimize its performance. This research studies the impact of applying normalization techniques as a pre-processing step to learning, as used by the IDSs. The impacts of pre-processing techniques play an important role in training neural networks to optimize its performance. This report proposes a Deep Neural Network (DNN) model with two hidden layers for IDS architecture and compares two commonly used normalization pre-processing techniques. Our findings are evaluated using accuracy, Area Under Curve (AUC), Receiver Operator Characteristic (ROC), F-1 Score, and loss. The experimentations demonstrate that Z-Score outperforms no-normalization and the use of Min-Max normalization.			
17. Key Words Data cleaning, Data models, Data sharing, Machine learning, Neural networks		18. Distribution Statement No restrictions. This document is available to the public through The National Technical Information Service, Springfield, VA 22161.	
19. Security Classif. (of this report) Unclassified	20. Security Classif. (of this page) Unclassified	21. No. of Pages 25	22. Price

Copyright © 2022
by **Mineta Transportation Institute**
All rights reserved.

DOI: 10.31979/mti.2022.2014

Mineta Transportation Institute
College of Business
San José State University
San José, CA 95192-0219

Tel: (408) 924-7560
Fax: (408) 924-7565
Email: mineta-institute@sjsu.edu

transweb.sjsu.edu/research/2014

ACKNOWLEDGMENTS

The author acknowledges the contributions of faculty collaborators Drs. Youngwook Kim, Matin Pirouz, Reza Raeisi, Aaron Stillmaker, and Aly Tawfik.

CONTENTS

Acknowledgments	vi
List of Figures	viii
Executive Summary	1
1. Introduction	2
2. Related Work.....	4
3. Methods	5
3.1 Pruning and Normalization Techniques	5
3.2 Neural Network Architecture Selection and Training.....	5
3.3 Metrics and Testbed.....	6
4. Results and Performanc	7
4.1 Accuracy.....	7
4.2 Cross-Entropy Loss.....	8
4.3 F-Score	9
4.4 Area Under the ROC Curve.....	9
4.5 Limitations of the Specific DNN Architecture.....	10
5. Summary & Conclusions.....	11
Endnotes	12
Bibliography	15
About the Author.....	18

LIST OF FIGURES

Figure 1. Vulnerability Surfaces of the Emerging Internet of Vehicles	3
Figure 2. Comparison of Accuracy in Training and Validation Datasets using Different Normalization Techniques.....	7
Figure 3. Comparison of Loss in Training and Validation Datasets using Different Normalization Techniques.....	8
Figure 4. Comparison of F-score in Training and Validation Datasets using Different Normalization Techniques.....	9
Figure 5. Comparison of AUC-ROC in Training and Validation Datasets using Different Normalization Techniques.....	10

Executive Summary

To stay ahead of novel attacks, cybersecurity professionals are developing new software programs and systems using machine learning techniques. Neural network architectures improve such systems, including Intrusion Detection System (IDSs), by implementing anomaly detection, which differentiates benign packets from malicious packets. For an IDS to best predict anomalies, the model's training dataset is typically pre-processed through normalization and feature selection/reduction. Pre-processing techniques play an important role in training a neural network to optimize its performance.

In this study, we extend the current research on the importance of data normalization through developing, training, and testing a Deep Neural Network on CIDDS network data. To this end, we evaluate the effect of Z-Score and Min-Max normalization on the model's accuracy, loss, F-Score, and AUC-ROC. Additionally, an analysis and comparison of the performance of the model on the NSL-KDD and CIDDS datasets are carried out.

1. Introduction

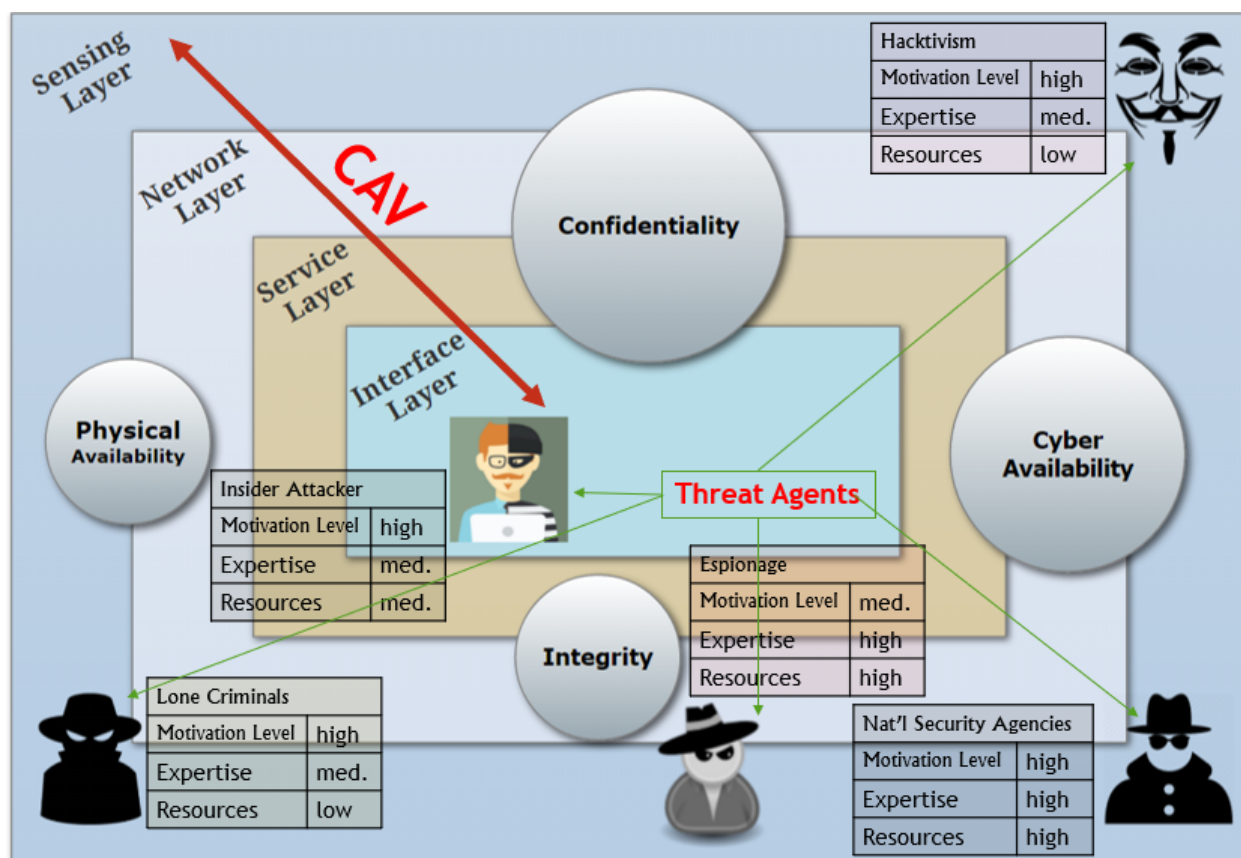
Anomaly-based Intrusion Detection System (IDSs) are at the forefront of research due to their ability to distinguish malignant from benign messages using machine learning and neural networks.¹ This provides a security advantage: specifically, the real-time identification of threats to prevent a system from becoming infected.² As neural networks have emerged as a focal area of research, the study of their application to security and IDSs has become more widespread.^{3 4} These are particularly important in the context of the Internet of Vehicles (IoV). Figure 1 provides an overview of the security vulnerabilities as they relate to IoV. In Figure 1, the potential intruders, also known as threat agents, are mapped onto the layered architecture of the IoV backbone. In parallel, the four major pillars of IoV security—namely, confidentiality, integrity, physical availability, and cyber availability—are plotted on the various layers.

Since IDSs seek to prevent threats from interrupting network communications, they have been proposed to solve security problems for larger network applications, such as the Internet of Things (IoT). These applications are characterized by having large volumes of continuous data streams that require rapid processing.⁵ Further, these data may be confidential—a consideration which mandates the security of network applications. The goal of taming such data presents a myriad of challenges, which can potentially be mitigated using self-adapting learning models. Deep-learning-based IDSs offer themselves as a solution to securing large-scale, network-based systems by predicting existing and novel attacks.⁶

In order for deep-learning-based IDS to effectively secure these systems, the process in which the model classifies different messages to be anomalies (or attacks) is important.⁷ Neural networks use binary or multi-class classification to determine if a data packet is malicious or benign. An IDS uses a neural network model to identify patterns and deeper characteristics in malicious and benign data to train in the classification of the different types of data.^{2 8} This pattern recognition allows IDSs to predict novel malicious data, which otherwise would need to be specifically identified to prevent it from infecting the system.^{9 10} Due to the use of patterns to classify data, neural networks require the pre-processing of data to properly classify intrusion.¹¹ Pre-processing of the training model also ensures the optimal usage of in-vehicle resources.

The scales in values may vary across attributes in a dataset, which can skew the influence of some features over others in packet prediction. For example, an attribute with binary values may not impact the prediction result as much as an attribute with values on the scale of hundreds would. Normalization modifies the attributes of data to a similar scale to enhance prediction accuracy.¹¹ Non-numeric attributes must also be reassigned discrete values during this data pre-processing stage. However, little research has been done to demonstrate the importance of normalization to the training and usage of neural networks and IDSs.¹³ This motivates the analysis of the effects of data normalization within pre-processing for anomaly-based IDSs, which serves as the problem statement for this research.

Figure 1. Vulnerability Surfaces of the Emerging Internet of Vehicles



2. Related Work

An Artificial Neural Network (ANN) attempts to replicate the human brain and its learning process using computational units called neurons.¹⁴ ANN architecture consists of an input layer, output layer, and at least one hidden layer. Taher et al.¹⁵ implement two ANN models, each with a different feature selection method. The model with filter selection used 35 features and produced an 83.68% accuracy; the model with wrapper selection used 17 features and produced a higher accuracy of 94.02%.

Deep Neural Networks (DNNs) form a subset of ANNs that can recognize complex patterns because of their ability to process higher-level features (e.g., images) with multiple layers of abstraction.^{16 17} Vigneswaran et al.'s proposed DNN with 3 hidden layers and 41 input neurons outperformed DNNs with other numbers of hidden layers with an accuracy of about 93% and precision of about 99.7%.¹⁸

A Convolutional Neural Network (CNN) is a variation of the ANN-based Multi-Layer Perceptron.¹⁹ In contrast to other neural networks, a CNN has preliminary convolutional layers to process raw data in the form of images and extract intermediate features. Blanco et al. proposed a CNN model whose convolutional filter had a dimension of 3x3 and a depth of 4 coupled with a genetic algorithm.²⁰ This design was chosen to address the layout of features in the raw data and produced a 94.47% accuracy.

Recurrent Neural Networks (RNNs) reuse information using cyclic connections that intertwine current input with previous hidden states.^{21 22 23} Yin et al. proposed a multi-class classification RNN with an accuracy of 81.29%.²⁴ Subsequent research utilized Long Short-Term Memory (LSTM) to improve upon RNNs by replacing neurons with memory cells to improve long-term dependency.^{25 26} Kim et al. proposed a LSTM-RNN IDS with an accuracy of 96.93% and a detection rate of 98.88%.²⁷

3. Methods

3.1 Pruning and Normalization Techniques

Results were collected using two validation datasets: KDD_Test+ and a pruned KDD_Test+. KDD_Test+ is an unaltered validation dataset provided by the NSL-KDD dataset that includes attacks not included in KDD_Train+ along with a slightly different distribution of attacks as a percentage of dataset entries. Thus, the pruned KDD_Test+ set eliminates attacks not included in KDD_Train+. Testing against both validation sets allowed for an analysis of normalization techniques against known attacks and novel attacks.

Normalization techniques (Z-Score and Min-Max) were applied to the training and validation datasets prior to passing them to the neural network model. Key values such as mean, standard deviation, minimum, and maximum were taken per feature column from the training dataset and applied to the training and validation datasets. “None,” as it appears in Figure 2, indicates that no normalization techniques were applied.

3.2 Neural Network Architecture Selection and Training

We opt to use a DNN as the base architecture for our approach because of its nature as a simple model forming complex relationships. DNNs are subsets of ANNs that have multiple hidden layers and added complexity in connections they may form, making ANNs seem simplistic and less desirable in comparison.^{28 29} Further, there is a research gap in implementations of DNNs for IDSs, which motivates the study of DNNs’ potential in the field.

CNNs specialize in artificial vision and image processing, deviating from the goal of IDSs and data encountered.^{30 31} This mismatch in compatibility deters us from selecting a CNN. RNNs and LSTMs are cyclically connected and are more complex than other neural networks. Their accuracy for IDS implementations, however, does not greatly improve compared to the others.^{32 33}

The DNN model tested consisted of: an input layer of 27 features; two dense, hidden layers of 128 neurons each; and a 1-neuron output layer. The two hidden layers used a ReLu activation function and the output layer used a sigmoid activation function.

The Adam optimizer was used to adjust weights with a learning rate of 0.0001, which is one-tenth of the default learning rate in TensorFlow. Reducing the default learning rate allowed for greater detail when plotting evaluated metrics per epoch, though it came at the cost of computation time.

3.3 Metrics and Testbed

Pre-processing techniques were evaluated based on accuracy, binary cross-entropy loss, AUC-ROC, True Positives (TPs), True Negatives (TNs), False Positives (FPs), and False Negatives (FNs) for each epoch on the training dataset (KDD_Train+) and the validation datasets (KDD_Test+ and KDD_Test+ with attacks not present in KDD_Train+ removed) using functions built into Tensorflow. A callback function was used to compare both validation sets with the same trained model.

Experimentation was conducted on a PC running Windows 10 and Python 3.7.4 with an AMD Ryzen 3900X CPU and an Nvidia GTX 1080 GPU. Pruning and normalization techniques were applied using Pandas v1.0.3. The dataset was imported and the NN model built, trained, and validated using TensorFlow-GPU v2.1.0.

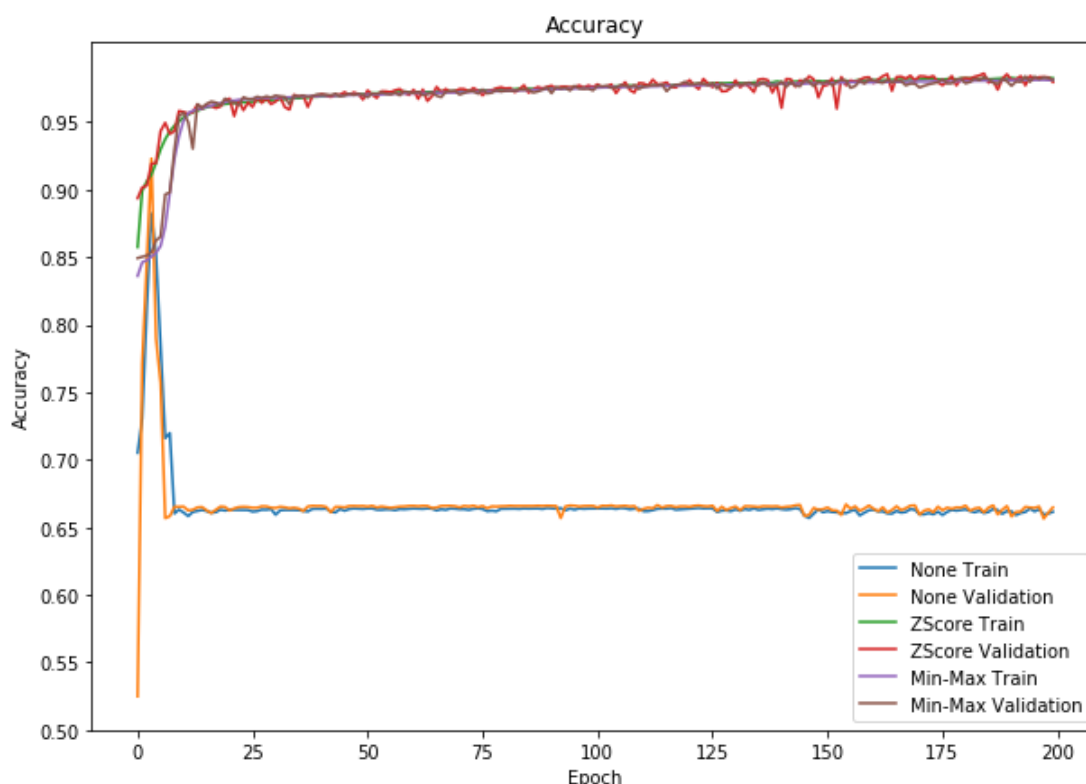
4. Results and Performance

4.1 Accuracy

This section compares the performance of the three approaches to normalization according to metric. For every metric, there is a graph analyzing the performance over 200 epochs of each normalization technique on the training and validation datasets. Results are recorded for the final performance metrics when the model is fully trained and we note its performance on the validation dataset. These metrics are then compared to results yielded from a similar experiment conducted on the NSL-KDD dataset.

For accuracy, the general trend in Figure 1 is a sharp increase in accuracy and then a stabilization of the curve. The trendline for no normalization, however, dips at the 75th epoch, indicating an overfitting of the training dataset. When a plateau is reached, the implication for implementation is significant additional overhead with little to no increase in accuracy, which is redundant in a lightweight application such as IoV security.

Figure 2. Comparison of Accuracy in Training and Validation Datasets using Different Normalization Techniques



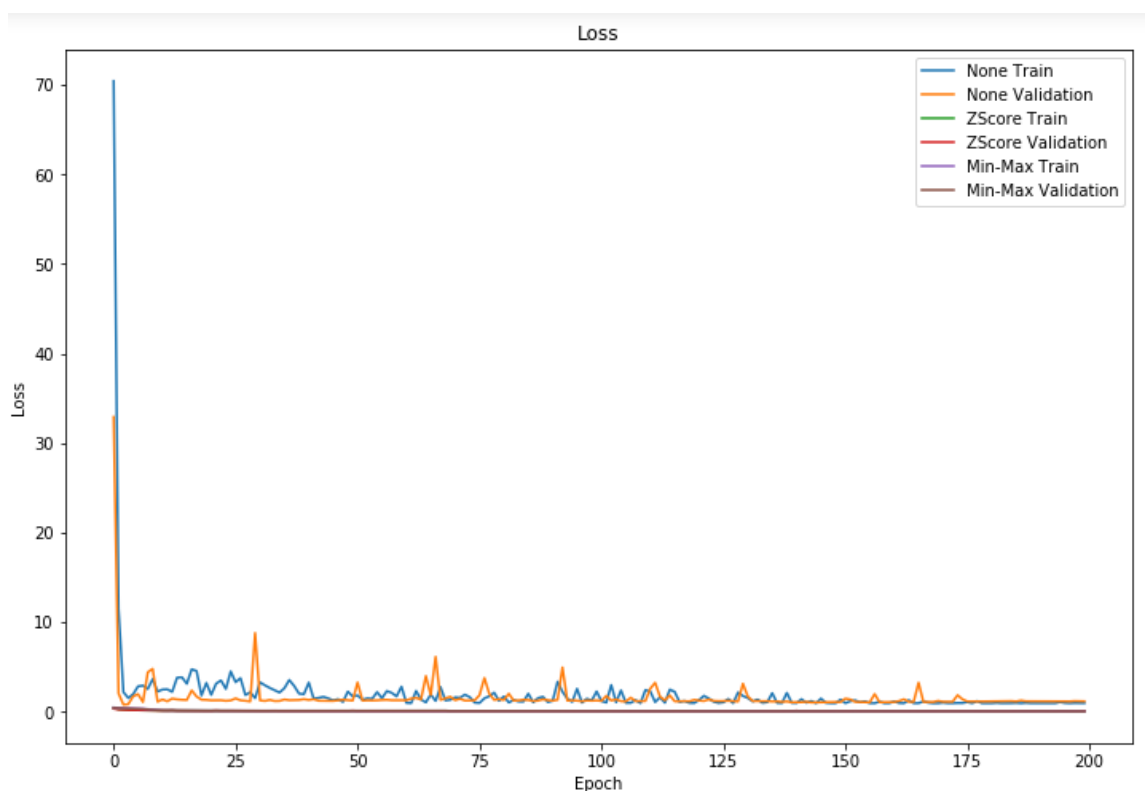
The training of the model resulted in an accuracy of 0.6626 for no normalization, 0.9824 for Z-Score, and 0.9816 for Min-Max. When tested on the validation dataset, the model produced similar results; no normalization and Z-Score achieved a minor increase in accuracy, while Min-Max yielded a slight decrease. The validation accuracy was 0.6681 for no normalization, 0.9833

for Z-Score, and 0.9779 for Min-Max. The percentage difference between the training accuracy and validation accuracy for each metric was below 1%. The model's performance was approximately 38% less effective when no normalization was used compared to when the data were transformed with either Z-Score or Min-Max. The difference in performance between Z-Score and Min-Max was 0.55%, a value too minimal to conclude which method provides better accuracy.

4.2 Cross-Entropy Loss

For cross-entropy loss, the overall pattern for Z-Score and Min-Max (see Figure 2) shows a sharp decrease in loss and then a stabilization of the curve beginning by or before the 25th epoch. The graph for no normalization differs from these two: loss remains stable at approximately 0.8 before a spike after the 75th epoch, where the curve stabilizes.

Figure 3. Comparison of Loss in Training and Validation Datasets using Different Normalization Techniques

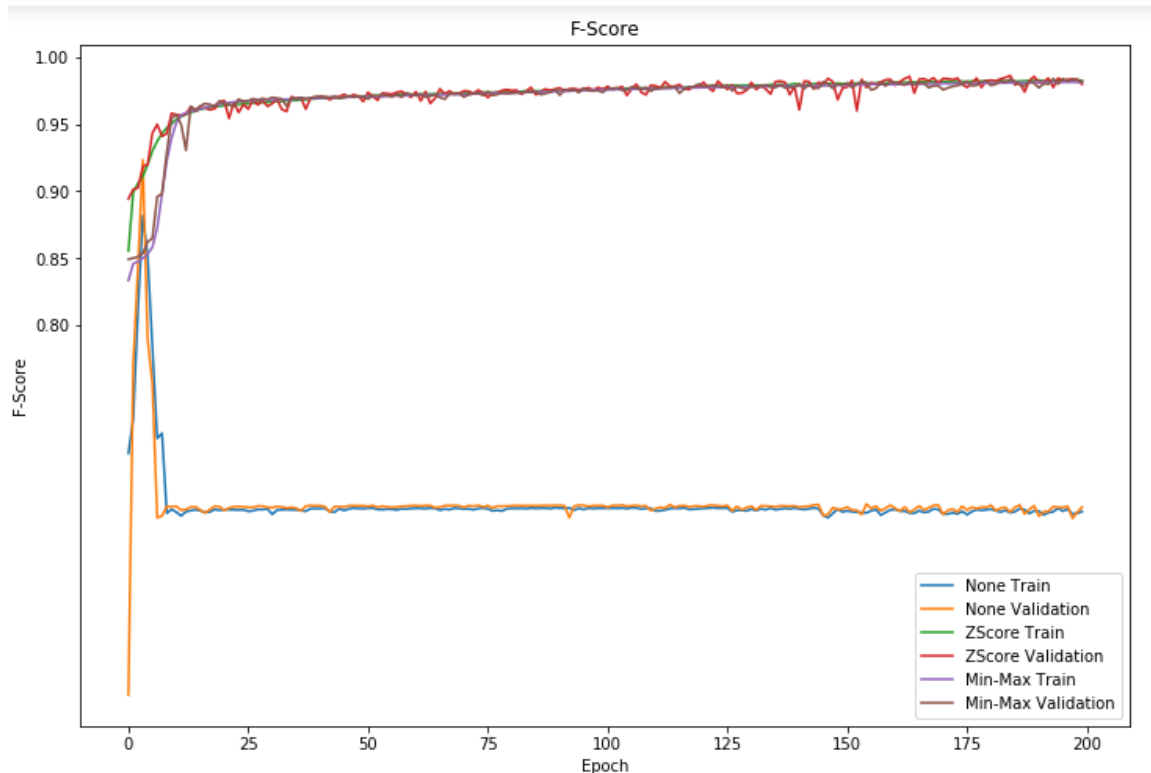


The training of the model produced a final loss value of 0.966 for no normalization, 0.0487 for Z-Score, and 0.0483 for Min-Max. On the validation dataset, no normalization and Z-Score decreased in loss while Min-Max showed an increase. The validation loss was 0.9411 for no normalization, 0.0462 for Z-Score, and 0.0577 for Min-Max. The percentage difference between the loss in training and in validation was 2.6% for no normalization, 5.27% for Z-Score, and 17.7% for Min-Max. The model performed significantly worse when no normalization was used compared to when either Z-Score or Min-Max was used in pre-processing the data. Z-Score outperformed Min-Max by 22.1%, demonstrating an advantage to using Z-Score.

4.3 F-Score

The general trend in Figure 3 illustrates a sharp increase in F-Score and then a stabilization of the curve by or before the 25th epoch for both Min-Max and Z-Score. The trendline for no normalization deviates from these results: there is a spike in F-Score followed by a levelling of the curve, and then a sharp dip at the 75th epoch.

Figure 4. Comparison of F-score in Training and Validation Datasets using Different Normalization Techniques



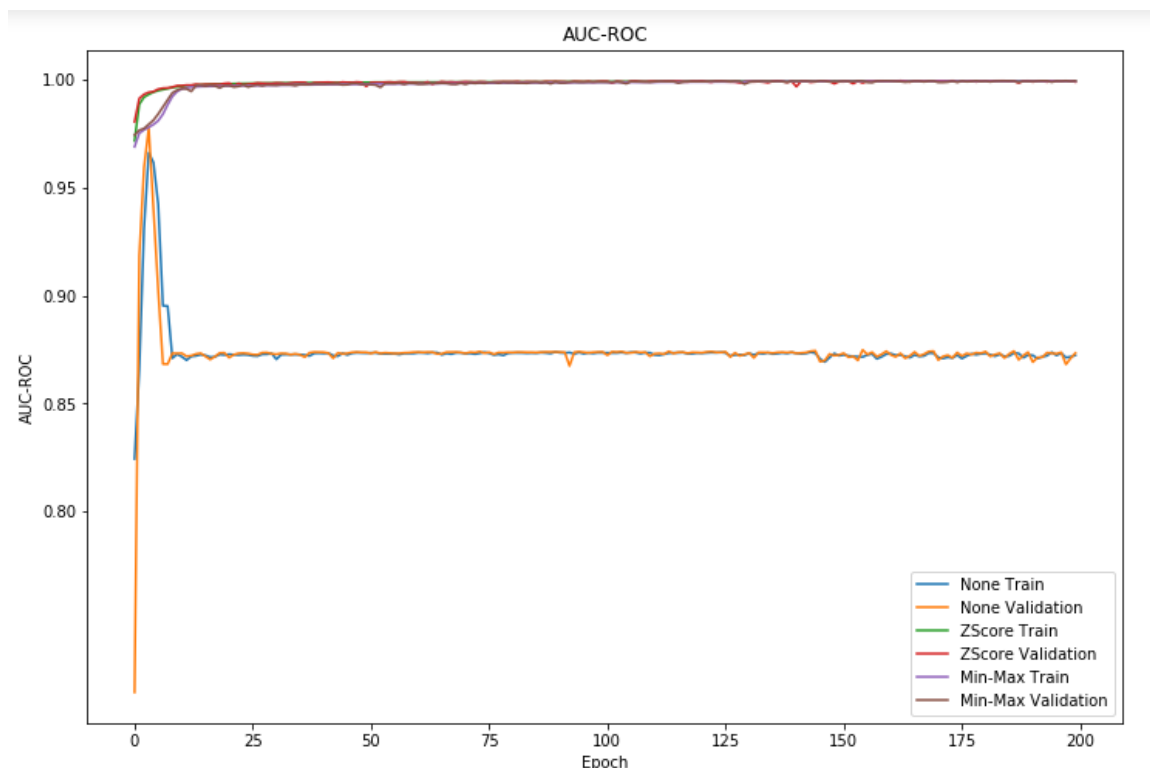
When tested on the training dataset, the model resulted in an F-Score of 0.6615 for no normalization, 0.9824 for Z-Score, and 0.9814 for Min-Max. For the validation dataset, the model behaved similarly. No normalization saw a slight increase and Z-Score showed a minor decrease, while Min-Max remained identical. The validation F-Score was 0.6646 for no normalization, 0.9796 for Z-Score, and 0.9814 for Min-Max. For each metric, the percentage difference between the training F-Score and validation F-Score stayed below 1%. The model was 38.4% less effective when using no normalization on the datasets instead of Min-Max or Z-Score. The percentage difference of the performance of Z-Score and Min-Max was 0.18%, a value too minute to support any conclusion regarding which method is more advantageous.

4.4 Area Under the ROC Curve

The overall trend for AUC-ROC shown in Figure 4 is a swift increase followed by a levelling out of the curve before the 25th epoch. For no normalization, the trendline follows a pattern similar

to its accuracy curve; the graph drops quickly at the 75th epoch, which can be attributed to overfitting.

Figure 5. Comparison of AUC-ROC in Training and Validation Datasets using Different Normalization Techniques



The training of the model yielded an AUC-ROC of 0.873 for no normalization and 0.9994 for both Z-Score and Min-Max. On the validation dataset, the model achieved very similar performance, with no normalization slightly increasing in AUC-ROC. The validation AUC-ROC was 0.8763 for no normalization, 0.9994 for Z-Score, and 0.9992 for Min-Max. For all metrics, the difference in training and validation AUC-ROC was less than 1%. When no normalization was performed on the dataset, the model's performance was 13.5%, less effective than if Z-Score or Min-Max had been used to pre-process the data. Z-Score and Min-Max produced almost identical results for AUC-ROC, so the results are inconclusive on which method outperformed the other.

4.5 Limitations of the Specific DNN Architecture

The studied model utilizes a lightweight layered architecture, in view of the resource-constrained nature of the IoV systems. Further studies are needed on the expansion of additional hidden layers and their added overhead in terms of memory and processing requirements.

5. Summary & Conclusions

Normalization and other pre-processing techniques applied to the data used for training an IDS are important for optimizing the performance metrics. We propose a DNN using 27 input features for binary classification trained using the NSL-KDD dataset. As expected, the experimentation on the pruned dataset outperforms the experimentation on the complete dataset across most metrics. Our proposed model determines for the complete dataset that Z-Score normalization, as well as Min-Max normalization to a lesser degree, improves the performance of the proposed IDS compared to no normalization. Our results demonstrate that Z-Score improves on no normalization by 4.46% for accuracy, 2.04% for loss, 4.70% for F-Score, and 23.64% for AUC-ROC. Min-Max normalization presents similar improvements with a 1.99% increase in accuracy, 1.00% decrease in loss, 0.32% increase in F-Score, and 23.65% increase in AUC-ROC. Implementing Z-Score normalization as a pre-processing step can improve the performance of DNN-based IDSs. Such pre-processing of the training model adds little to no overhead on the in-vehicle implementation of such a model, justifying the use of such pre-processing techniques. The study also highlights the optimal number of epochs for the training, after which little gain in accuracy is observed.

Endnotes

- ¹ Davis, Alexander, Sumanjit Gill, Robert Wong, and Shahab Tayeb. "Feature Selection for Deep Neural Networks in Cyber Security Applications." In 2020 IEEE International IOT, Electronics and Mechatronics Conference (IEMTRONICS), pp. 1-7. IEEE, 2020.
- ² Tidjon, Lionel N., Marc Frappier, and Amel Mammar. "Intrusion detection systems: A cross-domain overview." IEEE Communications Surveys & Tutorials 21, no. 4 (2019): 3639-3681.
- ³ Shone, Nathan, Tran Nguyen Ngoc, Vu Dinh Phai, and Qi Shi. "A deep learning approach to network intrusion detection." IEEE transactions on emerging topics in computational intelligence 2, no. 1 (2018): 41-50.
- ⁴ Nie, Laisen, Zhaolong Ning, Xiaojie Wang, Xiping Hu, Yongkang Li, and Jun Cheng. "Data-Driven Intrusion Detection for Intelligent Internet of Vehicles: A Deep Convolutional Neural Network-based Method." IEEE Transactions on Network Science and Engineering (2020).
- ⁵ Eskandari, Mojtaba, Zaffar Haider Janjua, Massimo Vecchio, and Fabio Antonelli. "Passban IDS: An intelligent anomaly based intrusion detection system for IoT edge devices." IEEE Internet of Things Journal (2020).
- ⁶ Zhong, Wei, Ning Yu, and Chunyu Ai. "Applying big data based deep learning system to intrusion detection." Big Data Mining and Analytics 3, no. 3 (2020): 181-195.
- ⁷ Hakim, Lukman, and Rahilla Fatma. "Influence Analysis of Feature Selection to Network Intrusion Detection System Performance Using NSL-KDD Dataset." In 2019 International Conference on Computer Science, Information Technology, and Electrical Engineering (ICOMITEE), pp. 217-220. IEEE, 2019.
- ⁸ Ma, Chencheng, Xuehui Du, and Lifeng Cao. "Analysis of Multi-Types of Flow Features Based on Hybrid Neural Network for Improving Network Anomaly Detection." IEEE Access 7 (2019): 148363-148380.
- ⁹ Ambusaidi, Mohammed A., Xiangjian He, Priyadarsi Nanda, and Zhiyuan Tan. "Building an intrusion detection system using a filter-based feature selection algorithm." IEEE transactions on computers 65, no. 10 (2016): 2986-2998.
- ¹⁰ Nisioti, Antonia, Alexios Mylonas, Paul D. Yoo, and Vasilios Katos. "From intrusion detection to attacker attribution: A comprehensive survey of unsupervised methods." IEEE Communications Surveys & Tutorials 20, no. 4 (2018): 3369-3388.
- ¹¹ Chiba, Zouhair, Noredine Abghour, Khalid Moussaid, and Mohamed Rida. "Intelligent approach to build a Deep Neural Network based IDS for cloud environment using combination of machine learning algorithms." Computers & Security 86 (2019): 291-317.
- ¹² Obaid, Hadeel S., Saad Ahmed Dheyab, and Sana Sabah Sabry. "The Impact of Data Pre-Processing Techniques and Dimensionality Reduction on the Accuracy of Machine Learning." In 2019 9th Annual Information Technology, Electromechanical Engineering and Microelectronics Conference (IEMECON), pp. 279-283. IEEE, 2019.
- ¹³ Obaid, Hadeel S., Saad Ahmed Dheyab, and Sana Sabah Sabry. "The Impact of Data Pre-Processing Techniques and Dimensionality Reduction on the Accuracy of Machine Learning." In 2019 9th Annual Information Technology, Electromechanical Engineering and Microelectronics Conference (IEMECON), pp. 279-283. IEEE, 2019.
- ¹⁴ Larriva-Novo, Xavier A., Mario Vega-Barbas, Víctor A. Villagrà, and Mario Sanz Rodrigo. "Evaluation of Cybersecurity Data Set Characteristics for Their Applicability to Neural Networks Algorithms Detecting Cybersecurity Anomalies." IEEE Access 8 (2020): 9005-9014.

- ¹⁵ Taher, Kazi Abu, Billal Mohammed Yasin Jisan, and Md Mahbubur Rahman. "Network intrusion detection using supervised machine learning technique with feature selection." In 2019 International Conference on Robotics, Electrical and Signal Processing Techniques (ICREST), pp. 643-646. IEEE, 2019.
- ¹⁶ Amarasinghe, Kasun, and Milos Manic. "Improving user trust on deep neural networks based intrusion detection systems." In IECON 2018-44th Annual Conference of the IEEE Industrial Electronics Society, pp. 3262-3268. IEEE, 2018.
- ¹⁷ Xiao, Huaxin, Jiashi Feng, Yunchao Wei, Maojun Zhang, and Shuicheng Yan. "Deep salient object detection with dense connections and distraction diagnosis." *IEEE Transactions on Multimedia* 20, no. 12 (2018): 3239-3251.
- ¹⁸ Vigneswaran, K. Rahul, R. Vinayakumar, K. P. Soman, and Prabakaran Poornachandran. "Evaluating shallow and deep neural networks for network intrusion detection systems in cyber security." In 2018 9th International Conference on Computing, Communication and Networking Technologies (ICCCNT), pp. 1-6. IEEE, 2018.
- ¹⁹ Sarvari, Samira, Nor Fazlida Mohd Sani, Zurina Mohd Hanapi, and Mohd Taufik Abdullah. "An Efficient Anomaly Intrusion Detection Method with Feature Selection and Evolutionary Neural Network." *IEEE Access* 8 (2020): 70651-70663.
- ²⁰ Blanco, Roberto, Pedro Malagón, Juan J. Cilla, and José M. Moya. "Multiclass network attack classifier using CNN tuned with genetic algorithms." In 2018 28th International Symposium on Power and Timing Modeling, Optimization and Simulation (PATMOS), pp. 177-182. IEEE, 2018.
- ²¹ Kim, Jihyun, Jaehyun Kim, Huong Le Thi Thu, and Howon Kim. "Long short term memory recurrent neural network classifier for intrusion detection." In 2016 International Conference on Platform Technology and Service (PlatCon), pp. 1-5. IEEE, 2016.
- ²² Chockwanich, Navaporn, and Vasaka Visoottiviseth. "Intrusion Detection by Deep Learning with TensorFlow." In 2019 21st International Conference on Advanced Communication Technology (ICACT), pp. 654-659. IEEE, 2019.
- ²³ Lynn, Htet Myet, Sung Bum Pan, and Pankoo Kim. "A deep bidirectional GRU network model for biometric electrocardiogram classification based on recurrent neural networks." *IEEE Access* 7 (2019): 145395-145405.
- ²⁴ Yin, Chuanlong, Yuefei Zhu, Jinlong Fei, and Xinzheng He. "A deep learning approach for intrusion detection using recurrent neural networks." *Ieee Access* 5 (2017): 21954-21961.
- ²⁵ Lynn, Htet Myet, Sung Bum Pan, and Pankoo Kim. "A deep bidirectional GRU network model for biometric electrocardiogram classification based on recurrent neural networks." *IEEE Access* 7 (2019): 145395-145405.
- ²⁶ Yin, Chuanlong, Yuefei Zhu, Jinlong Fei, and Xinzheng He. "A deep learning approach for intrusion detection using recurrent neural networks." *Ieee Access* 5 (2017): 21954-21961.
- ²⁷ Kim, Jihyun, Jaehyun Kim, Huong Le Thi Thu, and Howon Kim. "Long short term memory recurrent neural network classifier for intrusion detection." In 2016 International Conference on Platform Technology and Service (PlatCon), pp. 1-5. IEEE, 2016.
- ²⁸ Vigneswaran, K. Rahul, R. Vinayakumar, K. P. Soman, and Prabakaran Poornachandran. "Evaluating shallow and deep neural networks for network intrusion detection systems in cyber security." In 2018 9th International Conference on Computing, Communication and Networking Technologies (ICCCNT), pp. 1-6. IEEE, 2018.
- ²⁹ Vinayakumar, R., Mamoun Alazab, K. P. Soman, Prabakaran Poornachandran, Ameer Al-Nemrat, and Sitalakshmi Venkatraman. "Deep learning approach for intelligent intrusion detection system." *IEEE Access* 7 (2019): 41525-41550.
- ³⁰ Subba, Basant, Santosh Biswas, and Sushanta Karmakar. "A neural network based system for intrusion detection and attack classification." In 2016 Twenty Second National Conference on Communication (NCC), pp. 1-6. IEEE, 2016.

- ³¹ Blanco, Roberto, Pedro Malagón, Juan J. Cilla, and José M. Moya. "Multiclass network attack classifier using CNN tuned with genetic algorithms." In 2018 28th International Symposium on Power and Timing Modeling, Optimization and Simulation (PATMOS), pp. 177-182. IEEE, 2018.
- ³² Vigneswaran, K. Rahul, R. Vinayakumar, K. P. Soman, and Prabakaran Poornachandran. "Evaluating shallow and deep neural networks for network intrusion detection systems in cyber security." In 2018 9th International Conference on Computing, Communication and Networking Technologies (ICCCNT), pp. 1-6. IEEE, 2018.
- ³³ Yin, Chuanlong, Yuefei Zhu, Jinlong Fei, and Xinzheng He. "A deep learning approach for intrusion detection using recurrent neural networks." *Ieee Access* 5 (2017): 21954-21961.

Bibliography

- Amarasinghe, Kasun, and Milos Manic. "Improving User Trust on Deep Neural Networks Based Intrusion Detection Systems." In IECON 2018, 44th Annual Conference of the IEEE Industrial Electronics Society, pp. 3262–3268. IEEE, 2018.
- Ambusaidi, Mohammed A., Xiangjian He, Priyadarsi Nanda, and Zhiyuan Tan. "Building an Intrusion Detection System using a Filter-Based Feature Selection Algorithm." *IEEE Transactions on Computers* 65, no. 10 (2016): 2986–2998.
- Blanco, Roberto, Pedro Malagón, Juan J. Cilla, and José M. Moya. "Multiclass Network Attack Classifier using CNN Tuned with Genetic Algorithms." In 2018 28th International Symposium on Power and Timing Modeling, Optimization and Simulation (PATMOS), pp. 177–182. IEEE, 2018.
- Chiba, Zouhair, Noreddine Abghour, Khalid Moussaid, and Mohamed Rida. "Intelligent Approach to Build a Deep Neural Network Based IDS for Cloud Environment using Combination of Machine Learning Algorithms." *Computers & Security* 86 (2019): 291–317.
- Chockwanich, Navaporn, and Vasaka Visoottiviset. "Intrusion Detection by Deep Learning with TensorFlow." In 2019 21st International Conference on Advanced Communication Technology (ICACT), pp. 654–659. IEEE, 2019.
- Davis, Alexander, Sumanjit Gill, Robert Wong, and Shahab Tayeb. "Feature Selection for Deep Neural Networks in Cyber Security Applications." In 2020 IEEE International IOT, Electronics and Mechatronics Conference (IEMTRONICS), pp. 1–7. IEEE, 2020.
- Eskandari, Mojtaba, Zaffar Haider Janjua, Massimo Vecchio, and Fabio Antonelli. "Passban IDS: An Intelligent Anomaly Based Intrusion Detection System for IoT Edge Devices." *IEEE Internet of Things Journal* (2020).
- Hakim, Lukman, and Rahilla Fatma. "Influence Analysis of Feature Selection to Network Intrusion Detection System Performance Using NSL-KDD Dataset." In 2019 International Conference on Computer Science, Information Technology, and Electrical Engineering (ICOMITEE), pp. 217–220. IEEE, 2019.
- Kim, Jihyun, Jaehyun Kim, Huong Le Thi Thu, and Howon Kim. "Long Short Term Memory Recurrent Neural Network Classifier for Intrusion Detection." In 2016 International Conference on Platform Technology and Service (PlatCon), pp. 1–5. IEEE, 2016.
- Larriva-Novo, Xavier A., Mario Vega-Barbas, Víctor A. Villagrà, and Mario Sanz Rodrigo. "Evaluation of Cybersecurity Data Set Characteristics for Their Applicability to Neural

- Networks Algorithms Detecting Cybersecurity Anomalies.” *IEEE Access* 8 (2020): 9005–9014.
- Lynn, Htet Myet, Sung Bum Pan, and Pankoo Kim. “A Deep Bidirectional GRU Network Model for Biometric Electrocardiogram Classification based on Recurrent Neural Networks.” *IEEE Access* 7 (2019): 145395–145405.
- Ma, Chencheng, Xuehui Du, and Lifeng Cao. “Analysis of Multi-Types of Flow Features Based on Hybrid Neural Network for Improving Network Anomaly Detection.” *IEEE Access* 7 (2019): 148363–148380.
- Nie, Laisen, Zhaolong Ning, Xiaojie Wang, Xiping Hu, Yongkang Li, and Jun Cheng. “Data-Driven Intrusion Detection for Intelligent Internet of Vehicles: A Deep Convolutional Neural Network-Based Method.” *IEEE Transactions on Network Science and Engineering* (2020).
- Nisioti, Antonia, Alexios Mylonas, Paul D. Yoo, and Vasilios Katos. “From Intrusion Detection to Attacker Attribution: A Comprehensive Survey of Unsupervised Methods.” *IEEE Communications Surveys & Tutorials* 20, no. 4 (2018): 3369–3388.
- Obaid, Hadeel S., Saad Ahmed Dheyab, and Sana Sabah Sabry. “The Impact of Data Pre-Processing Techniques and Dimensionality Reduction on the Accuracy of Machine Learning.” In 2019 9th Annual Information Technology, Electromechanical Engineering and Microelectronics Conference (IEMECON), pp. 279–283. IEEE, 2019.
- Sarvari, Samira, Nor Fazlida Mohd Sani, Zurina Mohd Hanapi, and Mohd Taufik Abdullah. “An Efficient Anomaly Intrusion Detection Method with Feature Selection and Evolutionary Neural Network.” *IEEE Access* 8 (2020): 70651–70663.
- Shone, Nathan, Tran Nguyen Ngoc, Vu Dinh Phai, and Qi Shi. “A Deep Learning Approach to Network Intrusion Detection.” *IEEE Transactions on Emerging Topics in Computational Intelligence* 2, no. 1 (2018): 41–50.
- Subba, Basant, Santosh Biswas, and Sushanta Karmakar. “A Neural Network Based System for Intrusion Detection and Attack Classification.” In 2016 22nd National Conference on Communication (NCC), pp. 1–6. IEEE, 2016.
- Taher, Kazi Abu, Billal Mohammed Yasin Jisan, and Md Mahbubur Rahman. “Network Intrusion Detection using Supervised Machine Learning Technique with Feature Selection.” In 2019 International Conference on Robotics, Electrical and Signal Processing Techniques (ICREST), pp. 643–646. IEEE, 2019.

- Tidjon, Lionel N., Marc Frappier, and Amel Mammar. "Intrusion Detection Systems: A Cross-Domain Overview." *IEEE Communications Surveys & Tutorials* 21, no. 4 (2019): 3639–3681.
- Vigneswaran, K. Rahul, R. Vinayakumar, K. P. Soman, and Prabakaran Poornachandran. "Evaluating Shallow and deep neural networks for Network Intrusion Detection Systems in Cyber Security." In 2018 9th International Conference on Computing, Communication and Networking Technologies (ICCCNT), pp. 1–6. IEEE, 2018.
- Vinayakumar, R., Mamoun Alazab, K. P. Soman, Prabakaran Poornachandran, Ameer Al-Nemrat, and Sitalakshmi Venkatraman. "Deep Learning Approach for Intelligent Intrusion Detection System." *IEEE Access* 7 (2019): 41525–41550.
- Xiao, Huaxin, Jiashi Feng, Yunchao Wei, Maojun Zhang, and Shuicheng Yan. "Deep Salient Object Detection with Dense Connections and Distraction Diagnosis." *IEEE Transactions on Multimedia* 20, no. 12 (2018): 3239–3251.
- Yang, Jin, Tao Li, Gang Liang, Wenbo He, and Yue Zhao. "A Simple Recurrent Unit Model Based Intrusion Detection System with dcgan." *IEEE Access* 7 (2019): 83286–83296.
- Yin, Chuanlong, Yuefei Zhu, Jinlong Fei, and Xinzheng He. "A Deep Learning Approach for Intrusion Detection using Recurrent Neural Networks." *IEEE Access* 5 (2017): 21954–21961.
- Zhong, Wei, Ning Yu, and Chunyu Ai. "Applying Big Data Based Deep Learning System to Intrusion Detection." *Big Data Mining and Analytics* 3, no. 3 (2020): 181–195.

About the Authors

Shahab Tayeb, PhD

Dr. Shahab Tayeb is a faculty member with the Department of Electrical and Computer Engineering in the Lyles College of Engineering at California State University, Fresno. Dr. Tayeb's research expertise and interests include network security and privacy, particularly in the context of the Internet of Vehicles. His research incorporates machine learning techniques and data analytics approaches to tackle the detection of zero-day attacks. Through funding from the Fresno State Transportation Institute, his research team has been working on the security of the network backbone for Connected and Autonomous Vehicles over the past two years. He has also been the recipient of several scholarships and national awards including a US Congressional Commendation for STEM mentorship.

MTI FOUNDER

Hon. Norman Y. Mineta

MTI BOARD OF TRUSTEES

**Founder, Honorable
Norman Mineta***
Secretary (ret.),
US Department of Transportation

**Chair,
Will Kempton**
Retired Transportation Executive

**Vice Chair,
Jeff Morales**
Managing Principal
InfraStrategies, LLC

**Executive Director,
Karen Philbrick, PhD***
Mineta Transportation Institute
San José State University

Winsome Bowen
Vice President, Project Development
Strategy
WSP

David Castagnetti
Co-Founder
Mehlman Castagnetti Rosen &
Thomas

Maria Cino
Vice President, America & U.S.
Government Relations
Hewlett-Packard Enterprise

Grace Crunican**
Owner
Crunican LLC

Donna DeMartino
Managing Director
Los Angeles-San Diego-San Luis
Obispo Rail Corridor Agency

John Flaherty
Senior Fellow
Silicon Valley American Leadership
Forum

William Flynn *
President & CEO
Amtrak

Rose Guilbault
Board Member
Peninsula Corridor Joint Power
Board

Ian Jefferies*
President & CEO
Association of American Railroads

Diane Woodend Jones
Principal & Chair of Board
Lea & Elliott, Inc.

David S. Kim*
Secretary
California State Transportation
Agency (CALSTA)

Therese McMillan
Executive Director
Metropolitan Transportation
Commission (MTC)

Abbas Mohaddes
President & COO
Econolite Group Inc.

Stephen Morrissey
Vice President – Regulatory and
Policy
United Airlines

Dan Moshavi, PhD*
Dean
Lucas College and Graduate School
of Business, San José State
University

Toks Omishakin*
Director
California Department of
Transportation (Caltrans)

Takayoshi Oshima
Chairman & CEO
Allied Telesis, Inc.

Greg Regan
President
Transportation Trades Department,
AFL-CIO

Paul Skoutelas*
President & CEO
American Public Transportation
Association (APTA)

Kimberly Slaughter
CEO
Sysra USA

Beverley Swaim-Staley
President
Union Station Redevelopment
Corporation

Jim Tymon*
Executive Director
American Association of State
Highway and Transportation
Officials (AASHTO)

* = Ex-Officio

** = Past Chair, Board of Trustees

Directors

Karen Philbrick, PhD
Executive Director

Hilary Nixon, PhD
Deputy Executive Director

Asha Weinstein Agrawal, PhD
Education Director
National Transportation Finance Center Director

Brian Michael Jenkins National Transportation
Security Center Director

