*Thomas PAUL*
*Yana RAGOZINA*
*SI5 IA-ID FISA*

# Anomaly Detection in IP Traffic - Data Mining Project

## Context

This project aims to detect anomalous patterns of communication in IP traffic. The principal method is to build a profile of every IP communication in the form of a small graph, also called a graphlet. Each graphlet is passed to a Support Vector Machine model to automatically detect normal and anomalous communications.

## Results

To detect malicious communications, we have first attempted to do it without any data transformation in the data annotated-trace.txt. When testing the learnt SVM model on our non annotated dataset, we realized that the model was not able to detect any anomalous communications. This is a normal behaviour. Indeed, the class distribution of normal and anomalous communications is very unbalanced. We have around 99.3% of normal communications compared to 0.7% of anomalous communications. Yet, the accuracy of our model is very good (around 99.5%) and this is once again, a normal behaviour. Our model is right 99% of the time by predicting a normal communication, which is very good, but it is not the goal of the project.

To solve this problem, we have proceeded with data augmentation of anomalous classes. The oversampling method provided by the imblearn library will help us to do so. We have augmented the percentage of anomalous communications from 0.7% to 16.7%. While our accuracy has decreased from 99.5% to 83% for both dimensional space transformation and kernel trick methods, the model was able to detect some anomalous communications (around 0.1%).

Because we do not have the annotations of the communications detected by our SVM model, we are only able to make suggestions. The flow number 1628 is a potential true positive because there are many times the ports 53 and 80 in our malicious annotated communication traces. This could correspond to a DDoS attack, exploited as a DNS amplification attack (communication between ports 53 and 80).

*Thomas PAUL*
*Yana RAGOZINA*
*SI5 IA-ID FISA*

The flow number 7299 is a potential true positive as well as it has communications from ports 25 to 20, via protocol 17, which is not usually a standard behaviour in our annotated dataset. In fact, attackers sometimes use port 25 to conduct distributed denial-of-service (DDoS) attacks against web servers. We have cases of anomaly flows for this type of communication. This indicates that potential false negatives can have the same patterns.

Plus, there are many cases where the communications between the ports 25, 53 or 80 are annotated as anomalies, which corresponds to many cases of DDoS attacks in reality.

We could suggest that the remaining detected anomalies are false positives (i.e. normal communications). In fact, the communications between the detected as abnormal ports are always stated as normal.

Overall, our model relies a lot on data augmentation to detect anomalous communications, this is why we need to find the correct percentage of anomalous data to augment.