

AN EFFICIENT SPAM DETECTION TECHNIQUE FOR IOT DEVICES USING MACHINE LEARNING

ABSTRACT

The Internet of Things (IoT) is a group of millions of devices having sensors and actuators linked over wired or wireless channel for data transmission. IoT has grown rapidly over the past decade with more than 25 billion devices are expected to be connected by 2020. The volume of data released from these devices will increase many-fold in the years to come. In addition to an increased volume, the IoT devices produces a large amount of data with a number of different modalities having varying data quality defined by its speed in terms of time and position dependency. In such an environment, machine learning algorithms can play an important role in ensuring security and authorization based on biotechnology, anomalous detection to improve the usability and security of IoT systems. On the other hand, attackers often view learning algorithms to exploit the vulnerabilities in smart IoT-based systems. Motivated from these, in this paper, we propose the security of the IoT devices by detecting spam using machine learning. To achieve this objective, Spam Detection in IoT using Machine Learning framework is proposed. In this framework, five machine learning models are evaluated using various metrics with a large collection of inputs features sets. Each model computes a spam score by considering the refined input features. This score depicts the trustworthiness of IoT device under various parameters. REFIT Smart Home dataset is used for the validation of proposed technique. The results obtained proves the effectiveness of the proposed scheme in comparison to the other existing schemes.

CONTEXT

Chapter	list of content	Page no.
	ABSTRACT	i
	LIST OF FIGURES	ii
1	INTRODUCTION	1
1.1	PROPOSED SYSTEM	1
2	LITERATURE SURVEY	19
3	SYSTEM REQUIREMENTS	24
3.1	HARDWARE REQUIREMENTS	24
3.2	SOFTWARE REQUIREMENTS	24
3.3	LANGUAGE SPECIFICATIONS	24
3.4	HISTORY OF PYTHON	25
3.5	APPLICATION OF PYTHON	25
3.6	FEATURES OF PYTHON	26
3.7	FEASIBILITY STUDY	26
3.7.1	ECONOMICAL FEASIBILITY	27
3.7.2	TECHNICAL FEASIBILITY	27
3.7.3	SOCIAL FEASIBILITY	28

4	SYSTEM ANALYSIS	29
4.1	PURPOSE	29
4.2	SCOPE	29
4.3	EXISTING SYSTEM	29
4.4	PROPOSED SYETEM	30
5	SYSTEM DESIGN	32
5.1	INPUT DESIGN	32
5.2	OUTPUT DESIGN	32
5.3	DATA FLOW DIAGRAM	33
6	MODULES	39
6.1	LOGIN MODULE	39
6.2	DATA COLLECTION MODULE	39
6.3	PRE-PROCESSING MODULE	40
6.4	TRAIN AND TEST MODULE	40
6.5	DETECTION OF SPAM	41
7	SYSTEM IMPLEMENTATION	42
7.1	SYSTEM ARCHITECTURE	42
8	SYSTEM TESTING	43
8.1	TEST PLAN	43
8.2	VERIFICATION	43

8.3	VALIDATION	43
8.4	BASIC OF SOFTWARE TESTING	43
8.5	BLACK BOX TESTING	44
8.6	WHITE BOX TESTING	44
8.7	TYPES OF TESTING	44
8.7.1	UNIT TESTING	45
8.7.2	INTEGRATION TESTING	45
8.7.3	FUNCTIONAL TESTING	45
8.7.4	SYSTEM TESTING	45
8.7.5	STRESS TESTING	46
8.7.6	PERFORMANCE TESTING	46
8.7.7	USABILITY TESTING	46
8.7.8	ACCEPTANCE TESTING	46
8.7.9	REGRESSION TESTING	47
9	RESULT AND DISCUSSION	51
10	CONCLUSION	61
11	REFERENCE	62

LIST OF FIGURES:

FIGURE NO. NO.	FIGURE NAME	PAGE
1.1	EXPLAINING THE WORKING ALGORITHM OF THE RANDOM FOREST ALGORITHM	3
1.2	EXPLAINING THE RANDOM FOREST CLASSIFIER	5
1.3	EXPLAINING THE RANDOM FOREST CLASSIFIER CLASSIFIERALGORITHM WITH EXAMPLE	6
7.1	SYSTEM ARCHITECTURE	42
9.1	HOME PAGE	51
9.2	LOGIN PAGE	52
9.3	UPLOAD PAGE	54
9.4	PREVIEW PAGE	55
9.5	TRAINED DATA	56
9.6	PREDICTION NO SPAM	57
9.7	PREDICTION SPAM	58
9.8	PERFORMANCE ANALYSIS	59
9.9	CONFUSION MATRIX	59
9.10	CHART ANALYSIS	60

CHAPTER 1

INTRODUCTION

The safety measures of IoT devices depends upon the size and type of organization in which it is imposed. The behavior of users forces the security gateways to cooperate. In other words, we can say that the location, nature, application of IoT devices decides the security measures. For instance, the smart IoT security cameras in the smart organization can capture the different parameters for analysis and intelligent decision making. The maximum care to be taken is with web based devices as maximum number of IoT devices are web dependent. It is common at the workplace that the IoT devices installed in an organization can be used to implement security and privacy features efficiently. For example, wearable devices collect and send user's health data to a connected smartphone should prevent leakage of information to ensure privacy. It has been found in the market that 25-30% of working employees connect their personal IoT devices with the organizational network. The expanding nature of IoT attracts both the audience, i.e., the users and the attackers. However, with the emergence of ML in various attacks scenarios, IoT devices choose a defensive strategy and decide the key parameters in the security protocols for trade-off between security, privacy and computation. This job is challenging as it is usually difficult for an IoT system with limited resources to estimate the current network and timely attack status.

1.1 PROPOSED ALGORITHM

RANDOM FOREST ALGORITHM

Random forest algorithm can use both for classification and the regression kind of problems. In this you are going to learn, how the **random forest algorithm** works in machine learning for the classification task.

Random Forest is a popular machine learning algorithm that belongs to the supervised learning technique. It can be used for both Classification and Regression problems in ML. It is based on the concept of **ensemble learning**, which is a process of combining multiple classifiers to solve a complex problem and to improve the performance of the model.

A random forest algorithm consists of many decision trees. The ‘forest’ generated by the random forest algorithm is trained through bagging or bootstrap aggregating. Bagging is an ensemble meta-algorithm that improves the accuracy of machine learning algorithms.

As the name suggests, "Random Forest is a classifier that contains a number of decision trees on various subsets of the given dataset and takes the average to improve the predictive accuracy of that dataset." Instead of relying on one decision tree, the random forest takes the prediction from each tree and based on the majority votes of predictions, and it predicts the final output.

The below diagram explains the working of the Random Forest algorithm:

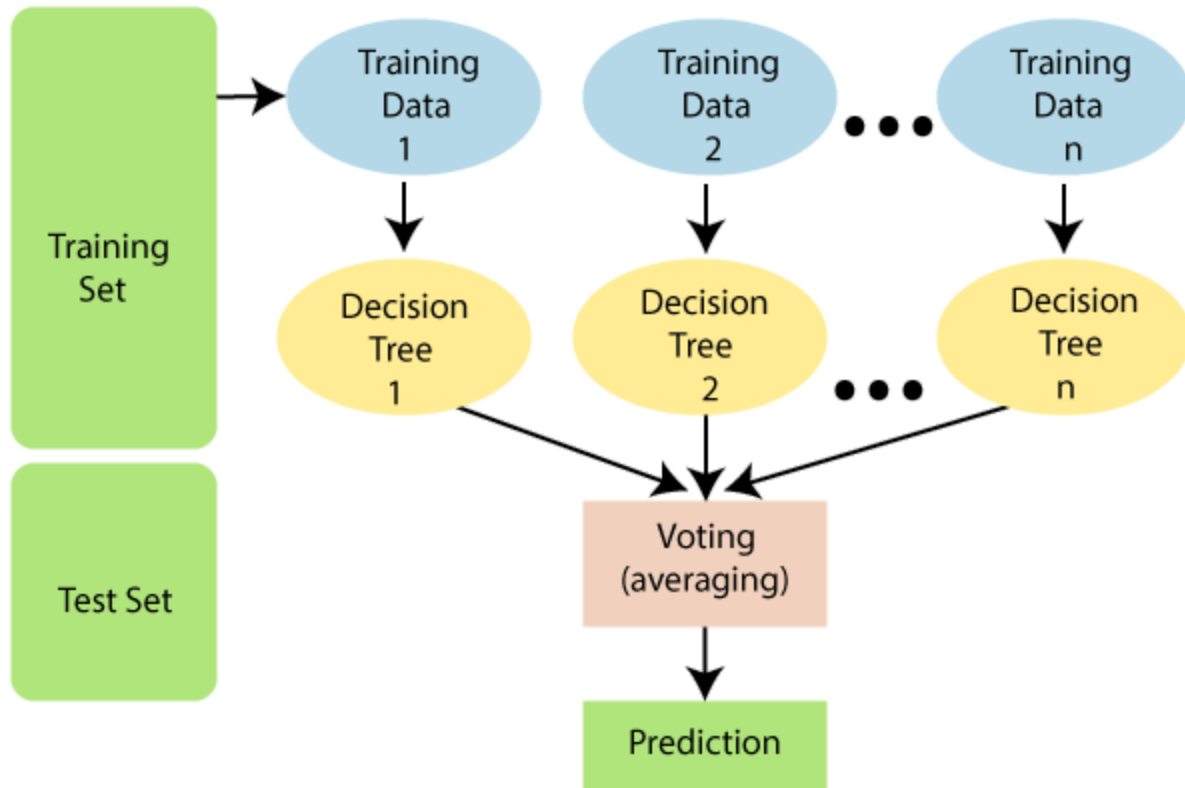


Fig 1.1: Explaining the working algorithm of the Random Forest algorithm

Below are some points that explain why we should use the Random Forest algorithm:

- o It takes less training time as compared to other algorithms.
- o It predicts output with high accuracy, even for the large dataset it runs efficiently.
- o It can also maintain accuracy when a large proportion of data is missing.

Features of a Random Forest Algorithm

- It's more accurate than the decision tree algorithm.
- It provides an effective way of handling missing data.

- It can produce a reasonable prediction without hyper-parameter tuning.
- It solves the issue of overfitting in decision trees.
- In every random forest tree, a subset of features is selected randomly at the node's splitting point.

Classification in random forests

Classification in random forests employs an ensemble methodology to attain the outcome. The training data is fed to train various decision trees. This dataset consists of observations and features that will be selected randomly during the splitting of nodes.

A rain forest system relies on various decision trees. Every decision tree consists of decision nodes, leaf nodes, and a root node. The leaf node of each tree is the final output produced by that specific decision tree. The selection of the final output follows the majority-voting system. In this case, the output chosen by the majority of the decision trees becomes the final output of the rain forest system. The diagram below shows a simple random forest classifier.

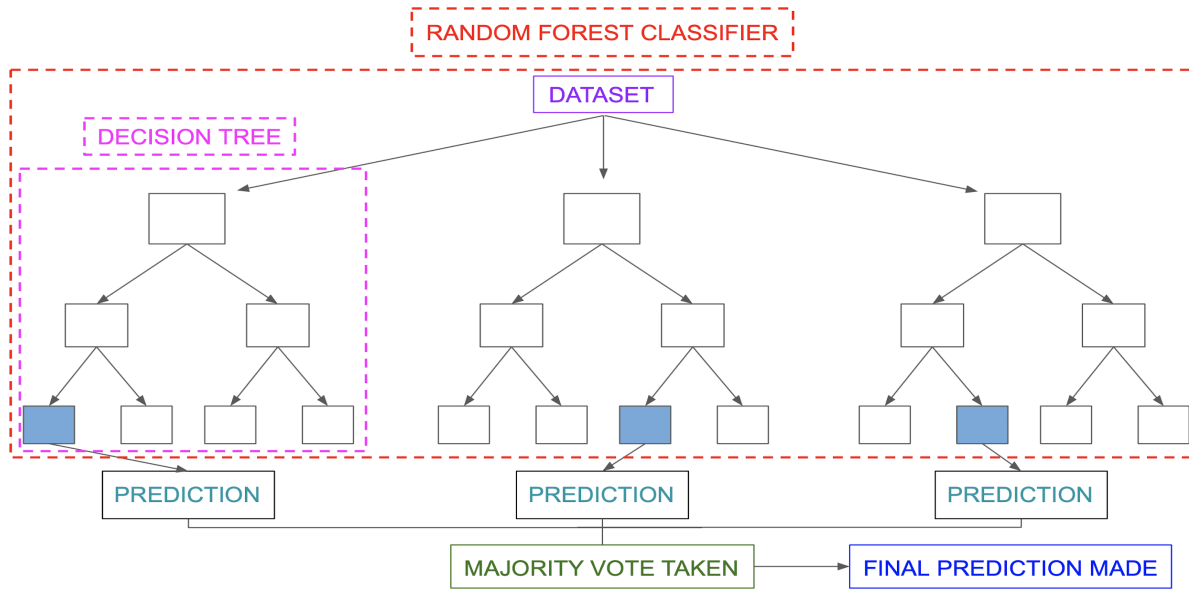


Fig 1.2: Explaining the Random Forest Classifier

Random Forest Steps

1. Randomly select “k” features from total “m” features.
 1. Where $k \ll m$
2. Among the “k” features, calculate the node “d” using the best split point.
3. Split the node into daughter nodes using the best split.
4. Repeat 1 to 3 steps until “l” number of nodes has been reached.
5. Build forest by repeating steps 1 to 4 for “n” number times to create “n” number of trees.

The beginning of random forest algorithm starts with randomly selecting “k” features out of total “m” features. In the image, you can observe that we are randomly taking features and observations.

The working of the algorithm can be better understood by the below example:

Example: Suppose there is a dataset that contains multiple fruit images. So, this dataset is given to the Random forest classifier. The dataset is divided into subsets and given to each decision tree. During the training phase, each decision tree produces a prediction result, and when a new data point occurs, then based on the majority of results, the Random Forest classifier predicts the final decision. Consider the below image:

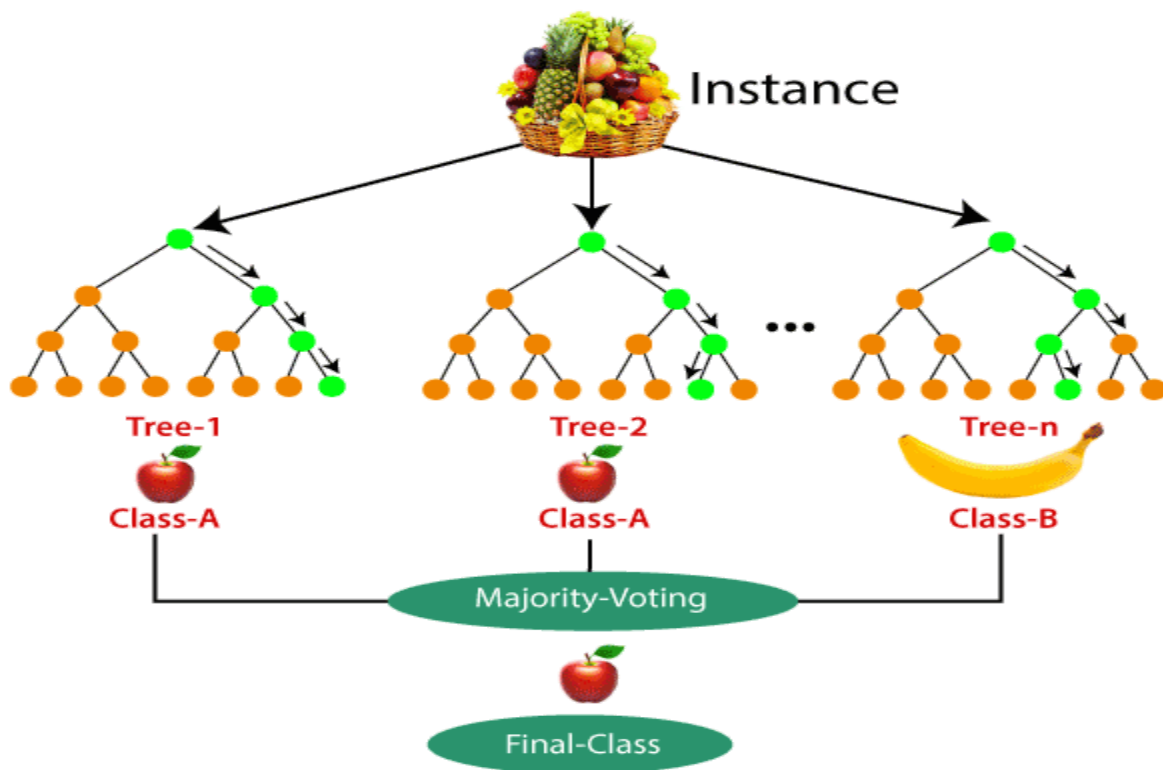


Fig 1.3: Explaining the Random Forest Classifier algorithm with example