

Intrusion Detection Prediction Using Data Science Technique

Thyagarajan C⁽¹⁾, Karthick Aravind B⁽¹⁾, Peniel Branham T⁽³⁾, Balamurugan S⁽⁴⁾

⁽¹⁾Assistant Professor, Department of CSE, Panimalar Engineering College, Chennai-600123, India

^(2,3,4)Student, Department of CSE, Panimalar Engineering College, Chennai-600123, India

Mail – thyaguwinner@gmail.com, thekarthickaravind@gmail.com, t.penielfranham2002@gmail.com, balamurugan1312002@gmail.com

Abstract— Enterprise networks are protected from cyberattacks via intrusion detection systems. However, the ability to produce a lot of poor-quality alerts continues to be a drawback. The most recent research on Intrusion Alert-based cyberattack prediction, including its models and shortcomings, has been examined in the current study. There has been a significant increase in the frequency and power of intrusion attempts globally, and systems for detecting and anticipating attacks are currently being developed in great detail. The design of attack prediction tools has typically been dominated by statistical methodologies. The study employed supervised machine learning technologies to analyze the dataset and extract information such as variable identification, univariate analysis, bivariate and multivariate analysis, and missing value treatments. (SMLT). To ascertain which machine learning algorithm was most effective at forecasting cyberattacks, a comparative analysis of the algorithms was conducted. In terms of accuracy, precision, recall, F1 Score, sensitivity, and specificity, the findings demonstrate that the suggested machine learning algorithm technique is on par with the best.

Keywords: Intrusion Detection Systems, Machine learning algorithms, Attack prediction, Statistical methods

I. INTRODUCTION

The objective of this study is to develop an efficient machine learning model for prediction of intrusion detection that can replace traditional supervised machine learning classification models by providing the highest accuracy in comparison to supervised techniques. The project concentrates on the scope of an Intrusion Prevention System (IPS), a useful device for preventing assaults by blocking malicious packets, dropping offending IPs, and warning security staff of potential dangers. The system may be trained to recognize assaults based on traffic and behavioral anomalies and is built on an existing database for signature identification. An Intrusion Detection System's (IDS) primary job is to look for abnormal activity and send out notifications when it does. After receiving these signals, the type of attack must be determined so that the security operations center (SOC) analyst or incident responder may investigate the situation and take the appropriate precautions to remove the threat.

Intrusion detection systems can keep an eye on networks for potentially hostile activity, but they can also produce a lot of false alerts. Therefore, at the initial installation stage, businesses must adjust their IDS products. Setting up intrusion detection systems correctly is essential to distinguish between legitimate network traffic and malicious behavior. This can improve the precision of identifying real security threats and decrease the number of false alarms. Therefore, by detecting potential security breaches, network intrusions, and other abnormal activities that may pose a serious risk to crucial assets and data, the development of an effective and accurate machine learning model for intrusion detection prediction can significantly contribute to the security of enterprise networks.

The proposed model for anomaly detection in this paper is focused on solving the difficulty of linking to and detecting fraudulent and suspicious behavior, network intrusions, and other anomalous events that are difficult to describe using conventional methods. A robust and reliable machine learning method will be created using the suggested model's colorful data wisdom techniques, such as variable identification, data preparation, and visualization. The system will be trained using a dataset that has already been collected informally, and it will use the data to improve its effectiveness in spotting anomalies.

The main objective of the proposed model is to compare and assess the performance metrics of several colored methods in order to select the optimal bone for anomaly detection. To do this, the model will evaluate each algorithm's performance using a variety of metrics, including delicacy, perfection, recall, F1 score, perceptivity, and particularity. The proposed approach may determine the most appropriate bone for directly identifying anomalies in a given dataset by comparing the results obtained from each technique.

The successful application of this model can aid in the creation of an anomaly discovery system that can assist associations in preventing security breaches and protecting their vital assets and data from hidden perils. In the current digital era, where cyberattacks are becoming less sophisticated and rarer, this is extremely crucial. Associations can take appropriate action to correct the dangers and prevent further harm to their systems by spotting irregularities in advance.

II. EXISTING SYSTEM

It might be difficult to incorporate supervised Machine Learning (ML) techniques into Network Intrusion Detection Systems (NIDS). Labelled datasets with both benign and malicious samples are necessary to train and assess an ML-NIDS. However, getting these labels is a difficult, pricey, and expert knowledge-required procedure, which has led to a dearth of actual deployments and a reliance on old information in research articles. Labelled datasets have recently started to become more widely available, however many earlier studies only used them as test subjects without fully realising their potential.

Despite a number of triumphs, supervised ML method integration in NIDS is still in its infancy. This is partly because it is challenging to gather the large quantities of labelled data required for ML-NIDS evaluation and training. However, the research community has welcomed the recent release of labelled datasets for ML-NIDS. Few articles have mentioned the idea of improving the state-of-the-art in intrusion detection systems, despite the fact that the availability of tagged datasets provide this option. By utilizing the capabilities of supervised machine learning, it is now possible to significantly increase the accuracy and efficacy of NIDS thanks to the availability of labelled datasets.

DISADVANTAGES:

1. They are not predicting the classification of attack types.
2. They are not mentioning the accuracy.
3. They are using machine learning technique only for analyzing purpose.

III. PROPOSED SYSTEM

The suggested model seeks to provide a machine learning model for anomaly detection that is both accurate and effective. An important tool for identifying difficult-to-detect anomalous events such as network intrusions, fraudulent activity, and suspicious activity is anomaly detection. The right data science methods are used to develop the machine learning model, including variable identification, data preparation, and data visualization. The machine learning model is built using a prior dataset, and as it learns and is trained using different techniques, the model's accuracy is increased. The objective is to assess the model's performance metrics and compare them in order to choose the best algorithm for anomaly detection. The suggested model intends to make a contribution to the creation of an effective anomaly detection system that can assist organizations in preventing security breaches and safeguarding their valuable assets and data from potential threats.

ADVANTAGES:

1. To perform classification, we are using a machine learning method.
2. To compare machine learning algorithms and find the one with the best accuracy, more than two are utilized.
3. Deployment is carried out to obtain results.

IV. MODULE DESCRIPTION

Data Pre-processing:

Machine Learning (ML) relies heavily on validation procedures to calculate a model's error rate. These methods aid in assessing how well the model will function in a real-world environment. The use of validation techniques might not be necessary when the data amount is substantial and accurately reflects the population. However, in the majority of real-world situations, we deal with data samples that might not accurately reflect the population of the given dataset. In these situations, validation procedures are used to precisely calculate the model's error rate.

Pre-processing data is essential in ML to get it ready for model training. This stage includes identifying duplicate and missing values as well as the data type, such as an integer or float variable. Data that has been preprocessed can then be separated into training and validation sets. After the training set has been used to fit the model, the validation set is used to evaluate the model's performance.

When adjusting model hyperparameters, machine learning experts employ the validation set to provide an unbiased evaluation of a model's fit to the training dataset. If skill from the validation dataset is incorporated into the model architecture, the evaluation may become more skewed. The validation set is frequently evaluated to adjust the model hyperparameters. It might take a lot of time to gather data, analyze it, and handle issues with data content, quality, and organization. However, the procedure for identifying data aids in comprehending the data and its characteristics. This information can be used to select the algorithm to create the model that most closely matches the properties of the data. We may create trustworthy and accurate models that perform well in real-world scenarios by using the right validation approaches and adhering to the data pre-processing stages.

Data visualization:

Data visualisation is a crucial ability for both machine literacy and applied statistics. Data visualisation offers a set of tools to develop a qualitative understanding of the data, in contrast to statistics, which largely concentrates on quantitative descriptions and estimates of data. When examining a dataset, data visualisation may be immensely beneficial because it makes it possible to spot trends, outliers, and implicit crimes in the data. Data visualisations are

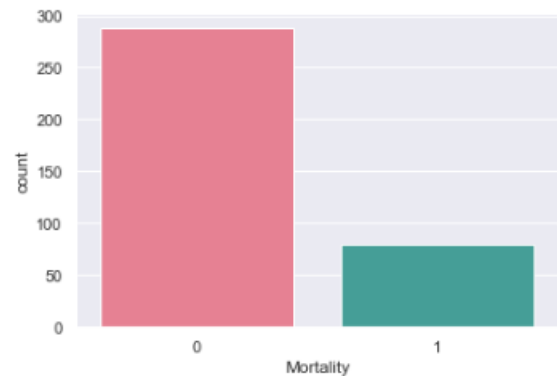


Fig 1: Data visualization

possible with certain sphere expertise compared to measurements of association or significance, can be applied to plot and map significant linkages in a way that is easier for stakeholders to comprehend. The fields of exploratory data analysis and data visualization provide a closer look into the topic. Sometimes data needs to be shown in a visual format, such maps or plots, in order to make sense. Being able to quickly fantasize data samples and others is a crucial ability for both applied machine literacy and applied statistics. When imaging data in Python, there are many different sorts of plots that you need to be familiar with. Knowing how to utilize them will help you better understand your data. You can develop critical perception and build more well-informed opinions in your statistical analyses and machine literacy models by studying and picturing your data.

V. ALGORITHM IMPLEMENTATION:

It's crucial to routinely assess how different machine learning algorithms work. This is essential because the performance of the model as a whole can be significantly impacted by selecting the best method for a given challenge. Each algorithm has unique strengths and disadvantages. To do this, a test harness that compares several methods using a common framework can be developed. The test harness can be expanded to incorporate more algorithms for comparison and used as a model for assessing machine learning algorithms on diverse challenges. The accuracy of a machine learning model, or its capacity to accurately forecast a result based on unobserved data, is one of the primary metrics used to assess a model. It can be difficult to determine a model's performance on unknown data, though. Utilizing resampling techniques, such as cross-validation, to estimate each model's potential accuracy on hypothetical data is one way to deal with this problem. This can provide a more accurate evaluation of model performance and be used to select the most suitable approach for a certain problem. Data visualization can offer meaningful information about the data and make it simpler to identify trends and outliers, making it a vital step in machine learning. The accuracy, variance, and other characteristics of the distribution of model accuracies can also be revealed by visualizing model performance. This can aid in choosing the optimal algorithm for a particular problem and also shed light on potential areas for model enhancement.

In conclusion, choosing the optimal machine learning algorithm for a particular task requires comparing the performance of many algorithms. A test harness can be used to consistently compare different algorithms, and resampling techniques like cross-validation can be used to calculate each model's accuracy on hypothetical data. By visualizing the data and model performance, relevant information about the properties of the data and the accuracy of the model can be discovered while selecting the best solution for a particular problem.

There are four possible results in True positives (TP), true negatives (TN), false positives (FP), and false negatives (FN). False positives (FP) happen when a person who will pay is anticipated to be a defaulter, even though their real class isn't at fault. False negatives (FN) happen when a defaulter is anticipated to be a

player, however, the real class is true and the projected class is false. A real affirmative occurs when a potential defaulter is forecasted and both the actual class and the expected class are yes. (TP). True negatives (TN) happen when a defaulter is supposed to be a player but both the actual class and the forecasted class are no.

To evaluate a binary classification model's efficacy, we can compute the true positive rate (TPR) and false positive rate (FPR). TPR, which is defined as $TP / (TP + FN)$, is the proportion of real positives that the model correctly classifies as positives. The proportion of true negatives that the model wrongly interprets as positives is determined as $FP / (FP + TN)$, or FPR. In assessing trade-offs between the cost of false positives and false negatives in a particular area, these indicators are crucial.

Accuracy is a popular performance indicator used to evaluate the effectiveness of machine learning models. It can be calculated as the proportion of accurately anticipated observations to all observations. A model is 90% accurate, for instance, if it accurately predicts the outcome of 90 out of 100 observations. Although accuracy is a simple and straightforward statistic, it may not always be the ideal option, especially when dealing with asymmetric datasets. When there are significantly more instances of one class than the other in a dataset, this is referred to as an asymmetric dataset. In certain situations, accuracy might not be the appropriate metric to assess the model's performance because it could give an inaccurate image of the model's efficacy. This is so that the model won't provide accurate forecasts for the minority class and won't be biased in favor of the dominant class. To solve this issue, one might use other performance indicators including precision, recall, F1-score, and area under the curve (AUC). These measures give a more complete picture of the model's performance because they take into consideration true positives, false positives, true negatives, false negatives, and both. Therefore, while choosing relevant performance measures to evaluate machine learning models, it is crucial to take the dataset's features and the particular problem at hand into account in addition to accuracy.

SVM:

The Support Vector Machine (SVM) is a well-known supervised machine learning technique that is typically used for classification tasks; however, it may also be used for regression. SVM's primary objective is to locate a hyperplane that can sharply delineate the data points in an N-dimensional space. SVM accomplishes this by mapping the input data to a high-dimensional feature space, which allows it to categorize the data points even when they are not linearly separable in the original feature space.

In order to represent the separator as a hyperplane, the input data must first be transformed in order to determine a separator between various categories of data points. SVMs are noted for their proficiency in handling both linear and non-linear data, which makes them appropriate for a wide range of uses.

A few applications for SVMs include handwriting recognition, intrusion detection, face identification, email

categorization, gene classification, and web page classification. SVMs are an effective classification are a few uses for SVMs. SVMs are a useful tool in machine learning due to their adaptability and capacity for both classification and regression problems.

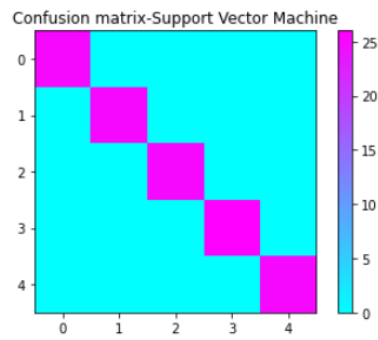


Fig 2: Support Vector Machine

Adaboost Classifier:

A meta-estimator called the AdaBoost classifier is used to improve the efficiency of machine learning algorithms. AdaBoost is especially helpful for weak learners, which are models that achieve accuracy slightly above random chance on a classification issue. In order for it to function, a classifier must first be fitted on the original dataset, and then other classifiers must be fitted on the same dataset. The AdaBoost algorithm most frequently employed is a one level decision tree.

The AdaBoost algorithm Is based on the Idea that learners develop in stages. Every consecutive learner, with the exception of the first, is grown from a prior learner, progressively transforming weak learners into powerful ones. Each weak learner concentrates on a certain aspect of the data, which aids in lowering the classification error.

AdaBoost modifies the weights of the incorrectly classified samples during the training phase such that the subsequent learner is compelled to pay more attention to the incorrectly classified data. This is done again until either a predetermined number of learners have been trained or the accuracy reaches a plateau. Finally, each weak learner's scores are merged to create a single strong learner, who outperforms each of the weak learners individually.

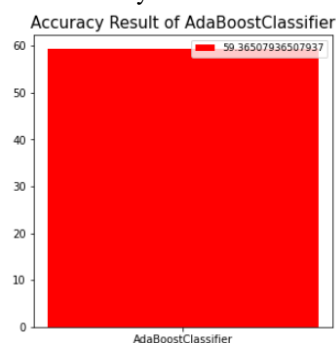


Fig 3: Adaboost Classifier

In conclusion, AdaBoost is a useful method for enhancing the performance of weak learners by gradually transforming them into strong learners by modifying the weights of the misclassified samples. The approach can be used to solve a variety of machine learning issues and works best when combined with decision trees that have just one level.

Random Forest Classifier:

For supervised learning, Random Forest is a potent and popular machine learning method. This approach is flexible and can be applied to machine learning tasks involving classification and regression. Ensemble learning, which mixes various decision trees to boost the accuracy of the model's predictions, is the basis of Random Forest.

The Random Forest approach is used to create a massive number of decision trees, and each tree is trained using a unique random subset of the original dataset. The programmed then combines all of the decision trees' predictions to provide a final forecast. The risk of overfitting the data, a typical problem in machine learning, is reduced by using this strategy of creating several decision trees.

The capability of Random Forest to handle high-dimensional datasets with numerous attributes is one of its main advantages. Furthermore, Random Forest can handle missing data points, which is a problem that frequently arises in real-world datasets. The Random Forest algorithm's accuracy is closely correlated with the number of trees in the forest. The likelihood of overfitting is reduced as the model's accuracy increases along with the number of trees in the forest.

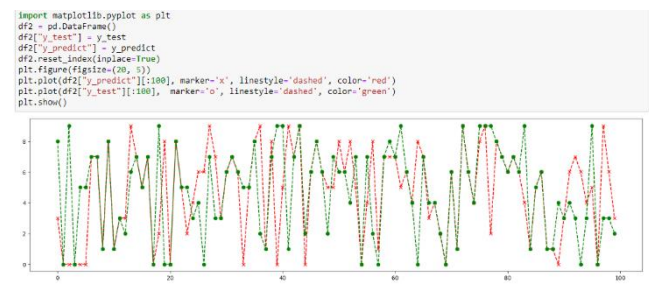


Fig 4: Random Forest Classifier

To sum up, Random Forest is a very effective machine learning technique that can be used for many different tasks, including text analysis, predictive modeling, and image categorization. As a result of its ability to handle high-dimensional datasets, missing data points, and the avoidance of overfitting, it is a favored choice for data scientists and machine learning practitioners.

Voting Classifier:

A voting classifier is a sort of machine literacy model that combines several classifiers to predict the class based on the class for which they have the highest probability. It is a form of ensemble literacy where a difficult issue is divided into manageable parts to enhance the performance of the model.

Using the combined maturity of these models' voting for each affair class, a single model is created that trains on the affair of these models and anticipates the affair. as opposed to

training separate models and changing the outcome for each. Instead of relying solely on one model, this can produce predictions that are more precise.

A voting classifier distinguishes between hard voting and soft voting. The prognosticated affair class in hard voting is the class with the highest voting maturity, or the class most likely to be anticipated by each of the classifiers. For instance, if three classifiers (A, A, and B) predicted the affair class and the maturity projected A as the affair, then A would be the final vaticination.

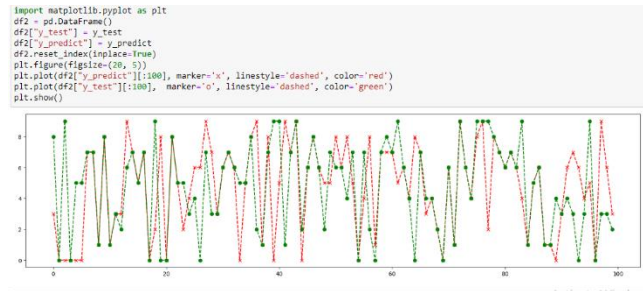


Fig 5: Voting Classifier

The affair class is the voting group that is based on the average likelihood provided to that group in soft voting. For example, the average likely for class A to occur is 0.4333, and the average likelihood for class B to occur is 0.3067 if three models predict that class A will occur with a probability of (0.30,0.47,0.53) and class B will occur with a probability of (0.20,0.32,0.40). Class A came out on top because it was most likely to be matched by all classifiers.

Voting Classifiers are frequently employed in machine literacy for tasks such as sentiment analysis, face recognition, and image bracketing, among others. They provide a means of reducing prognostications' fragility while avoiding the drawbacks of utilizing a single model.

VI. CONCLUSION

In conclusion, a crucial component of cybersecurity is intrusion discovery vaccination employing data wise methods. The identification and assistance of implicit security breaches, network incursions, and other anomalous conditioning that may constitute a serious threat to an organization's vital resources and data can be achieved through anomaly detection employing machine literacy models. The colorful approach of the proposed methodology includes data collection, preprocessing, exploratory data analysis, point selection, model selection, training, assessment, optimization, and deployment. The successful implementation of the suggested model can aid in the creation of trustworthy and efficient intrusion discovery systems that can assist associations in protecting their resources from hidden dangers. Associations can stay ahead of the changing problem geography and guarantee the security of their networks and systems by continuously monitoring and updating the model. The logical steps in the procedure were data preparation and cleaning, missing value analysis, exploratory analysis, and lastly model structure and evaluation. The fashionable delicacy will be revealed by the

sophisticated delicacy score algorithm's public test set. The inventive one is used in the operation that can help determine the type of invasions.

VII. REFERENCES

- [1] Tchakoucht TA, Ezziyyani M. "Building a fast intrusion detection system for high-speed-networks: probe and DoS attacks detection",2018.
- [2] Zuech R, Khoshgoftaar TM, "Wald R. Intrusion detection and big heterogeneous data",2015.
- [3] Debar H. "An introduction to intrusion-detection systems. In: Proceedings of Connect",2000.
- [4] Ferhat K, Sevcan A. "Big Data: controlling fraud by using machine learning libraries on Spark. Int J Appl Math Electron Comput",2018.
- [5] Manzoor MA, Morgan Y. "Real-time support vector machine based network intrusion detection system using Apache Storm",2016.
- [6] Vimalkumar K, Radhika N. "A big data framework for intrusion detection in smart grids using Apache Spark.",2017.
- [7] Dahiya P, Srivastava DK. "Network intrusion detection in big dataset using Spark"2018.
- [8] Vapnik, "The Nature of Statistical Learning Theory", Springer-Verlag, New York, 1995.
- [9] Bremner D, Demaine E, Erickson J, Iacono J, Langerman S, Morin P, Toussaint G (2005). "Output-sensitive algorithms for computing nearest-neighbor decision boundaries".
- [10] Domingos, Pedro & Michael Pazzani (1997) "On the optimality of the simple Bayesian classifier under zero-one loss".