

INTRODUCTION

The need for efficient Intrusion Detection System (IDS) is a modern problem that aims to keep communication networks operational and secure. The present study focuses on improving IDS by using Machine Learning (ML) that can be trained with limited data. Thus, it intends to remove the burden of gathering large dataset that are difficult to create in this field.

A Neural Network (NN) aims to imitate the human brain. To classify any type of objects, some patterns need to be recognized such as a shape, a colour, a number of occurrences... When classifying network traffic, such features can be the number of packets exchanged, their length or the flags used.

METHODOLOGY

- [1]:** The objective is to **extract features** from files (pcap format) that contain raw network captures. The **38 extracted features** are then sorted in two different dataframes, one in which the labels are either malicious or benign traffic, the other specifying 12 different types of attacks.
- [2]:** First dataset is used for **binary classification** which is a simple way to start working with NN:
- ⇒ It takes 38 inputs on first layer.
 - ⇒ It goes through 2 hidden layers (the size of which can be optimized).
 - ⇒ The final layer has 1 neuron which returns the probability of the sample to belong in one of the two classes.
- [3]:** Evolution of [2] can predict the type of traffic between **13 different classes**
- ⇒ Major change is the size of the output layer that now contains 13 neurons.

- [4]: Data limitation** is integrated in the project which causes a drop in the performances of the previous models. **Siamese Network** is a solution that aims to override the problems caused by limited training data.
- ⇒ Instead of predicting the class of one sample, the network takes two different samples and gives it to two exact same Neural Networks.
 - ⇒ It computes the Euclidean distance between the two output vectors.
 - ⇒ It thresholds the result value and predicts if the two samples belong to the same class.
- Siamese Network is powerful with limited data because it can create a large number of pairs from few samples. An evolution of this Siamese Network then integrates a **K-Nearest Neighbours Algorithm** during the testing process. Indeed, to be classified, a sample is not just compared to random other samples but to its nearest neighbours which are more likely to provide accurate results.

RESULTS AND DISCUSSION

The performance of each model is given in term of accuracy, which is the percentage of samples that have been well classified in the testing set. The results also compare how each model performs on both large dataset available and with data limitation.

	Binary NN	Multi-class NN	Siamese NN
Large dataset	96%	85%	-
Small dataset	80%	68%	75%

Figure 1: Evaluation of different models performances depending on the size of data available

While binary classification performs well, the aim is to get better results for many types of attack. The limits of feed forward NN are reach fast as soon as the dataset becomes smaller than 100 samples per class. However, figure 2 shows that Siamese NN can still be efficient down to 25 samples per class.

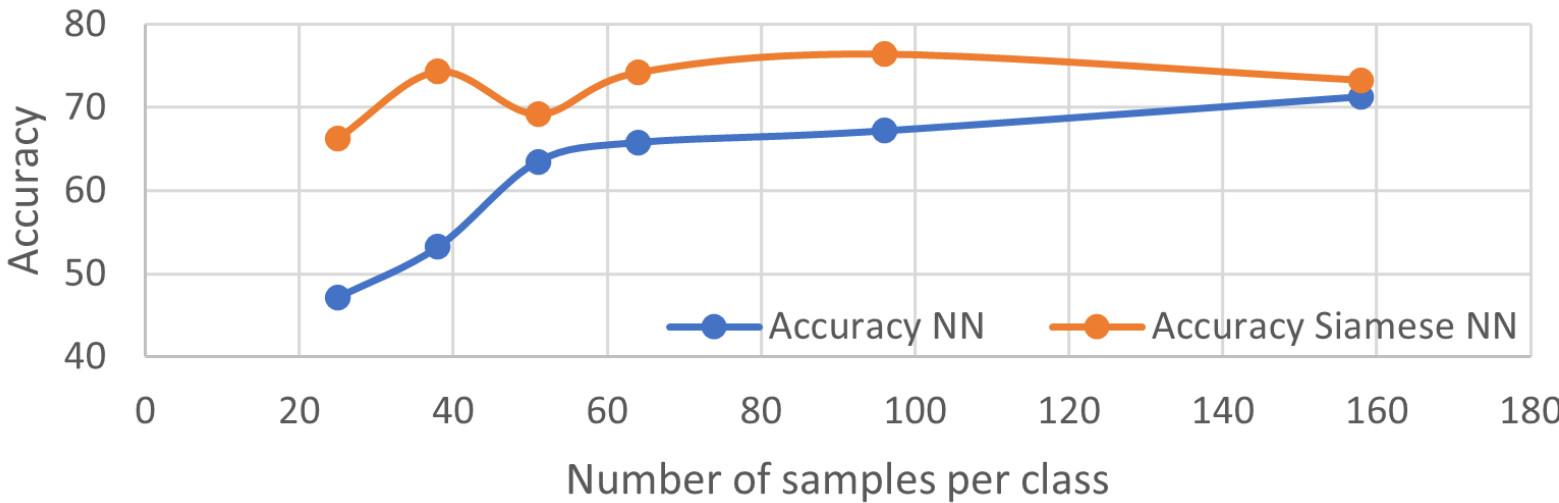


Figure 2: Evolution of accuracy depending on the dataset size

CONCLUSION AND FURTHER WORK

The problem has been approached step by step, from binary classification with large dataset to multi-class classification with few samples available.

However, even if Siamese network can be more efficient than feed forward NN when dealing with limited data, the accuracy did not exceed 75% which is a major limitation to the implementation of an Intrusion Detection System. Further research on NN architectures and on features selection should lead to improvement in the classifier performances.