Accepted Manuscript

Classification of ransomware families with machine learning based on N-gram of opcodes

Hanqi Zhang, Xi Xiao, Francesco Mercaldo, Shiguang Ni, Fabio Martinelli, Arun Kumar Sangaiah



 PII:
 S0167-739X(18)30732-5

 DOI:
 https://doi.org/10.1016/j.future.2018.07.052

 Reference:
 FUTURE 4370

To appear in: Future Generation Computer Systems

Received date :31 March 2018Revised date :9 July 2018Accepted date :25 July 2018

Please cite this article as: H. Zhang, X. Xiao, F. Mercaldo, S. Ni, F. Martinelli, A.K. Sangaiah, Classification of ransomware families with machine learning based on *N*-gram of opcodes, *Future Generation Computer Systems* (2018), https://doi.org/10.1016/j.future.2018.07.052

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

Classification of Ransomware Families with Machine Learning Based on N-gram of Opcodes

Hanqi Zhang^{a, b}, Xi Xiao^b, Francesco Mercaldo^c, Shiguang Ni^{b*}, Fabio Martinelli^c, Arun Kumar Sangaiah^d

^aCollege of Physical Science and Technology, Central China Normal University, Wuhan, China
 ^bGraduate School at Shenzhen, Tsinghua University, Shenzhen, China
 ^cInstitute for Informatics and Telematics, National Research Council of Italy, Pisa, Italy
 ^dSchool of Computing Science and Engineering, VIT University, Vellore, India

Abstract: Ransomware is a special type of malware that can lock victims' screen and/or encrypt their files to obtain ransoms, resulting in great damage to users. Mapping ransomware into families is useful for identifying the variants of a known ransomware sample and for reducing analysts' workload. However, ransomware that can fingerprint the environment can evade the precious work of dynamic analysis. To the best of our knowledge, to overcome this shortcoming, we are the first to propose an approach based on static analysis to classifying ransomware. First, opcode sequences from ransomware samples are transformed into N-gram sequences. Then, Term frequency-Inverse document frequency (TF-IDF) is calculated for each N-gram to select feature N-grams so that these N-grams exhibit better discrimination between families. Finally, we treat the vectors composed of the TF values of the feature N-grams as the feature vectors and subsequently feed them to five machine-learning methods to perform ransomware classification. Six evaluation criteria are employed to validate the model. Thorough experiments performed using real datasets demonstrate that our approach can achieve the best Accuracy of 91.43%. Furthermore, the average F1-measure of the "wannacry" ransomware family is up to 99%, and the Accuracy of binary classification is up to 99.3%. The proposed method can detect and classify ransomware that can fingerprint the environment. In addition, we discover that different feature dimensions are required for achieving similar classifier performance with feature *N*-grams of diverse lengths.

Keywords: ransomware classification, static analysis, opcode, machine learning, N-gram

1. Introduction

According to Internet Security Threat Report [1], more than 357 million new malware samples were disclosed in 2016. Malicious software can erode or steal data and destroy computer systems [2-4]. Ransomware is one of the most harmful types of malware, and it forces victims to pay ransoms in exchange for encrypted data [5]. Unlike traditional malware, it is difficult to kill ransomware even when it is discovered, and the loss is irreversible even its removal [6]. The earliest known ransomware "AidsInfo" was discovered in 1989. Thereafter, ransomware developed rapidly, and in recent years, many new families of ransomware have emerged. Last year, "wannacry" ransomware infected more than 200,000 computers in 150 countries in less than a day [7], posing enormous danger to users.

Ransomware classification can help identify the variants of a known ransomware sample and effectively reduce the workload of analysts. However, most studies on ransomware pertain to its

^{*} Corresponding author at: Graduate School at Shenzhen, Tsinghua University, Shenzhen, China Tel: +86 18038153875. E-mail address: ni.shiguang@sz.tsinghua.edu.cn

Download English Version:

https://daneshyari.com/en/article/6872779

Download Persian Version:

https://daneshyari.com/article/6872779

Daneshyari.com