

IoT Intrusion Detection Competition using Machine Learning

Applied Machine Learning – Group Practical Project

Timothy Chan – *MSc in Data Science*, Ian Dickerson – *MSc in Advanced Computing Technologies*,
Mike Jun Ming – *MSc in Computing for Financial Services*, Michelangelo Rubino – *MSc in Data Science* and Cosmin Stanciu – *MSc in Data Science*

STRUCTURE

This document has been subdivided into ten parts: an introduction, two group components, five individual components and two appendixes. They have been listed with the responsible and the word count:

- I. Overview;
- II. Planning – Group;
- III. Pre-processing – Michelangelo Rubino, Student No. 13140302 – 600 words;
- IV. Selecting features – Ian Dickerson, Student No. 12407190 – 600 words;
- V. Exploring and selecting ML algorithms - Timothy Chan, Student No. 13160525 – 600 words;
- VI. Refining algorithms - Cosmin Stanciu, Student No. 13179534 – 600 words;
- VII. Evaluating model and analyzing the results - Mike Jun Ming, Student No. 13126984 – 600 words;
- VIII. Future Work – Group;
- IX. Appendix A (Figures);
- X. Appendix B (Tables).

I. OVERVIEW

THE Internet of Things (IoT) is growing continuously in the last years for a number of aspects: the fields it is applied to (healthcare, agriculture, cars...), the number of devices, the network traffic. At the same time the cyber-attacks against the 802.11 networks are increasing steeply

and becoming more and more aggressive. In the most recent literature on the subject the cyber-attacks have been often subdivided into four categories, based on the execution mode [1]: injection, flooding, impersonation, passive attack. These are the reasons why the Intrusion Detection System (IDS) is one of the critical components of a wireless network: in the last decade the machine learning techniques, from more traditional ones [1] to neural networks [2], [3], have been applying to IDS as more static prevention systems have shown unsatisfying rates in detecting false positives and false negatives and in adapting to new threats. A clear explanation of the differences between IDSs and the other detection systems can be found in [4].

II. PLANNING

This work uses a reduced version of the AWID dataset and focuses only on the impersonation type of attack. This because the impersonation attacks have been identified as the most dangerous threats to the 802.11 networks [1]. The goal is to build a binary classifier able to identify real intrusions in the network traffic using the machine learning techniques.

It has been agreed to use a Gantt chart (Fig. 1) to schedule the workload, even though this kind of chart is more suitable to a traditional waterfall methodology. The Machine Learning (ML) projects instead use an iterative lifecycle: several iterations to find the best settings. The idea in this project is then to spend the first six weeks in building a model with the simplest techniques to have a baseline and then the following twenty days in testing a certain number of more complex algorithms and in adjusting parameters and hyperparameters.

There are several challenges in this research area which are described in the following subsections along with the approach followed in this work.

A. *Known threats vs Unknown threats*

The traditional prevention systems are often reliable for the known threats but not for new ones. The method a threat is detected can be categorized into three main types [2]: misuse-based (or signature-based), anomaly-based, specification-based. The first group is considered reliable in detecting threats already known, though not the same can be said for the new ones. The third group relies on a set of rules manually written and updated frequently. These two methods have been discarded, then the chosen one for this work is the anomaly-based one, which uses the machine learning techniques to build the user behavior considered “normal”, while the others are labelled as “threat”.

B. *Feature selection*

This is one of the most delicate stages in the ML process. The general question is whether to use every feature in the dataframe to identify the threat or just part of them, typically those considered more relevant for the class. In the previous literature, the prevalent choice is to use only a certain number of features, like 12-22 in [2], 10-20 in [3] or 36 in [5]: the reasons are summarized in [6]: “The objective of variable selection is three-fold: improving the prediction performance of the predictors, providing faster and more cost-effective predictors, and providing a better understanding of the underlying process that generated the data”. In this work it has been decided to use from 20 to 40 features based on a mix of criteria (correlation, Mutual Information, PCA).

C. *The algorithm to use*

Due to the limited amount of time it is not possible to test every machine learning technique, the focus will be on those more frequently used for binary classifiers.

D. *Model Size*

The traffic of networks 802.11 is generated by portable devices having a small memory where the IDS can be implemented, then the need for light models is crucial and is strictly linked to the feature selection and the algorithm to use. A light model seems more appropriate for this kind of projects for the reasons previously described, provided that this is not to detriment of the performances: not necessarily a lighter model has low performances and a more complex one has better performances, then it is necessary to find a trade-off between the complexity of the model and the results.

E. *Precision vs Recall*

It is commonly accepted that for binary classifiers the Accuracy is misleading and may lead to wrong conclusions about a model [7] as it detects only the correct predictions, positive or negative, while it says nothing about the false positives and negatives. For this reason, it is important to focus on two more relevant metrics, Precision and Recall. Precision and Recall show respectively the proportion of relevant instances out of the total number of instances and the proportion of relevant instances out of the total number of relevant instances. High Precision means a low false positive rate, while high recall means a low false negative rate. These metrics are like two opposing forces as increasing the first one decreases the second one and vice versa. This is known as the Precision/Recall trade-off and it is often necessary to establish the metric to focus on. For a binary classifier like the one proposed, the question should be: is it desirable to have a low false positive rate or a low negative rate? Is it desirable to see a lower number of normal behaviors wrongly classified as threats or a lower number of real threats wrongly classified as normal behaviors? It is clearly better to focus on the recall as this would limit the number of real threats misclassified as non-threats. **What have we focused on?**

F. *ML knowledge vs Domain knowledge*

This is a wide topic and is out of the scope of this paper. It can be said however that a certain

knowledge of the domain where the ML model is applied is beneficial to the accuracy ([8], [9]).

III. PRE-PROCESSING

AWID is a public collection of standard data sets often used to evaluate the performances of an IDS, for this project it will be used the reduced version AWID-CLS, comprising a training and a test set. The class is the column 155:

Dataset	Obs.	Features	Classes
AWID-CLS-Train	97,044	152	1
AWID-CLS-Test	40,158	152	1

After loading the data through the pandas library and creating a dataframe, it is possible to see that both the datasets are perfectly balanced, 50% of the observations classified as 1, 50% as 0. This is an important aspect as unbalanced datasets may lead to wrong predictions. The dataframe contains no nulls, then it is not necessary to handle them. The first exploration of the raw data can be done through the pandas functions `info()` and `head()`. The `describe()` function shows the descriptive statistics and the first thing to notice is that many features contains only zeros, then the features are useless. They can be safely removed (74), ending up with 78 from the original 152. In the following figure three features with no values have been highlighted:

	1	2	3	5	6
count	97044.0	97044.0	97044.0	97044.000000	97044.000000
mean	0.0	0.0	0.0	0.006252	0.006252
std	0.0	0.0	0.0	0.015541	0.015541
min	0.0	0.0	0.0	0.000003	0.000003
25%	0.0	0.0	0.0	0.001442	0.001442
50%	0.0	0.0	0.0	0.003706	0.003706
75%	0.0	0.0	0.0	0.005916	0.005916
max	0.0	0.0	0.0	0.978440	0.978440

By combining the pandas dataframe functions `loc()` and `duplicated()` it is possible to detect duplicated features, having exactly the same values. These features can be removed as well, leading to a dataframe with only 64 features and the class. Decreasing the number of features is an important stage for some reasons:

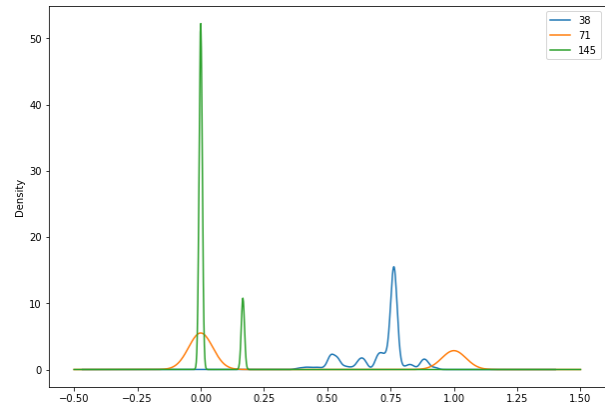
- it reduces the processing time;
- it simplifies the model, avoiding overfitting;

- focusing only on relevant features, it improves the results.

The real selection of features will happen in the next stage, though removing useless features can be done in the preprocessing stage. The unique values show that the features can be categorized into three main groups:

1. continuous values between 0 and 1 (ex: 38);
2. only two values, 0 or 1 (ex: 71);
3. only a few values, integer or floating, between 0 and 1 (ex: 145).

As every feature show only values between 0 and 1, the data does not need to be rescaled. This term refers to the techniques used to give the features the same scale (typically between 0 and 1), which may give the model some improvements [10]. There is often confusion when using the terms scaling, standardizing and normalizing as authors tend to use them interchangeably, then the meanings explained in [11] will be used. The three abovementioned groups have different distributions and none is normal:



Rather than using the correlation matrix plot or the scatter plot matrix, which are not easy to visualize due to the number of features, it has been computed the correlation among pairs of features and established a threshold of 0.9: for each pair having a correlation equal or higher than the threshold, the second feature will be removed. This will avoid multicollinearity: in [12] multicollinearity is defined as “a phenomenon in which one feature variable in a regression model is

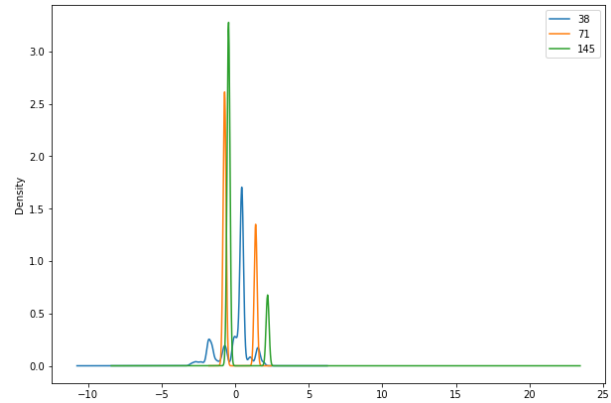
highly linearly correlated with another feature variable” and this may affect the model and lead to misleading results. Multicollinearity will never improve the results, on the contrary avoiding it may bring advantages to the model. There are several techniques used to avoid multicollinearity: for this project it will be removed the second feature of the highly correlated pairs. When dealing with dataframes showing a gaussian distribution, the Pearson correlation is preferred, here instead the non-parametric correlations will be used. The pandas library contains correlation methods different from Pearson’s such as the Kendall rank correlation coefficient (Kendall’s tau coefficient) or the Spearman rank correlation coefficient (Spearman’s rho). These methods do not rely on the normal distribution assumption and seem more appropriate. A more detailed explanation of the non-parametric methods can be found in [13]. The Kendall and Spearman correlations are quite similar and indeed the features exceeding the 0.9 threshold are the same, apart from one. The 19 features in common will be removed (Tab. 1X.).

After dropping the unnecessary features, the dataframe has 45 features and 1 class. Plotting the features and computing their skewness confirm that no features have a normal distribution: the histogram, the density and the box and whisker plots of the features can be seen in the Appendix A (Fig. 2-4). To improve the data distribution some techniques have been tested using the scikit-learn functions: yet, it has been arranged to use only the standardization through the StandardScaler function, which removes the mean and scales to unit variance. In the following table a short explanation for each technique and why it has been discarded:

Function	Description	Explanation
Binarizer	Continuous values are converted into only two values, 0 and 1, based on if less or more than 0.50	No improvement, on the contrary the results are slightly worse.
Power Transformer	It helps the data to look more Gaussian-like by using two methods: Box-Cox or Yeo-Johnson	Similar to StandardScaler
MinMaxScaler	It scales data to values between 0 and 1	Not necessary as the feature values are all included between 0 and 1

RobustScaler	Similar to MinMaxScaler, in addition it is useful to deal with outliers	Not necessary
Normalizer	It helps the data to look more Gaussian like	It works at row level. Not recommended [11]

The standardization through the StandardScaler will help the dataframe to show a distribution closer to the normal one (mean = 0, standard deviation = 1). It will not be a perfect normal distribution, but this is better than the original one as shown by the following figure with the same features previously analyzed:



IV. SELECTING FEATURES

It is possible to train a classification model on a large number of features, but this has clear disadvantages. The time taken for training can be prohibitive. And there is a danger that the model will fit the noise in irrelevant features, reducing generalizability. A selection of fewer, demonstrably relevant features should improve classification performance and processing time while in some cases also yielding better interpretability.

To avoid feature selection becoming a processing bottleneck in this project, a broadly filter-based approach was adopted. That is, testable characteristics of individual features were used as the primary selection method. A logistic regression classifier, trained and tested on a 60/40 split of the training data, was used to sense-check the feature selection and fine tune the thresholds used. The recall of the classifier was used as the principal metric, as the priority was to make sure that as

many intrusions are detected as possible, but in practice precision was equally high. The disadvantage here is that a set of features that leads to good recall from a logistic regression model might not transfer well to a more complex classification model like Random Forest or a neural network. The filter methods employed were:

- *Statistical hypothesis testing* - If the distribution of values in a feature does not change under the two classes, then the feature will be of no value for the classifier model. Independent t-tests are a computationally fairly cheap process to eliminate such features. Tests were carried out for all the features remaining after preprocessing and the features ordered by test statistic. Since almost all of the features showed significantly -- that is, with vanishingly low p-value -- different distributions under the two conditions, this method was of limited value in the present case;
- *Correlation* - A second common and quite efficient method for assessing relevance is to measure the strength of correlation between feature and target classification. Either a strong positive or strong negative correlation with the target suggests that the feature will have predictive value. The linear correlation measures employed here was Kendall's tau and Pearson's rho which gave sufficiently different results to warrant the inclusion of both. There are two obvious limitations to these measures: firstly, that as statistical measures they are not designed to be applied to a binary variable, and secondly that they measure linear correlation, and may miss a relationship with the target variable that is not linear. The first weakness appears not to have too great an impact in practical applications. To account for the second possibility, another measure of relationship is needed.
- *Mutual Information* - To capture relationships with the target variable that are not linear, or otherwise not captured by

correlation, the Mutual Information between individual features and the target class was calculated, following Parker et al. One possible weakness of this approach may be that the highest-ranked features for MI may share information with each other as well as the target and that lower ranked features may encode genuinely different information. In the present case, measurement of MI between so many features was found to be impractical.

- *PCA* - To compensate for the possible weaknesses in MI and linear correlation identified above, Principal Component Analysis was used to capture any remaining information from features that had not so far been identified. This unsupervised technique captures as much variance as possible from high dimensional data into a lower dimensional space.

Good recall was found by exploring combinations of A) five to ten features with high absolute value of Kendall's tau or Pearson's rho with the target, B) five to ten features with highest mutual information with the target, and C) the first ten to twenty principle components of the continuous features. A Python function was written to allow the developer of the model to easily vary the thresholds to obtain slightly different feature sets for training.

V. EXPLORING AND SELECTING ML ALGORITHMS

A. Random Forests

Previous research [14] and alert detection studies [15] suggests trees perform well for classification with mixed data. RandomForests provide a model with enough flexibility without making assumptions about the data, with tools to reduce high variance on large data sets.

Untrimmed trees tend to overfit and produce many leaf nodes. RandomForests improve over traditional trees, through tools to reduce overall error and variance; bagging and random feature sampling. Generating multiple trees, through random sampling with replacement, of subsets of test data and features. This reduces the overall

reliance of prediction on one single model.

RandomForests benefit from identifying feature importance, unlike other models, but are relatively blackbox in nature; there is little scope to amend parameters.

B. *K-Nearest Neighbors (K-NN)*

K-NN was selected for its relative simplicity compared to other blackbox models. The algorithm works by grouping datapoints, based on other datapoints which are ‘close’ to one another’. Each new datapoint is classified by its K nearest partners. The model therefore only requires setting a distance calculator and K.

Due to classification being solely reliant on a simple distance calculation, there is a likelihood that the model is not flexible enough for a large data set, with a large feature space. A small change in a binary feature across data rows, may lead to only a small change in distance between datapoints. This binary feature, however, may be critical in the true categorization model.

The initial use of Euclidean distance poses issues in measuring distance between continuous and categorical variables. Selection of a distance calculator such as a Gower [16], may reduce the disadvantages outlined.

C. *Binary Linear Logistic Regression*

LR was selected due to its probabilistic nature. The calculation of logits makes the model relatively easy to interpret and highlight feature importance.

The assumptions within LR however, poses issues for the large data set with a high number of features. LR requires continuous independent variables to have linear relationships with their logits. Secondly, there cannot be multi-collinearity between independent variables. Both can lead to high standard errors, and therefore poor prediction reliability [17].

Due to the data set being unlabeled, and assumptions made about the types data contained within the set, we cannot ensure that these assumptions are met.

D. *MLP Neural-Network (Stochastic Gradient Descent)* [18]

A NN was selected due to its flexibility on large data sets with a higher number of features. NN and Deep NN has been particularly successful in classification accuracy tasks [19]. Further and recent research has shown the increased capabilities of CNN and ANN in intrusion detection [20].

Neural-Networks complexity, and computational cost lead to significant disadvantages. Owing to backpropagation, the algorithm classifies based on patterns identifying patterns between layers, and patterns between neurons. Many neurons and layers can lead to a computationally expensive model, with a tendency to overfit. Owing to how these patterns are formed, NN are blackbox; little information can be gained from viewing these patterns.

The ability to balance model overfitting, variance and computational cost, therefore, requires time spent iteratively finetuning the selection of the number of layers and neurons used in each layer.

E. *Support Vector Machine (SVM)*

SVM was selected due to its relative simplicity compared to blackbox models. SVM transforms the data into N dimensional space (N=Number of features). The transformed data is separated by the maximal margin hyperplane, with threats and non-threats falling on either side.

SVM works under the assumption the transformed data is separable by a hyperplane. This has benefits; It is easy to understand and tends to not overfit.

However, difficulties arise in a large data set with selecting the correct kernel (shape) of the hyperplane and gamma (smoothness of the hyperplane) – requiring numerous iterations to find the correct parameters. The assumption of the data being separable, in itself, also limits the models flexibility compared to blackbox models; RF and NN.

A model initialization and parameter justification can be found in Fig. 5.

VI. REFINING ALGORITHMS

Cosmin here...

VII. EVALUATING MODEL AND ANALYSING THE RESULTS

Mike here...

VIII. FUTURE WORK

Data vs Model... more data means better accuracy?

IX. APPENDIX A (FIGURES)

GANTT CHART

[illegible]

Fig. 1

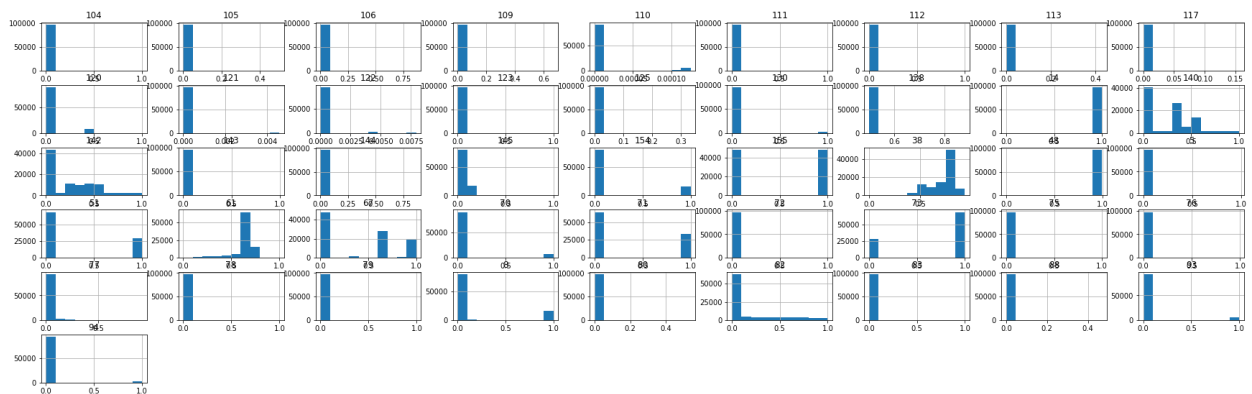


Fig. 2

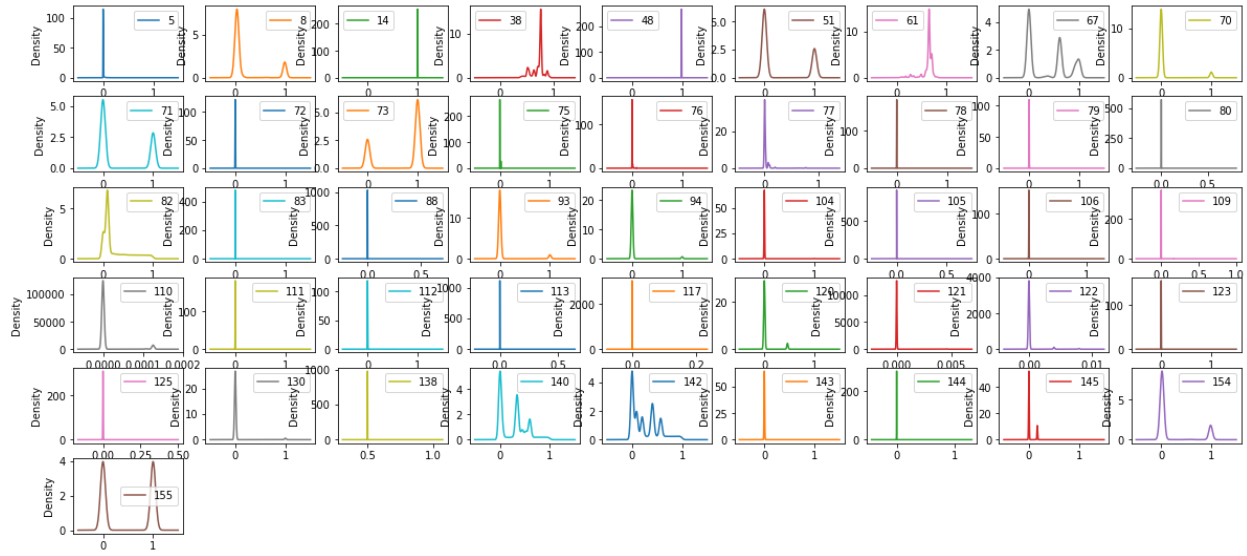


Fig. 3

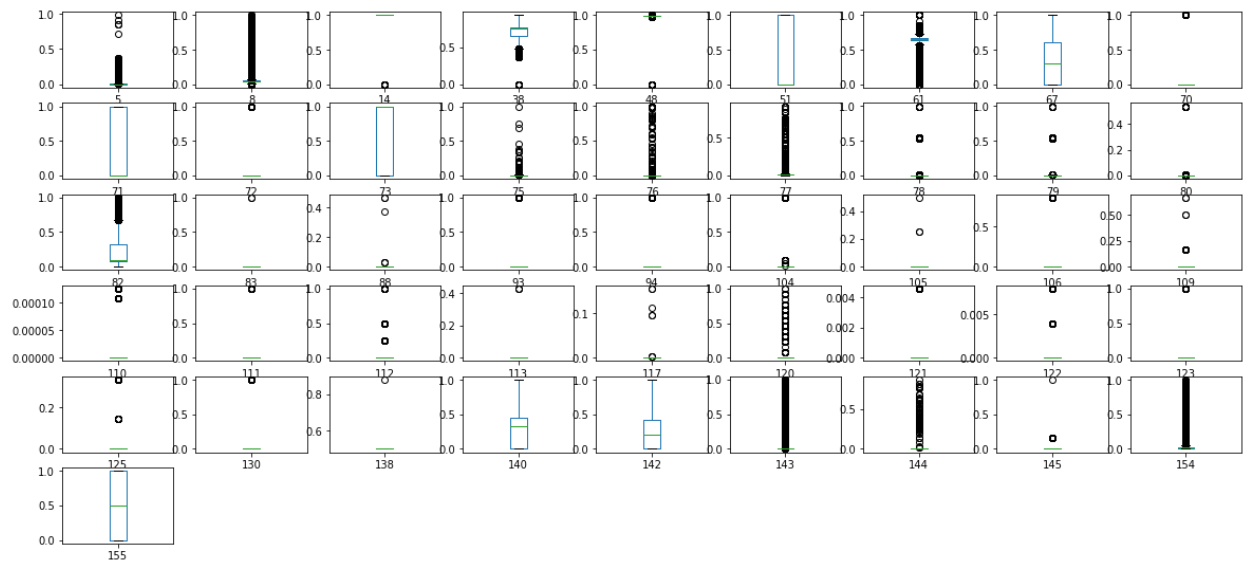


Fig. 4

Model	Parameters	Justification
Random Forest	(1) Untrimmed, (2) Depth limited to 6	(1) Test the expectation of a very accurate model and explore feature importance. (2) To initially access how trimming reduces the precision and recall.
K-NN	Euclidean Distance, 5 neighbours	Euclidean distance for it's simplicity. 5 Neighbours, square root of number of features.
Binary Logistic Regression		Explore feature feature importance.
Neural Network	2 Layers, 5 Neurons. 1000 Iterations of BackPropogation	Small number of layers and neurons to access suitability and reduce overfitting.
SVM	Kernel: RBF, Gamma = $1/n$, Cost = 1	Flexible kernel due to high dimensionality. Low cost to assess separability of data and low gamma to reduce flexibility of separating hyperplane

Fig. 5

X. APPENDIX B (TABLES)

TABLE 1

Features having high Kendall and Spearman correlation

Feature	Kendall corr ≥ 0.9	Spearman corr ≥ 0.9
47	Yes	Yes
50	Yes	Yes
64	Yes	Yes
66	Yes	Yes
67	No	Yes
68	Yes	Yes
84	Yes	Yes
86	Yes	Yes
90	Yes	Yes
97	Yes	Yes
98	Yes	Yes
107	Yes	Yes
108	Yes	Yes
118	Yes	Yes
119	Yes	Yes
126	Yes	Yes
127	Yes	Yes
128	Yes	Yes
129	Yes	Yes
141	Yes	Yes

XI. REFERENCES

- [1] C. Kolias, G.Kambourakis, A.Stavrou and S. Gritzalis, "Intrusion Detection in 802.11 Networks: Empirical Evaluation of Threats and a Public Dataset," *IEEE Communications Surveys and Tutorials*, vol. 18, no. 1, 2015.
- [2] M.E.Aminanto, R.Choi, H.C.Tanuwidjaja, P.D.Yoo, "Deep Abstraction and Weighted Feature Selection for Wi-Fi Impersonation Detection" *IEEE Transactions on Information Forensics and Security*, vol. 13, no. 3, 2018.
- [3] L.R.Parker, P.D.Yoo, T.A.Asyhari, L.Chermak, Y.Jhi, K.Taha, "DEMise: Interpretable Deep Extraction and Mutual Information Selection Techniques for IoT Intrusion Detection" in *Proceedings of the 14th International Conference on Availability, Reliability and Security*, Article no. 98, 2019.
- [4] Rao U.H., Nayak U. – "Intrusion Detection and Prevention Systems" in *The InfoSec Handbook*. Apress, Berkeley, CA, 2014.
- [5] S. Rezvy, M. Petridis, T. Zebin, Y. Luo, A. Lasebae "An efficient deep learning model for intrusion classification" in *53rd Annual Conference on Information Sciences and Systems (CISS)*, 2019.
- [6] I. Guyon, A. Elisseeff – "An Introduction to Variable and Feature Selection" in *Journal of Machine Learning Research* 3, 2003.
- [7] Y. Jankay – "The 3 pillars of Binary Classification: Accuracy, Precision & Recall" (2018, April). <https://medium.com/@yashwant140393/the-3-pillars-of-binary-classification-accuracy-precision-recall-d2da3d09f664>. [Online].
- [8] Anand – "Why domain knowledge is important in Data Science" (2019, March). <https://medium.com/@anand0427/why-domain-knowledge-is-important-in-data-science-anand0427-3002c659c0a5>. [Online].
- [9] T. Blanchard – "The importance of domain knowledge – A healthcare data science perspective" (2017, November). <https://data-science-blog.com/blog/2017/11/10/the-importance-of-domain-knowledge-a-healthcare-data-science-perspective/>. [Online].
- [10] S. Raschka – Department of Statistics, University of Wisconsin-Madison. "About Feature Scaling and Normalization" (2014, July). https://sebastianraschka.com/Articles/2014_about_feature_scaling.html. [Online] - e-mail: sraschka@wisc.edu.
- [11] J. Hale. "Scale, Standardize or Normalize with Skikit-Learn" (2019, March). <https://towardsdatascience.com/scale-standardize-or-normalize-with-scikit-learn-6ccc7d176a02> [Online].
- [12] J. Saslow. "Collinearity – What it means, why it is bad and how does it affect other models?" (2018, July). <https://medium.com/future-vision/collinearity-what-it-means-why-its-bad-and-how-does-it-affect-other-models-94e1db984168> [Online].
- [13] J. Brownlee. "How to calculate Nonparametric Rank Correlation in Python" (2018, July). <https://machinelearningmastery.com/how-to-calculate-nonparametric-rank-correlation-in-python> [Online].
- [14] Intrusion Detection in 802.11.
- [15] Deep Abstraction and Weighted Feature Selection.
- [16] Gower, J. (1971). A General Coefficient of Similarity and Some of Its Properties. *Biometrics*, 27(4), pp.857-871.
- [17] Stoltzfus, J. (2011). Logistic Regression: A Brief Primer. *Academic Emergency Medicine*, 18(10), p.1101.
- [18] Kingma, D. and Ba, J. (2017). Adam: A Method for Stochastic Optimization.
- [19] Moradi, M. and Zulkernine, M.,(2014). A Neural Network Based System for Intrusion Detection and Classification of Attacks. Conference: *Proceedings of 2004 IEEE International Conference on Advances in Intelligent Systems - Theory and Applications*.
- [20] H., B. and Tindo, G. (2017). Network Intrusion Detection Systems based Neural Network: A Comparative Study. *International Journal of Computer Applications*, 157(5), pp.42-47.