

BABEȘ BOLYAI UNIVERSITY, CLUJ NAPOCA, ROMÂNIA  
FACULTY OF MATHEMATICS AND COMPUTER SCIENCE  
SPECIALIZATION COMPUTER SCIENCE IN GERMAN

# Cloud based malicious PDF Detection using Machine Learning

– Diploma thesis –

**Author**  
Viorel GURDIS

2020

## **Abstract**

# Contents

## List of Tables

## List of Figures

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Context . . . . .	1
1.2	Motivation . . . . .	1
1.3	Paper structure and original contributions . . . . .	2
<b>2</b>	<b>Related work</b>	<b>4</b>
2.1	Cloud based malware detection . . . . .	4
2.2	Detection of malicious PDF . . . . .	5
<b>3</b>	<b>Scientific Problem</b>	<b>7</b>
3.1	Problem definition . . . . .	7
3.2	Background processes in Microsoft Windows . . . . .	8
3.3	Filesystem monitoring . . . . .	8
3.4	Analyzing PDF File Structure . . . . .	8
3.5	Malware in PDF . . . . .	8
3.6	Machine Learning for Malware Detection . . . . .	8
3.7	Benefits of Cloud Computing . . . . .	8
<b>4</b>	<b>Proposed approach</b>	<b>9</b>
4.1	Dataset . . . . .	9
4.2	Proof of Concept . . . . .	9
4.2.1	Feature Selection . . . . .	9
4.2.2	Classification Techniques . . . . .	9
4.2.3	Performance Evaluation . . . . .	9
4.2.4	Experiment . . . . .	9
4.3	Used technologies . . . . .	9
4.3.1	Microsoft .NET Framework . . . . .	9
4.3.2	PDF Tools and Metasploit Framework . . . . .	9
4.3.3	Python Flask . . . . .	9
4.3.4	Scikit-learn Machine Learning Library . . . . .	9
4.3.5	ReactJS Framework . . . . .	10

---

<b>5</b>	<b>Application</b>	<b>11</b>
5.1	Design . . . . .	11
5.2	Windows Service . . . . .	11
5.3	Cloud API . . . . .	11
5.4	Dashboard Interface . . . . .	11
<b>6</b>	<b>Conclusion and future work</b>	<b>12</b>
	<b>Bibliography</b>	<b>13</b>

# List of Tables

# List of Figures

# List of Algorithms

# Chapter 1

## Introduction

### 1.1 Context

The informational technology progress, that is in continuous growth, brings a lot of benefits along with new responsibilities. There are plenty of applications that we use everyday across the entire World Wide Web. We don't even realize how much of our personal information is transferred to the virtual environment. That being said, it's important to be prepared for cybersecurity threats that can misuse our sensitive data. The problem is that the cyber attackers develop a lot of *hacking*<sup>1</sup> techniques, which are becoming increasingly difficult to detect. The popular software applications are the best target to inject malicious behavior. This happened with the Adobe PDF format. PDF documents are well known, trustworthy files and they are a global solution for sharing information. There are a lot of PDF readers and even browsers and email applications have support to open these files for viewing. It became so convenient to work with this format, that users interact with PDF documents without noticing any possible danger. However PDF is a often used attack vector. The large number of discovered PDF vulnerabilities and also the support of embedding Javascript code into documents are just some of the most exploited methods.

### 1.2 Motivation

Under the guise of seeming harmless, PDF documents are used on a daily basis across numerous public institutions, private companies and for personal purposes. Most of the people don't consider that these files could be dangerous and just copy the documents to their computers and access them. It is enough for one PDF to be malicious and the entire network affiliated to an institution becomes compromised. The cyber attackers succeed in achieving their goals, but

---

<sup>1</sup>attempt to gain unauthorized access to data in a system



the victim institution requires huge resources, both financially and time wise, to restore the integrity and security of their infrastructure. A measure to combat the described situation is having a real time protection installed on the computer. The most common solution, antivirus applications, work by signature matching, which is effective for detecting previously identified malware<sup>2</sup>. This means that all the antivirus applications require permanent updates in order to keep their malware databases up-to-date. Many users opt out of security solutions because of the multitude of hardware requirements they need. In this thesis we will go through a new approach that implies transferring complex detection algorithms from user's computers to a remote service, whose only task is analyzing suspicious uploaded PDF documents. Consequently the computers that will use this solution get rid of additional workload while still maintaining the system secure.

### 1.3 Paper structure and original contributions

The main contribution of this thesis is the research, design and implementation of alternative security solution oriented on multiplatform use and performance efficiency. For achieving this, the work was splitted into three important parts:

1. **Real time monitoring system on Windows**, whose goal is to notify the user when a specific file type, in our case PDF, is downloaded. The main focus when designing this software was to keep its functionality basic, so that it would require minimal computer's resources. As soon as it detects a new such file, it submits it to the Cloud for further analysis.
2. **Development of intelligent algorithm for PDF classification** implies studying the PDF format skeleton, researching popular attack vectors using PDFs, extracting the most relevant features from a document and feeding them to the most fit Machine Learning classification algorithm. A part of this effort was also spent for searching the dataset, based on which the classification model was trained to correctly detect malicious PDF files.
3. **Cloud based application for remote scanning**, which is a generic application that can integrate models such as the developed one for PDF classification and can use them for analyzing the uploaded files. It offers an user friendly dashboard<sup>3</sup> for visualizing the scanned files and also provides the possibility of submitting an URL containing a PDF file. The application will care about downloading and analyzing the file, showing the user

---

<sup>2</sup>any software intentionally designed to cause damage to a computer, server, client, or computer network

<sup>3</sup>type of graphical user interface which provides relevant informations

the scanning results. This is a perfect option for users that requires the Clouds features from a smartphone.

The remaining parts of this thesis are structured in the following manner: Chapter 2 contains the research made in the past years on malicious PDFs and Cloud Computing as an antimalware solution. Chapter 3 introduces the reader into fundamental mechanisms used for implementing the final application, as well as a brief history of malware in PDF. In Chapter 4 we present a Proof of Concept for developed classification model based on recreating a pseudo attack using malicious PDF. The 5th Chapter is putting each component together and presents the overall design of the application. Chapter 6 is meant for our final conclusions and we also show some of the future ideas for improving the application.

# Chapter 2

## Related work

Beginning with the first reported occurrences of malware, the security researchers have made a huge effort to prevent the harmful behavior. Over the years malware has evolved in different forms and of course the antivirus industry has developed more complex detection solutions. Since PDF documents became a good target for cybercriminals, more and more specialists pay attention to the analysis of dangerous PDFs. Integration of so many analysis tools in a single antivirus software has a negative impact on computers performance. For this reason, the cybersecurity field attaches importance to Cloud Computing. In the following sections we will analyze other academic and industrial approaches for transferring antimalware engines to the cloud and some documented detection techniques of malicious PDFs.

### 2.1 Cloud based malware detection

In the field of Cloud solutions for malware detection, there is a sparse amount of shared academic works. As compensation, in the antivirus industry there is a fast growing interest for Cloud Computing, which has higher computational power and could therefore run more advanced and complex detection algorithms.

An important example is **VirusTotal** [7], a popular Cloud application developed by Hispasec Sistemas, that aggregates more than 50 antivirus engines and makes them publicly available for scanning uploaded files. From the user perspective, this is an excellent application, that can extract metadata of the submitted file and can identify any dubious signal. The provided result represents a comparison between analysis verdicts of cybersecurity market leaders. Of course this is an advantage in terms of the scanning result precision and respectively a gain regarding spent computational time. On average it takes aprox. 55 seconds to upload and analyze a 400KB file. VirusTotal is also helping to maintain the global cybersecurity at a high level, by sharing all the submitted suspicious files with the security researchers. Thereby the antivirus engines will be permanently improved.

**Sandbox Analyzer** is another antimalware cloud solution developed by Bitdefender [1]. Its approach is a bit different, because it ensures the security on a private network, where the Bitdefender product is installed, thus not being available for public access. Its operating principle is also specific, by preventing the execution of threats on an endpoint and automatically sending of harmful files to the Cloud. After extra analysis in the Cloud, the Sandbox Analyzer can take remediation action based on the verdict. In that way, malicious files get disinfected, quarantined or deleted. One of the benefits for this solution is in-depth analysis of malicious files in an isolated environment, rather than on user's machine. Thereby the risk for performance implication, as well as the risk of accidental run of malware on an endpoint machine are eliminated. At the core of the Sandbox Analyzer there are Machine Learning algorithms and dynamic behavior analysis techniques that are constantly improved to detect fresh threats.

## 2.2 Detection of malicious PDF

**PJScan**, presented in the paper of Laskov and Srndic [5], is a Machine Learning approach that trains One-Class Support Vector Machine (*OCSVM*) to classify PDF files. This approach is focused on static analysis of embedded Javascript code, as it is known for the ability of integrating malicious behavior in PDF documents. The authors used *n-gram* analysis to extract lexical features, such as Javascript operators and other tokens. The obtained sequence of features served later as input for the machine learning algorithm. The trained model can correctly classify malicious samples containing Javascript code. However PJScan has a lower accuracy, because it is not able to detect obfuscated parts and there are also some samples containing other types of malicious payload, such as SWF<sup>1</sup>.

Another example of static analysis of PDF Structure is **PDF Tools** by Didier Stevens [8]. It represents a suite of tools for scanning, parsing and dumping PDF files. This approach focuses on identifying the fundamental elements of the format, such as Streams, Cross Reference Tables etc. The advantage of PDF Tools is their ability of name obfuscation handling and their simplicity. Because of their high speed performance, PDF Tools are largely used in cybersecurity research.

A completely new approach mentioned in the research paper of Fettaya et al. [3] describes an algorithm that uses a Convolutional Neural Network (*CNN*) to detect malicious PDF files. The trained model, based on a single convolutional layer with a global max pool and a linear layer doesn't require any data preprocessing. Instead of this, the model is trained using as input the binary representations of the files. It is worth noting that the described algorithm could be efficiently used for distinguishing various families of malware. The rate of correct detections achieves 94%, the algorithm being also capable to classify approx. 80% of the malware into

---

<sup>1</sup>Adobe Flash file format used for multimedia

different categories.

# Chapter 3

## Scientific Problem

### 3.1 Problem definition

The purpose of the research in this thesis consists in approaching some alternative methods to replace classic file scanning techniques. Particularly, our target is to demonstrate that Machine Learning algorithms can be used as an efficient mechanism for malicious PDF files detection. Additionally, we want to present a framework that can remotely analyze suspicious PDF files by applying earlier mentioned performant algorithms. This framework has several advantages:

- Takes on the complicated task of applying classification model for uploaded files and giving the analysis result to the users.
- It is designed to be deployed as a Cloud Application on a powerful server that can handle multiple requests at the same time, lacking users care about having high specifications for their computers.
- Guarantees the privacy of the scanned documents. The algorithms doesn't require the documents content in order to decide whether they are malicious or benign. Also, after analysis process, none of the documents are stored on the Cloud.
- Keeps the detection algorithms up-to-date. It won't be the users responsibility anymore to regularly update the software antivirus installed on their computers.
- Provides a generic implementation for analyzing uploaded files. The support for a new file extension to be scanned, can be effortlessly added by adding a classifier Machine Learning model for the wanted file format.

- 3.2 Background processes in Microsoft Windows
- 3.3 Filesystem monitoring
- 3.4 Analyzing PDF File Structure
- 3.5 Malware in PDF
- 3.6 Machine Learning for Malware Detection
- 3.7 Benefits of Cloud Computing

# Chapter 4

## Proposed approach

### 4.1 Dataset

### 4.2 Proof of Concept

#### 4.2.1 Feature Selection

#### 4.2.2 Classification Techniques

#### 4.2.3 Performance Evaluation

#### 4.2.4 Experiment

### 4.3 Used technologies

#### 4.3.1 Microsoft .NET Framework

[\[10\]](#)

#### 4.3.2 PDF Tools and Metasploit Framework

[\[6\]](#) [\[11\]](#) [\[4\]](#)

#### 4.3.3 Python Flask

#### 4.3.4 Scikit-learn Machine Learning Library

[\[2\]](#) [\[9\]](#)



#### 4.3.5 ReactJS Framework

# Chapter 5

## Application

### 5.1 Design

### 5.2 Windows Service

### 5.3 Cloud API

### 5.4 Dashboard Interface

## Chapter 6

### Conclusion and future work

# Bibliography

- [1] Bitdefender. Sandbox analyzer. <https://download.bitdefender.com/resources/files/News/CaseStudies/2017-TechnicalBrief-SandBoxAnalyzer-crea2103-A4-en-EN-2-GenericUse.pdf>. [Online; Accessed 1-April-2020].
- [2] Clarence Chio and David Freeman. *Machine Learning and Security*. O'Reilly Media, Inc., 2018.
- [3] R. Fettaya and Y. Mansour. Detecting malicious pdf using cnn. *International Conference on Learning Representations*, 2020.
- [4] Adobe Systems Incorporated. *PDF Reference, sixth edition: Adobe Portable Document Format Version 1.7*. Adobe, 2006.
- [5] P. Laskov and N. Srndic. Static detection of malicious javascript-bearing pdf documents. In *Proceedings of the 27th Annual Computer Security Applications Conference*, pages 373–382. ACM, 2011.
- [6] Xakep Magazine. Looking for exploits in pdf-documents on our own. <https://xakep.ru/2014/09/26/search-document-exploit/>, 2014.
- [7] Hispasec Sistemas. Virustotal. <https://www.virustotal.com/>. [Online; Accessed 1-April-2020].
- [8] Didier Stevens. Pdf tools. <https://blog.didierstevens.com/programs/pdf-tools/>, 2008.
- [9] Emmanuel Tsukerman. *Machine Learning for Cybersecurity Cookbook*. Packt Publishing, 2019.
- [10] Pavel Yosifovich, Alex Ionescu, Mark E. Russinovich, and David A. Solomon. *Windows Internals*. Microsoft Press, 2017.
- [11] Lenny Zeltser. <https://zeltser.com/information-security/>. [Online; Accessed 1-April-2020].