

# **Deep Real-time GDPR Compliance & Malware Auditing**

## Objectives

The objective of this project is to manage and track the storage and access of personal data of an enterprise by looking for malware and privilege misuse.

- Identifying and understanding data regulations relating to user data and privacy (e.g. GDPR)
- Identifying and understanding enterprise level security and privacy requirements
- Building and training a deep learning neural network to predict whether a sequence of events leads to a security and/ or privacy violation
- Enterprise compliance reporting

### Introduction

Deep Learning is a growing field of study which already has many applications in the industry, for example image recognition, drug discovery and speech recognition.

For this project, the focus lies on using artificial neural networks (ANN) to classify all activity in an enterprise system as either normal or suspicious behaviour and subclasses of that. The ANNs output is the probability that a sequence of actions is malicious and in violation of legal or enterprise policies. It should also classify malicious activity by attack type (e.g. user misuse, botnet).

The dataset being analysed is gathered and preprocessed by another part of the DICE project. It represents chains of events that are then classified. Examples of suspicious causality chains are:

- A regular communication initiated by an enterprise host to an external domain. (potential botnet)
- The copy/ pasting of sensitive data and sending an email to an external address in a short timeframe on the same host. (potential misuse)



#### Methods

Deep Learning (DL) was chosen as the basis of this project because it requires no manual selection to parse a dataset and extract relevant features as opposed to Machine Learning (ML). This is called the training phase. The ANN parses a large dataset of labelled data (supervised learning), extracts relevant features and calculates the weighting of the edges between the neural nodes.

Once that is done, the ANN is initialized and can be used to classify any data in the same format as the training data based on its observations during training. The network can be retrained with new data at any time to optimize its performance or to apply it to a different problem.

Performance depends on the selected structure of the ANN and the training algorithm, each of which has advantages and disadvantages.

Julia Meister

Supervised by: Konstantinos Markantonakis and Raja Naeem Akram Information Security Group, Smart Card and IoT Security Centre

### Challenges

One of the challenges of this project is to find the right balance between the True Positive Rate (TPR) and the False Negative Rate (FNR). The parameters that configure the ANN should be enterprise specific and easy to modify to offer greater functionality.

It is also important that the implementation is efficient to allow for analysis of data in real-time. This challenge can be further broken down into these sections:

- Choosing a suitable ANN structure and identifying the optimal training algorithm. These affect both the network's efficiency in training and classification and prediction accuracy.
- Potential dimension reduction of the dataset to make parsing faster. This can help reduce the chance of overfitting by making exceptions specific to the dataset less influential.



The Smart Card and Internet of Things Security Centre

# Additional Information

This project is related to the EPSRC funded project "Data to Improve Customer Experience (DICE)". The project is particularly interested in personal data, and is using rail passengers as a specific focus of interest. The overall aims of the project are:

• Understand the role that personal data plays in enhancing the user experience of rail passengers

• To develop technical solutions to data privacy

• To develop an evaluation framework that can be implemented so passengers can understand how their data is used and how they can control and verify its use.

The project started in October 2016, and runs for three years to September 2019. For more information the about please project, visit <u>http://www.dice-project.org</u>.

#### Acknowledgements

We acknowledge the support of the ISG-SCC for the summer internship program and EPSRC funded project. The views and opinions expressed in this poster are those of the authors and do not necessarily reflect the position of DICE project or any of partners associated with this project.

#### Contact

Web: https://scc.rhul.ac.uk/

Email: Julia.Meister.2016@live.rhul.ac.uk