

IEEE Copyright Notice

© 2021 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

Trends in Attack Cases and Vulnerability Candidates From Large Vulnerability Database Using Topic Model Analysis

Hiroki Koyanagi

*Electrical and Information Eng., Shonan Inst. of Tech.
Cyber Physical Security Research Center, AIST
Kanagawa, Japan
20t2006@sit.shonan-it.ac.jp*

Sven Wohlgemuth

*Yokohama Research Laboratory, Hitachi, Ltd.
Kanagawa, Japan
sven.wohlgemuth.kd@hitachi.com*

Kazuo Takaragi

*Cyber Physical Security Research Center, AIST
Tokyo, Japan
kazuo.takaragi@aist.go.jp*

Katsuyuki Umezawa

*Dept. of Information Science, Shonan Inst. of Tech.
Cyber Physical Security Research Center, AIST
Kanagawa, Japan
omezawa@info.shonan-it.ac.jp*

Abstract—Security cases have been increasing in recent years. Security problems are caused by human factors related to existing software vulnerabilities, which are often caused by insufficient testing. These problems can be resolved by software improvements using knowledge gained by previous experience and artificial intelligence that mechanically minimizes the number of human errors. To this end, we propose a method that informs the designer of the vulnerabilities inherent in a design document. The vulnerabilities are expressed as an attack tree built from the documents of past attack cases. Specifically, we searched for and narrowed down the vulnerabilities that can be exploited in attacks by matching past case documents with the natural languages of vulnerability information. Although the narrowed data contained the correct answer, the number of retrieved cases was several thousand. In this paper, we analyze trends in vulnerability type on the top 20 cases in the narrowed dataset.

Index Terms—Threat analysis, Vulnerability information, Attack tree, Natural language processing, Topic model

I. INTRODUCTION

In recent years, the number of serious security breaches has increased. Many of these incidents are often the result of attacks by malicious external actors, and these are caused by the existence of software and system vulnerabilities. Eradicating these vulnerabilities should solve many of these incidents, but this is, in reality, difficult, the reason being that vulnerable areas are mistakenly created by human error, such as poor coding and oversight of test items. Therefore, it is necessary to reduce vulnerabilities from human error by using mechanical methods such as AI. In many cases, new attacks are similar to previously launched attacks; therefore, it is necessary to utilize the vast amount of vulnerability information that is available about past attacks. We can systematically present an attack procedure that was performed in the past by using a tree and applying it to future cases or to a system currently being built to present what a possible attack would look like against

a similar system. We have proposed a method to manage this, and have used topic modeling techniques, such as latent Dirichlet allocation (LDA), to match the natural language of past case documents with large-scale vulnerability databases. The matching process guaranteed the target vulnerabilities within the narrowed-down document group, but the cases still numbered in the thousands. In this paper, we analyze the types of vulnerabilities among the top 20 vulnerabilities in the narrowed number of cases. We believe that by clarifying the types of vulnerabilities using the knowledge gained from past cases, we can guide the development of effective countermeasures and guidelines for necessary future tests.

II. PREVIOUS WORK

We used the BROWSER HACKING section (hereafter referred to as the “BH document”) and the LOCAL PRIVILEGE ESCALATION section (hereafter referred to as the “LPE document”) of the paper [1] as case sentences from past studies [2] [3], and the vulnerability information contained therein. These were collated with a large-scale vulnerability database [4]. LDA [5] was one of the topic modeling methods used to perform matching. The cosine similarity was calculated from the results, and the results shown in Table I are obtained when the calculated cosine similarity is viewed with the actual vulnerability information used in the case text as the lower limit of the threshold value.

III. METHOD

We first focus on the top 20 cases among the 1514 and 5180 cases shown in Table I, which were narrowed down in a previous study [3]. Next, we manually confirm which category of the Common Weakness Enumeration (CWE) [6] of the MITRE Corporation in the United States includes all 40 of these cases. The confirmation was made through a security

TABLE I
RESULTS OF MATCHING PROCESSING IN PAST STUDIES

| Case document section title | Number of topics | Total number of cases | Applicable number | Applicable ratio |
|-----------------------------|------------------|-----------------------|-------------------|------------------|
| BH document | 7 | 119479 | 1514 | 1.3 |
| LPE document | 11 | 119479 | 5180 | 4.3 |

vulnerability database [7]. We also examined the CWE IDs of the vulnerabilities used in the case documents (CVE-2011-3928 and CVE-2013-6282) and assessed their relevance. The results for CVE-2011-3928 and CVE-2013-6282 are given in Table II.

TABLE II
VULNERABILITY DATA IN CASE STUDIES

| CVE name | CWE ID | Products Affected |
|---------------|--------|-------------------|
| CVE-2011-3928 | 399 | Google Chrome |
| CVE-2013-6282 | 20 | Linux Kernel |

IV. RESULT AND ANALYSIS

First, we focused on the BH document and CVE-2011-3928 for analysis. When the top 20 cases were extracted, there were five types of CWE IDs: 19, 119, 189, 264, and 399. CWE ID 399 was the number to which CVE-2011-3928 belonged. In other words, CVE-2011-3928 could not be directly shown, but was a result that could present similar vulnerabilities elsewhere. Looking at the LPE document and CVE-2013-6282 from the same perspective, there were nine types of CWE IDs for the top 20 vulnerabilities: 20, 77, 78, 125, 200, 264, 362, 399, and 476. Furthermore, CWE ID 20 was the number to which CVE-2013-6282 belonged, and CVE-2013-6282 could not be directly indicated, but a similar vulnerability could be presented. From this result, it has been shown that the necessary information used as a guideline for the effective measures and tests in question in this proposal could be provided.

V. CONSIDERATION

In this proposal, the CVE result was manually linked to the CWE ID number. It was thus confirmed that it is possible to provide information similar to trends in the top 20 cases and the vulnerabilities actually used in the case documents. However, we were not able to verify what results would be obtained if we used a more refined classification than the one used for CWE. For example, CWE ID 399, to which CVE-2011-3928 belonged, was a category of Resource Management Errors. It is not known what causes these in the future when classified by the CWE ID. Therefore, manpower is required to understand the details. In order to solve this problem, it is sufficient to provide classification using a more refined database, but if this does not exist or is insufficient, it will be necessary to create one separately. For example, grouping

is performed for CVEs by using the same CWE ID, and classification by parts and functions can be considered.

VI. CONCLUSION

In this proposal, the vulnerabilities extracted by matching past case sentences and a large-scale vulnerability database were narrowed down to the top 20, and further analyzed. Specifically, we investigated the correspondence with CWE IDs for the top 20 CVE data and compared this with the CWE IDs of the vulnerabilities actually used in the case sentences. As a result, we were able to confirm the vulnerabilities with the same CWE ID as the vulnerabilities actually used in the case document in the top 20 matching results for each sentence. From this, it is possible to present vulnerabilities similar to those that are intended for presentation by matching past case documents with the vulnerability database and presenting the vulnerabilities that have the highest number of hits. We were able to verify this. However, because the range of classification was wide, it was also necessary to categorize the vulnerability information, as well as the vulnerability information that is classified more closely.

TRADEMARK

- CAPEC™ and the CAPEC logo are trademarks of The MITRE Corporation.
- CWE™ and the CWE logo are trademarks of The MITRE Corporation.
- Google Chrome™ is a trademark of Google LLC.
- Linux® is the registered trademark of Linus Torvalds in the U.S. and other countries.

ACKNOWLEDGMENTS

A part of this work was supported by Council for Science, Technology and Innovation (CSTI), Crossministerial Strategic Innovation Promotion Program (SIP), “Cyber-Security for Critical Infrastructure” (funding agency: NEDO).

REFERENCES

- [1] S. Nie et al.: FREE-FALL: HACKING TESLA FROM WIRELESS TO CAN BUS, Briefing,Black Hat USA 2017, July 2017.
- [2] Hiroki Koyanagi, Kazuo Takaragi, Yusuke Mishina, Sven Wohlgenuth, and Katsuyuki Umezawa, “Threat Analysis Method using Vulnerability Database: Attack case and large-scale vulnerability DB matching by topic model analysis using LDA classifier and cosine similarity,” 2020, 2020-CSEC-88(38), pp. 1-6.(in Japanese)
- [3] Hiroki Koyanagi, Kazuo Takaragi, Yusuke Mishina, Sven Wohlgenuth, and Katsuyuki Umezawa, “Threat analysis using the vulnerability database: Application of Attack Cases and Large-Scale Vulnerability DB Collation Methodology to Multiple Cases Using Topic Model Analysis,” 2020, 2020-SPT-38(16), pp. 1-6.(in Japanese)
- [4] MITRE Corporation, “CVE - Common Vulnerability and Exposure,” <https://cve.mitre.org/> (Last accessed 16 Dec. 2020).
- [5] D. Blei, A. Ng, and M. Jordan,: Latent Dirichlet Allocation,in Journal of Machine Learning Research(2003), pp. 1107-1135.
- [6] MITRE Corporation, “CWE List - Common Weakness Enumeration,” <https://cwe.mitre.org/data/> (Last accessed 16 Dec. 2020).
- [7] CVE Details The ultimate security vulnerability datasource, <https://www.cvedetails.com/> (Last accessed 16 Dec. 2020).