

# NETWORK SECURITY

# Project

---

KOCKAN Melisa / CANIN Raphaël / BIBILA Jodry / BACHA Eliès  
E4 AIC

## Introduction

The purpose of this lab is to use machine learning approaches in order to detect intrusion on the use case of IoT.

To do so, the dataset we are going to use is a wireshark capture which depicts the traffic of a real network receiving attacks.

The dataset at stake is available using this link : [click here](#). In this dataset, each line corresponds to an element from the traffic received.

After identifying the different types of attacks in our dataset, we will analyze how they can impact an IoT system in order to select the relevant features enabling us to predict if a line in the traffic corresponds to an attack.

## Identification of the attacks

After exploring our dataset on our ipynb file, we identified 4 different types of attacks :

- Mirai
- DoS
- Scan
- Spoofing

According to the definition of Oracle, the Internet of Things (IoT) refers to « the network of physical objects—“things”—that are embedded with sensors, software, and other technologies for the purpose of connecting and exchanging data with other devices and systems over the internet ».

To understand how these 4 attacks can affect an IoT system, let's define each of them beforehand.

- Mirai  
Mirai is a malware looking for IoT systems using default usernames and passwords in order to infect them
- DoS  
Fact of sending indefinitely requests to a service in order to make it unavailable
- Scan  
Analyzing ports to establish a list of open ports to attack

- Spoofing

Attackers disguise into well known and safe company in order to retrieve informations

Mirai does not only make IoT devices unavailable ; it leads to the creation of BotNet, which can lead to launching a DoS attack. By making an IoT device unavailable, DoS will indirectly make unavailable any other IoT devices depending on them. Scan permits to pave a path to spoofing

## Using a machine learning model

Using python, We train a machine learning model upon these data and evaluate its performances. We used several algorithms to compare them and choose the one with the highest performance

- KNN (93.7% accuracy)
- Logistic regression (79.1%)
- Linear discriminant analysis (76.8%)
- Support vector machine (79.3%)
- Trees (96.6%)
- Random forest (94.2%)
- A neural network for classification (84.0%)
- Gradient Boosting (97.2%)
- **Stacking classifier (97.6%)**

Please find our *Project\_Network\_security* notebook to see the result.

### ❖ Explanation of the different steps

First, after importing the dataset, we need to take a very important step for the results :

#### ➤ **PREPROCESSING**

Preprocessing is an important step in Machine Learning because it strongly impacts the results of our model.

- In datasets, we rarely have a clean dataset. The first step in preprocessing is often to deal with null values. Indeed, these zero/nan values do not allow us to train our model correctly.

- In this dataset we have different types of data. However, to train a model we need only numeric features. So, we have to search for the different categorical features that we have and convert them into numeric features. The categorical data is identified with the type Object. We therefore isolate this type of data. The LabelEncoder method allows each categorical data to be coded into numerical labels. Then, we convert each integer into one-hot encoded (i.e. one binary number) hence the method name one-hot encoded. In this dataset, the only categorical data is the column CLASS.
- To learn more about the dataset, we will also use another tool: the correlation matrix. This tool allows you to find the possible dependencies between the variables of the dataset.

#### ➤ TRAINING AND TESTING DIFFERENT MODELS

First, we divide the dataset into two parts: one for the training and the other for the test.

We implement a function that calculates the average precision over n different training sessions for a given model as parameter.

Then, we calculate the accuracy for different algorithms on normal data and normalized data.

## Relevance of each features

Sometimes, the columns of a dataset aren't all relevant to take into account for our analysis. Before studying each column individually, we printed the matrix correlation of our data set.

On this matrix correlation, we had no surprise noticing that some columns such as 'Forward\_Packets\_Length\_Max', 'Forward\_Packets\_Length\_Min' and 'Forward\_Packets\_Length\_Mean' were linked together.

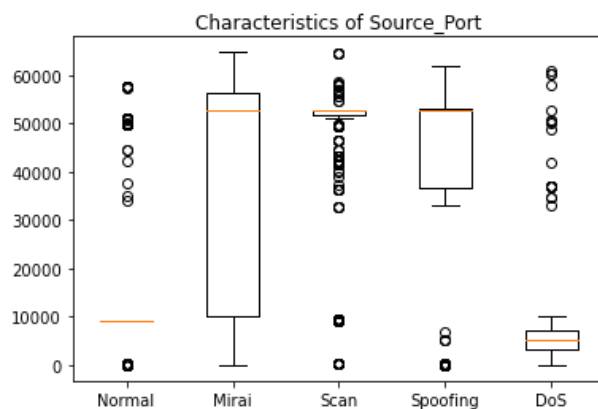
To take a further look at each column, we draw a box plot for each of them owing to a function we implemented during our Machine Learning project.

In the 32 columns of our dataset , we highlighted in green those who seemed the most meaningful to distinguish our attacks :

'Source\_Port', 'Destination\_Port', 'Protocol', 'Flow\_Duration',  
'Total\_Forward\_Packets', 'Total\_Backward\_Packets',  
'Total\_Length\_Forward\_Packets', 'Total\_Length\_Backward\_Packets',  
'Forward\_Packets\_Length\_Max', 'Forward\_Packets\_Length\_Min',  
'Forward\_Packets\_Length\_Mean', 'Backward\_Packets\_Length\_Max',  
'Backward\_Packets\_Length\_Min', 'Backward\_Packets\_Length\_Mean',  
'Packets\_Flow\_IAT\_Mean', 'Packets\_Flow\_IAT\_Standard\_Deviation',  
'Packets\_Flow\_IAT\_Max', 'Packets\_Flow\_IAT\_Min', 'Forward\_IAT\_Total',  
'Forward\_IAT\_Mean', 'Forward\_IAT\_Min', 'Backward\_IAT\_Total',  
'Backward\_IAT\_Mean', 'Backward\_IAT\_Min', 'Forward\_Header\_Length',  
'Backward\_Header\_Length', 'Forward\_Packets\_per\_Second',  
'Backward\_Packets\_per\_Second', 'Packets\_Length\_Min',  
'Packets\_Length\_Max', 'Packets\_Length\_Mean',  
'Packets\_Length\_Standard\_Deviation',

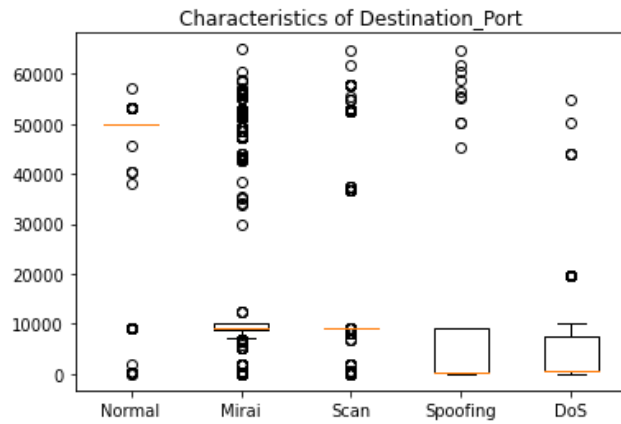
Let's take a further look into these most impactful columns individually.

### 'Source Port'



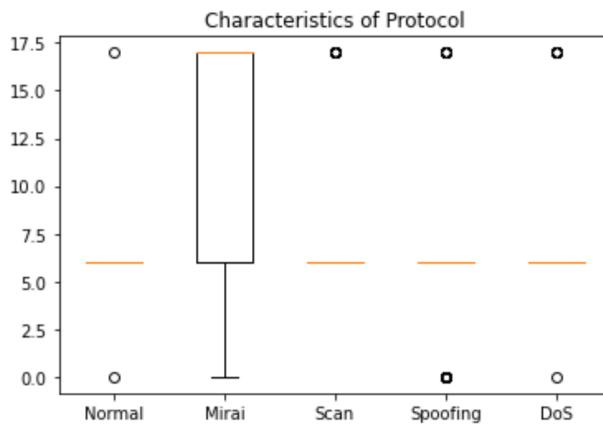
This column is one of the more useful to detect each attack. As we can see, each type of attack emits from a pretty different segment of ports.

## 'Destination Port'



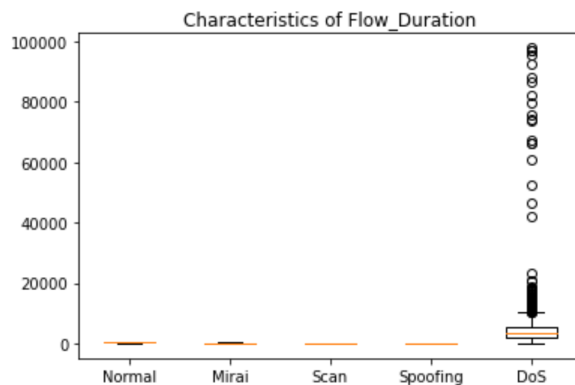
Like the previous column, this one has a strong influence on the accuracy of predictions.

## 'Protocol'



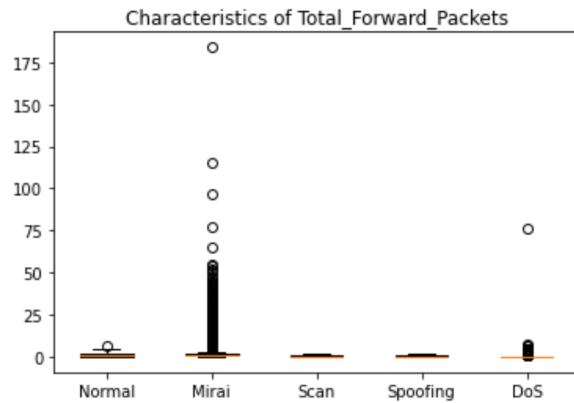
This parameter is relevant to detect a Mirai attack. Mirai is a malware, compared to the three other attacks ; hence, it seems coherent to say that the protocol used by this attack is suspicious.

## 'Flow\_Duration'



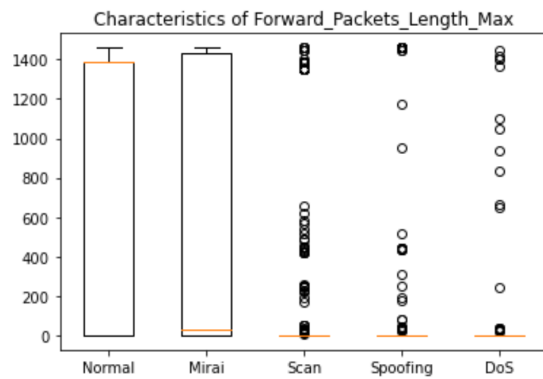
This parameter is characteristic of DoS. During DoS, requests are sent succinctly and indefinitely ; hence, the flow duration is drastically affected.

## 'Total\_Forward\_Packets'



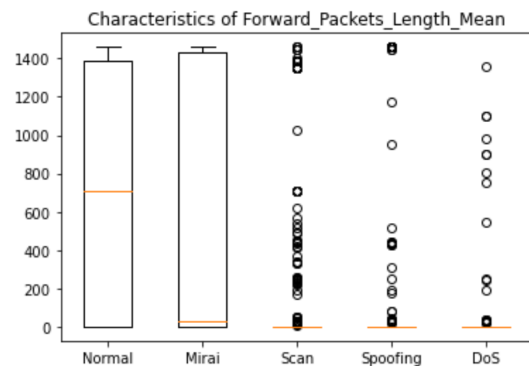
This parameter is relevant to detect Mirai attacks.

## 'Forward\_Packets\_Length\_Max'



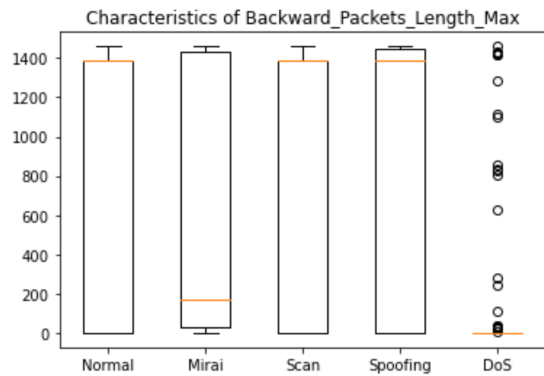
This parameter enables us to detect the fourth of our attacks.

## 'Forward\_Packets\_Length\_Mean'



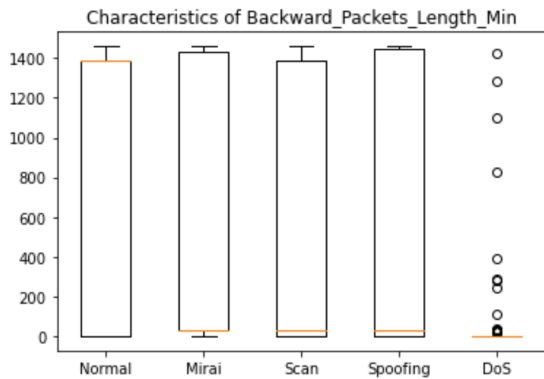
This parameter enables us to detect the fourth of our attacks.

## 'Backward\_Packets\_Length\_Max'



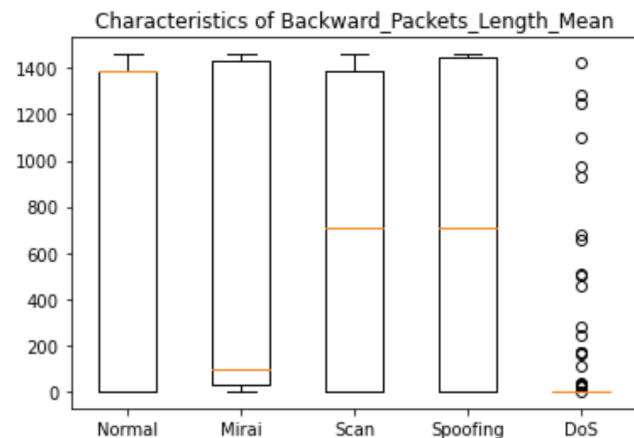
This parameter helps to detect Mirai and DoS.

## 'Backward\_Packets\_Length\_Min',



This parameter enables us to detect the fourth of our attacks.

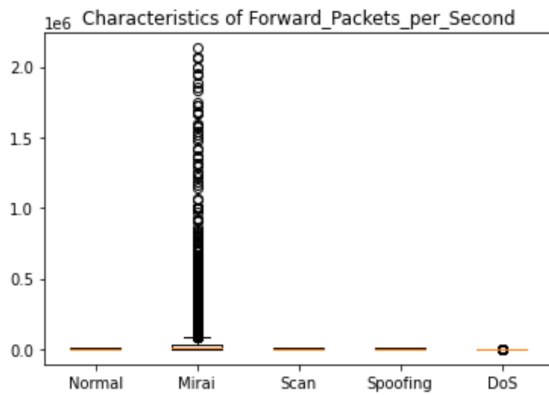
## 'Backward\_Packets\_Length\_Mean'



This parameter enables us to detect the fourth of our attacks, but especially for Mirai and DoS.

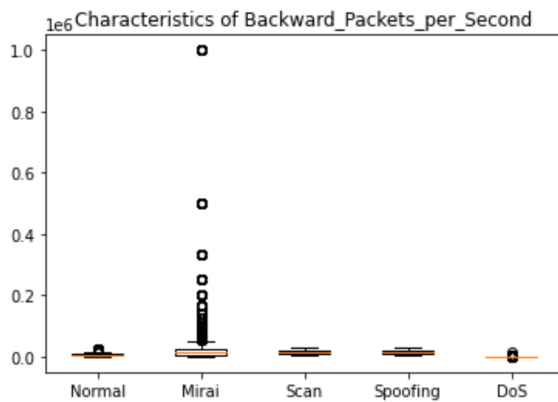


### 'Forward\_Packets\_per\_Second'



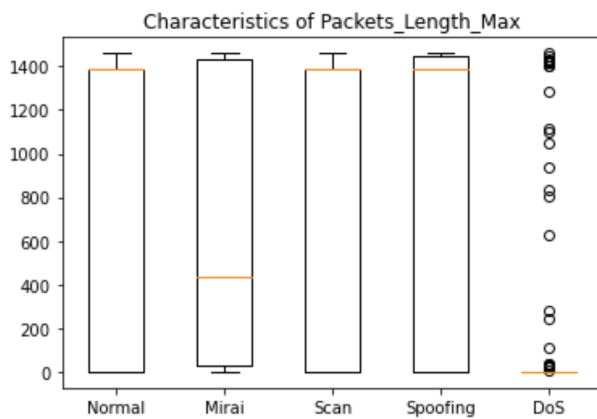
This parameter is meaningful for Mirai attacks.

### 'Backward\_Packets\_per\_Second'

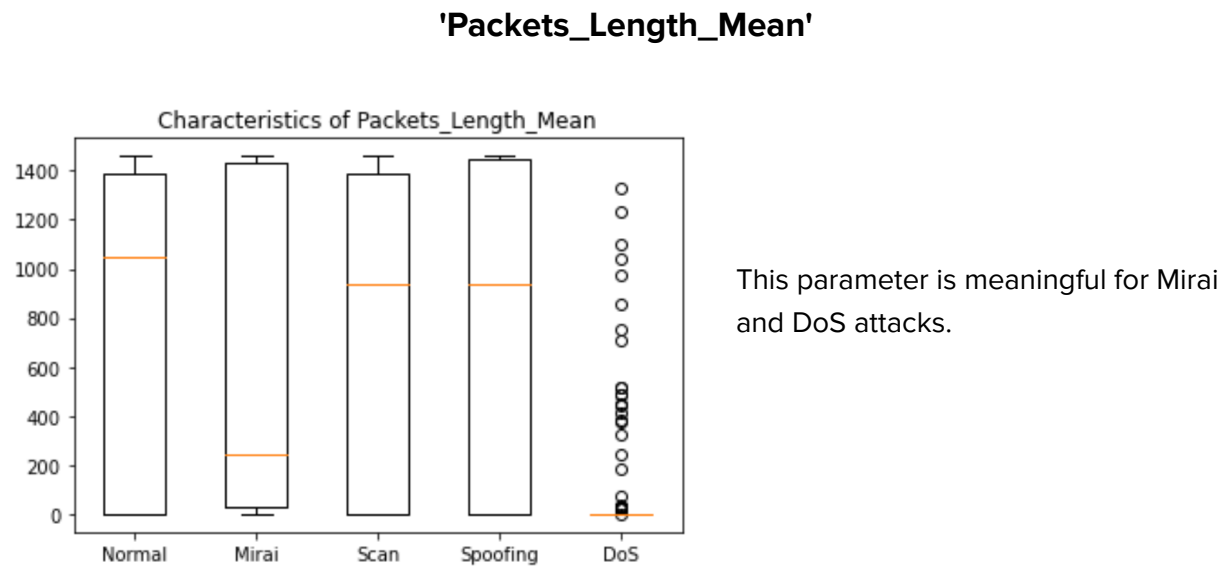


This parameter is meaningful for Mirai attacks.

### 'Packets\_Length\_Max'



This parameter is meaningful for Mirai and DoS attacks.



We reduced our dataset to these columns, and retrained our models.

Finally we obtained the following performance:

- KNN (93.7% accuracy)
- Logistic regression (77.97%)
- Linear discriminant analysis (75.0%)
- Support vector machine (79.06%)
- Trees (97.17%)
- Random forest (95.8%)
- A neural network for classification (83.36%)
- Gradient Boosting (97.27%)
- **Stacking classifier (97.73%)**

The accuracy is a little better than before reducing our dataset. We can also point out the fact that reducing the dataset enabled us to reduce the calculation time.