

CVE 동향을 반영한 3-Step 보안 취약점 위험도 스코어링

임지혜¹ · 이재우^{2*}

3-Step Security Vulnerability Risk Scoring considering CVE Trends

Jihye Lim¹ · Jaewoo Lee^{2*}

¹Graduate Student, Department of Convergence Security, Chung-Ang University, Seoul, 06974 Korea

^{2*}Assistant Professor, Department of Industrial Security, Chung-Ang University, Seoul, 06974 Korea

요 약

보안 취약점 수가 해마다 증가함에 따라 보안 위협이 지속해서 발생하고 있으며 취약점 위험도의 중요성이 대두되고 있다. 본 논문에서는 보안 취약점 위험도 판단을 위해 동향을 반영한 보안 위협 스코어링 산출식을 고안하였다. 세 단계에 따라 공격 유형과 공급업체, 취약점 동향, 최근 공격 방식과 기법 등의 핵심 항목 요소를 고려하였다. 첫째로는 공격 유형, 공급업체와 CVE 데이터의 관련성 확인 결과를 반영한다. 둘째로는 LDA 알고리즘으로 확인된 토픽 그룹과 CVE 데이터 간 유사성 확인을 위해 자카드 유사도 기법을 사용한다. 셋째로는 최신 버전 MITRE ATT&CK 프레임워크의 공격 방법, 기술 항목 동향과 CVE 간의 관련성 확인 결과를 반영한다. 최종 보안 취약점 위험 산출식 CTRS의 활용성 검토를 위해 공신력 높은 취약점 정보 제공 해외 사이트 내 데이터에 제안한 스코어링 방식을 적용하였다. 본 연구에서 제안한 산출식을 통하여 취약점과 관련된 일부 설명만으로도 관련성과 위험도가 높은 취약점을 확인하여 신속하게 관련 정보를 인지하고 대응할 수 있다.

ABSTRACT

As the number of security vulnerabilities increases yearly, security threats continue to occur, and the vulnerability risk is also important. We devise a security threat score calculation reflecting trends to determine the risk of security vulnerabilities. The three stages considered key elements such as attack type, supplier, vulnerability trend, and current attack methods and techniques. First, it reflects the results of checking the relevance of the attack type, supplier, and CVE. Secondly, it considers the characteristics of the topic group and CVE identified through the LDA algorithm by the Jaccard similarity technique. Third, the latest version of the MITRE ATT&CK framework attack method, technology trend, and relevance between CVE are considered. We used the data within overseas sites provide reliable security information to review the usability of the proposed final formula CTRS. The scoring formula makes it possible to fast patch and respond to related information by identifying vulnerabilities with high relevance and risk only with some particular phrase.

키워드 : 사이버 위협 인텔리전스, 위협 산출식, 취약점 위험도, LDA 토픽 모델링, 자카드 유사도

Keywords : Cyber Threat Intelligence, Threat Scoring, Vulnerability Risk, LDA Topic Modeling, Jaccard Similarity

Received 19 October 2022, Revised 25 October 2022, Accepted 7 November 2022

* Corresponding Author Jaewoo Lee(E-mail: jaewoolee@cau.ac.kr, Tel: +82-2-820-5935)

Assistant Professor, Department of Industrial Security, Chung-Ang University, Seoul, 06974 Korea

Open Access <http://doi.org/10.6109/jkiice.2023.27.1.87>

print ISSN: 2234-4772 online ISSN: 2288-4165

© This is an Open Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License(<http://creativecommons.org/licenses/by-nc/3.0/>) which permits unrestricted non-commercial use, distribution, and reproduction in any medium, provided the original work is properly cited.
Copyright © The Korea Institute of Information and Communication Engineering.

I. 서론

CVE 보안 취약점 데이터베이스, 정보 소스인 CVE Details에 따르면 2021년도 한 해 동안 20,171개의 취약점(CVE)이 발표되었다[1]. 최근 10년간 연도별 CVE 건수는 꾸준히 증가하였다. 해마다 취약점의 수가 증가함에 따라, 하루에도 수많은 취약점이 새롭게 발표되고 있으며 취약점과 관련된 여러 보안사고가 발생하고 있다. 이러한 보안사고의 경우 아직 패치 또는 업데이트가 적용되어 있지 않은 시스템이 취약 대상으로 공격의 대상이 되기 때문에 개인정보 유출, 시스템 손상 등의 큰 피해로 이어질 수 있다.

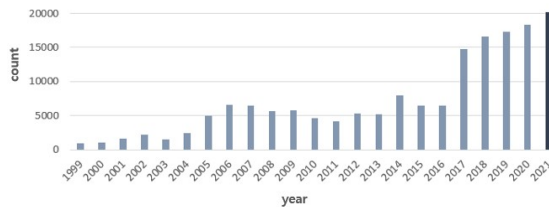


Fig. 1 Number of CVEs per year (CVE Details)

변화하는 보안 흐름에 따라, CVE 취약점에 동향을 반영한 보안 취약점 위험도 스코어링 시스템을 제안하고자 한다. 보안 취약점 위험에 대한 여러 연구가 이루어지고 있으나, 위험 산출에 있어 고려해야 하는 항목이 다양하기에 어떠한 요소를 어떻게 반영해 위험도를 확인할 것인지가 중요하다.

취약점 위험도 산출에 앞서 취약점 생명주기에 대해 유의 깊게 살펴볼 필요가 있다. 취약점 생명주기를 그림으로 시각화하여 나타낸 것은 그림 2와 같다. 취약점이 처음 발견되고 취약 부분에 대한 공격이 활발하게 발생한다. 이후, NVD(National Vulnerability Database)라는 공식 취약점 번호를 부여하는 미국 표준 기관에서 식별자를 승인하고 공개적으로 취약점과 CVSS점수가 발표된다. 취약점 발표 이후, 해당 취약점에 영향을 받는 여러 소프트웨어, 하드웨어 공급 및 유관 기관 공급업체(Vendor)에서 대응책 및 패치를 배포 및 발표한다.

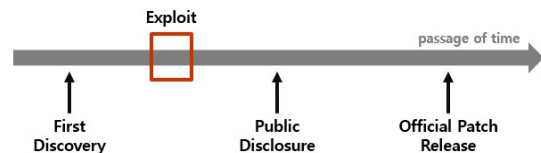


Fig. 2 The life cycle of a vulnerability

이 과정에서 취약점 발견 이후부터 악용되는 시기, 즉, 공식 패치 이전까지의 단계가 매우 중요하다. 공식적인 발표 또는 분석 이전까지 어떠한 취약점인지 위험도가 높은지 판단하기 어렵기 때문이다.

본 논문에서는 이러한 초기 발견된 간단한 정보만으로 위험 관련도가 높은 과거의 취약점들까지 보여줄 수 있는 취약점 위험 스코어링 방식을 연구하였다.

본 논문의 공헌은 다음과 같다.

- 취약점 동향 반영을 위해 발표된 모든 CVE와 MITRE ATT&CK(Adversarial Tactics, Techniques and Common Knowledge) 프레임워크 정보[2]를 활용하였다.
- LDA(Latent Dirichlet Allocation) 토픽 모델링[3], 자카드 유사도[4] 등의 연구 방법을 이용하여 동향을 반영한 취약점 위험 산출식 CTRS를 제안하였다.
- 공신력 높은 해외 뉴스 취약점 정보에 제안한 스코어링 방식을 적용하여 활용성을 검토하였다.

본 논문은 2장에서 보안 취약점의 위험 산출, LDA 알고리즘 관련 기존 연구들을 살펴본다. 3장에서는 연구 방법과 활용하는 연구 데이터에 관해 설명하고, 4장에서는 제안하는 위험도 스코어링 과정을 설명한다. 또한 최종 결과를 실제 데이터에 적용한 과정과 활용 방법을 설명한다. 이후, 5장에서는 연구 결과를 정리한 결론과 개선 방안을 포함한 향후 연구 방향을 제시한다.

II. 관련연구

2.1. 취약점 위험 산출 관련 연구

보안 취약점이 빠르게 증가하면서, 관련 취약점 대응 및 위험도 산출에 관한 연구가 활발하게 진행되고 있다. 취약점 위험도 확인을 위하여 취약점 자체 공격 사이클 내 위험 요소 또는 평가 척도 항목들을 고려한 연구들이 있다. Haipeng Chen[5]은 취약점 악용 예상 위험 점수를 반영한 시스템인 VEST(Vulnerability Exploit Scoring & Timing)를 통해 사이버 취약 부분에 대한 사전 분석 및 경고에 활용할 수 있는 정보를 제공하고자 하였다. 김민철[6]은 증가하는 소프트웨어 취약점의 대응 방법으로 취약점 자체 정보와 함께 MITRE에서 제공하는 CWSS의 평가 척도를 활용하여 취약점에 대한 위험도 스코어링 시스템을 제안하였다.

이 요소들 외에 취약점의 환경, 흐름, 악용의 시기 등의 중요성을 강조한 연구도 있다. 진희훈[7]은 대상, 환경에 따라 달라질 수 있는 실제 웹 취약점 위험도가 사전 평가된 위험도와 다를 수 있는 점을 개선하기 위해 위험도 평가 모델을 제시하였고, 결과 도출을 위해 사이버 킬체인 중심의 공격 시도 결과를 활용하였다. 이러한 요소와 유사하게 취약점 개념 증명(Proof Of Concept, POC)을 중요 요소로 보는 연구도 있다. 박찬일[8]은 네트워크 보안 위협 분석에 사용하기 위한 공격 그래프 분석을 위해 취약점 착취 여부에 따른 위험도를 제시하였고, 등급 가중치를 부여하여 최종 위험도 레벨을 계산하였다. 하지만 이미 해당 취약점이 널리 공개되어 공격 시도가 다수 발생한 이후이기 때문에 POC가 없는 취약점의 경우 한계가 있다.

본 논문에서는 변화하고 있는 보안 취약점 위협에 동향을 반영할 수 있는 연구 데이터와 연구 방법을 활용하여 위험도와 관련 취약점을 알아보고자 한다.

2.2. LDA 모델링 관련 연구

토픽 모델링 기법의 하나로 대표적인 LDA 기법은 다양한 분야에서 동향 파악을 위해 사용되고 있다. 채호근[9]은 국내외 금융 보안 분야 연구 동향을 분석하기 위하여 LDA 분석을 통해 10개의 토픽을 선정하였고, 국외 금융 보안 연구 동향을 파악해 국내 금융 보안 연구의 보완점을 찾고자 하였다. 이선우[10]는 국내 클라우드

보안 동향 변화 분석을 위하여 LDA 기법을 통해 코로나19 대유행 전후 동향을 비교해 보았으며 클라우드 보안에 대한 관심 향상에 대비한 보안 대책 수립의 중요성을 확인하였다. Choi, Hyo Shin[11]는 개인정보 프라이버시의 학술 동향을 살펴보기 위하여 Scopus DB 내 저널 기사, 학회 논문 등 다양한 문서 2,356개 문서 대상으로 LDA 알고리즘을 적용하였다. 실증적이고 객관적 주제들을 확인하여 기존 분류 체계를 보완할 수 있는 개선점을 시사하였다. 이처럼 LDA 기법은 보안 분야에서의 동향 파악에도 활발하게 사용되고 있다. 본 연구는 기존 연구에서는 적용하지 않았던 취약점 데이터에 LDA 토픽 모델링을 접목해 CVE 데이터 내 동향을 반영한 토픽 그룹 도출에 활용하고자 한다.

III. 연구 데이터 및 방법

본 논문에서는 보안 취약점 위험 산출에 있어 동향을 효과적으로 반영하기 위하여 현재까지 발표된 전체 취약점 데이터 대상 분석을 통한 흐름을 적용한다. 또한 기존 연구들에서의 주관적 평가 요소 항목들을 대신해 사이버 공격 전술, 절차에 있어 표준 모델로 자리 잡고 있는 MITRE ATT&CK 매트릭스를 활용한다. 본 논문의 보안 취약점 위험도 산출 과정은 그림 3과 같으며, 연구 방법으로 토픽 모델링 방식, LDA 알고리즘, 자

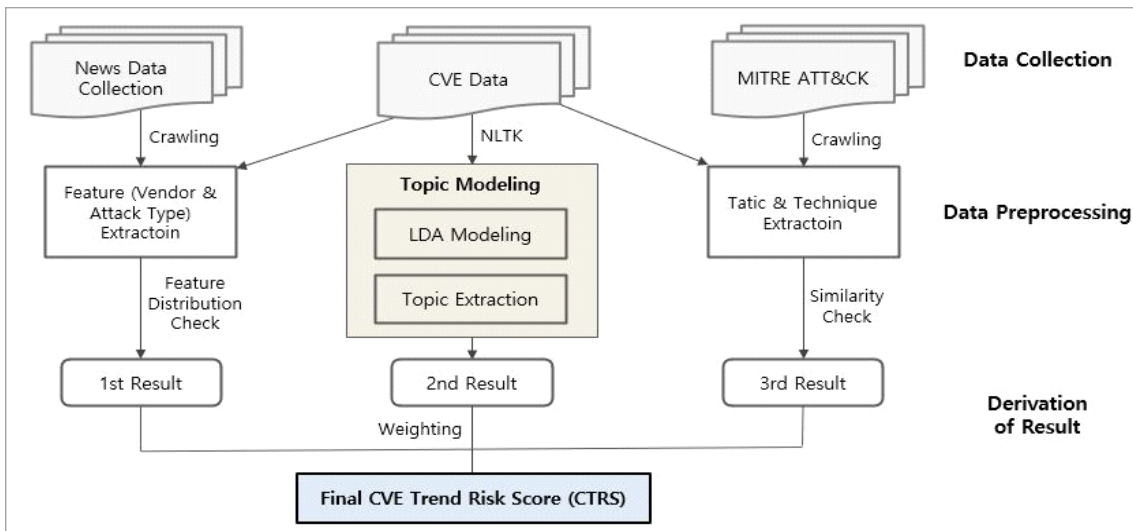


Fig. 3 Security vulnerability risk scoring calculation process - CTRS(CVE Trend Risk Score)

카드 유사도 기법을 사용한다.

각 단계에서 서로 다른 형태의 데이터를 적용해 동향을 반영하고자 하였으며 단계별 결과에 가중치를 부여하여 최종 결과를 도출한다.

3.1. 연구 데이터

본 논문에서는 취약점 목록과 위험도 스코어링 예측에 필요한 정보를 수집하기 위하여 국내외 20여 개 사이트의 보안 및 취약점 관련 뉴스 제공 정보, MITRE의 CVE와 MITRE ATT&CK 프레임워크를 활용하였다.

3.1.1. 보안 뉴스 데이터

Twitter, NIST News, CISA Alerts, 보안뉴스, 데일리시큐, 전자신문 등 20여 개의 사이트에서 최근 한 달 보안 뉴스 데이터를 수집하여 워드 카운팅 작업을 수행하였다. 단어 빈도수를 확인하여 시각적으로 나타난 워드 클라우드를 그림4와 같다. 워드 카운팅 결과 각 단어의 개수가 높게 집계될수록 워드 클라우드에서 크기가 큰 단어로 표기된다.



Fig. 4 Latest Security News Word Cloud

본 연구에서는 워드 카운팅 결과에서의 상위 항목으로 집계된 Vendor와 공격 유형 정보를 추출하여 이후 CVE 취약점 데이터와 함께 활용한다.

3.1.2. CVE 취약점 데이터

MITRE에서 CVE가 초기 발표된 시점인 1999년부터 2022년 6월까지의 모든 CVE 239,898개를 수집한다 [12]. 취약점 원시 데이터에서 다중 단어를 추출하고 중복 단어를 제거한다. Vendor, 공격 유형과 관련된 항목을 최종 선별하고 정규화 작업을 거쳐 연구 데이터를 마련한다.

3.1.3. MITRE ATT&CK 데이터

MITRE ATT&CK 프레임워크의 엔터프라이즈용 Tactic(공격 방법)과 Technique(기술) 데이터를 크롤링해 수집한다. MITRE ATT&CK은 공격자가 네트워크에 침투하고 데이터 유출 및 특정 목표를 달성하기 위해 사용하는 여러 기술을 추적하기 위한 모델로, 다양한 진술에 따라 분류된 사이버 공격 기술 매트릭스이다. 체계적으로 분류 및 고도화와 함께 버전 또한 업데이트되고 있으며, 침입 탐지, 보안 엔지니어링, 위협 예상 및 인텔리전스, 위험 관리 등 여러 분야에서 전 세계적으로 널리 활용되고 있다.

Tactic은 공격자가 공격하고자 하는 목표인 정보의 탐색이나 파일의 추출 및 실행 등에 따른 행동을 나타내며 상황에 따라 달라지는 각 Technique를 묶어주는 범주 역할을 한다. Technique은 공격자가 목표한 Tactic을 달성하기 위한 세부 방법을 나타내고, 해당 Technique를 사용함으로써 발생하는 결과나 피해 사항 등을 명시한다. 본 논문에서는 가장 최신 버전으로 확인된 MITRE ATT&CK for Enterprise v11의 14가지 Tactic과 222개의 Technique 항목을 수집하여 활용한다.

3.2. 연구 방법

3.2.1. 토픽 모델링

토픽 모델링은 구조화되어 있지 않은 많은 양의 전체 문서 집합에서 숨겨져 있는 의미 구조를 발견하고 주제를 찾기 위한 텍스트 마이닝 기법이다. 전체적인 맥락과 관련도가 높은 정보들을 통해 유사한 의미를 뜻하는 단어들을 군집화하여 클러스터링하는 방식으로 문서의 주제를 추론할 수 있다. 또한 대량의 텍스트에서 사용된 단어들을 분석하여 발견한 주제들이 어떠한 관련성을 가지고 연결 지어지는지를 분석해 주는 통계적 방법이다. 이를 통해 기존의 키워드들만으로 찾아내기 어려운 의미를 탐색할 수 있어 다양한 분야에 대한 동향 분석에 주로 사용되고 있다.

3.2.2. LDA(Latent Dirichlet Allocation)

본 논문에서는 전처리 과정을 거친 약 24만 개의 CVE 연구 데이터 대상인 각 CVE를 설명하는 텍스트로부터 단어의 교환성(Exchangeability)을 고려하지 않고, 단어의 동시 발생 빈도를 바탕으로 내재되어 있는 주제들을 발견하기 위하여 통계적 알고리즘인 LDA를 활용

하였다.

LDA 알고리즘은 확률 기반 모델링 기법이다. 대량의 문서 데이터를 분석하여, 방대한 문서 데이터 내에 어떤 토픽이 있는지, 그 토픽 내 단어가 어떠한 비율로 구성되어 있는지 분석한다. 분석을 위해 구조화되지 않은 전체 데이터 집합에서 비슷한 데이터를 클러스터링하여 패턴이나 형태를 찾아내기 위한 방법인 비지도 학습 방식을 사용한다. LDA 과정을 거쳐 여러 단어 사이의 동시 발생 빈도를 바탕으로 문서의 토픽을 얻어내고 각 토픽이 어떤 키워드들로 구성되어 있는지 확인할 수 있다.

3.2.3. 자카드 유사도

자카드 유사도는 두 개의 텍스트 문서와 같은 집합 사이에서 집합 간의 유사한 정도를 측정해서 보여주는 방법 중 하나이다. 두 집합 A, B가 있을 경우, 두 집합의 자카드 유사도는 $J(A, B)$ 로 표기할 수 있고, 이는 두 집합 사이의 교집합 크기를 합집합 크기로 나누어 준 값을 뜻하며 이를 식으로 나타낸 것은 식 1과 같다. 이를 활용하여 문자열 사이의 유사도를 계산할 수 있다.

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|} \quad (1)$$

자카드 유사도 알고리즘은 두 개의 데이터 집합의 시퀀스 크기와 관계없이 잘 동작하기 때문에 시퀀스 길이의 가변성에 영향을 받지 않는다. 본 논문에서는 CVE 데이터에 따라 리스트의 길이가 가변적으로 변화하는 특성이 있어 유사도 측정을 위한 적합한 알고리즘으로 선정하였다.

IV. CTRS 위험 산출식 도출

본 연구에서는 CVE 보안 취약점을 활용한 위험 산출식 CTRS(CVE Trend Risk Score)를 제안한다. CVE 전체 동향을 반영하기 위해 LDA를 활용한 토픽 모델링 기법을 적용하고, MITRE ATT&CK 프레임워크의 Tatic 키워드들과의 유사도 관련도를 반영하여 현재 시점에서 주어진 상황에 가장 위협이 될 수 있는 CVE 후보를 도출해 내고자 한다. CVE 위험 평가에 아래와 같은 새로운 관점들을 고려하고, 효과적인 위험도 산출을 위하여 세 단계 과정으로 나누어 식을 고안하였다.

4.1. WordCount 빈도 수 확률 (1단계)

첫 번째 식에서 각 CVE 설명 정보들과 Vendor, Attack Type(공격 유형) 중심의 데이터 관련성을 확인한다. Vendor와 Attack Type 기반으로 CVE 전체 동향을 파악한 결과를 수식에 반영하고자 한다.

3장에서 연구 데이터 단계에서 마련한 범주 (Vendor 16종, Attack Type 40종)를 활용하였고 관련성 확인을 위해 범주별 빈도수 확률을 이용하였다. 위험 산출식에 앞서, 전체 CVE 데이터 대상으로 범주별 항목이 CVE 설명에 포함되어 있는지 비교를 위한 워드 카운팅 과정을 거친다. 각 항목이 포함된 CVE의 개수를 확인하였고, 아래의 표 1, 표 2와 같다.

Table. 1 Word counting by Attack Type

No.	Attack Type	Count	No.	Attack Type	Count
1	remote	87820	21	spoofing	664
2	unauthenticated	5552	22	smuggling	73
3	web	25107	23	crafted file	443
4	local	19165	24	memory leak	434
5	injection	12872	25	crafted image	310
6	authenticated	12696	26	virus	169
7	bypass	8350	27	trojan	842
8	buffer overflow	6345	28	encoding	235
9	code execution	6215	29	redirection	213
10	access control	1525	30	scanning	173
11	privileged	4292	31	read data	71
12	configuration	4004	32	cryptography	68
13	race condition	547	33	backdoor	121
14	security bypass	66	34	logic error	81
15	modify data	141	35	file write	77
16	denial of service	26899	36	integer overflow	1317
17	memory corruption	3024	37	crafted application	277
18	command injection	1473	38	session hijacking	51
19	directory traversal	1580	39	unauthorized access	1461
20	authentication bypass	365	40	privilege escalation	898

Table. 2 Word counting by Vendor

No.	Vendor	Count	No.	Vendor	Count
1	oracle	18910	9	google	2781
2	windows	17014	10	apache	1836

No.	Vendor	Count	No.	Vendor	Count
3	cisco	7604	11	hp	1137
4	ibm	5492	12	huawei	675
5	microsoft	5431	13	debian	208
6	linux	5260	14	opensuse	173
7	apple	3916	15	asus	173
8	adobe	2904	16	redhat	5

각 CVE 리스트별로 취약점을 설명하는 요약 데이터 들 대상으로 자연어 처리 과정을 수행한다. 문장에서 단어를 추출한 후, 불용어(Stopwords) 목록[13]에 추가로 예외 처리할 단어 목록을 추가한다. 이후 해당 단어들을 제외한 각 CVE 별 목록을 생성한다.

Vendor 16종, Attack Type 40종에 대해서도 리스트를 생성한다. 범주(ex. Vendor)에 해당하는 특정 항목(ex. Microsoft)별 전체 CVE 데이터에 해당 항목이 등장하는 개수를 확인한다. 즉, CVE 리스트별로 Attack Type과 Vendor 리스트 목록을 비교하도록 하여 CVE에 해당하는 항목이 있을 경우, 워드 카운팅 결과값을 누적하여 합산한다. 각 범주 항목들의 총합을 분모 값으로 고정한다. 다음 각 CVE 별 단어 목록에 각 범주 항목들이 포함되어 있는지를 확인하고, 포함되어 있을 경우 항목별 CVE 카운팅 값들을 모두 합한다.

범주에 따른 결과값을 합한 최종값(W)을 통해 Vendor, Attack Type과 CVE 데이터 사이의 관련 점수를 확인할 수 있다. 점수가 높을수록 관련도가 높으며, 최종 위험도 산출식 도출을 위한 항목으로 반영한다.

$$W = \frac{\sum_{n=1}^{16} V_n}{V_{count}} + \frac{\sum_{n=1}^{40} A_n}{A_{count}} \quad (2)$$

- V : Vendor's total sum of hits
- A : Attacktype's total sum of hits
- Vcount : Sum of Vendor hits in CVE
- Acount : Sum of Attacktype hits in CVE

4.2. LDA토픽과 CVE 유사도 (2단계)

두 번째 식에서는 LDA 알고리즘을 활용한 토픽 그룹과 CVE 데이터 간의 유사도를 확인한다. 약 24만 건의 CVE 설명에서의 토픽 모델링을 통해 확인된 토픽의 특징을 반영하고자 한다.

LDA 모델에 들어갈 객체를 만들기 위하여, 문서 집

합(corpus)을 만들어 자연어 처리 과정을 거쳐야 한다. 전체 문서 데이터를 문자열 단위(token)로 나누고 문자열별 빈도수를 확인하여 BoW(Back of Words) 형식으로 바꿔준다. 모든 단어를 문자열로 토큰화한 후 객체를 생성하고 토큰화된 단어 목록 확인한다.

LDA 결과 확인을 위해 최적 토픽 그룹의 수를 알아야 한다. 전체 CVE 데이터 대상 LDA 모델링 결과, Coherence(일관성) 값과 Perplexity(혼란도) 값을 활용하여 최적 토픽 그룹을 확인하였다. Perplexity 값은 확률 모델 자체가 모델링 결과를 얼마나 정확하게 예측하는지 판단하는 지표이고 낮을수록 정확하다. Coherence 값은 토픽이 얼마나 의미론적으로 일관성 있는지 판단하는 지표이고 높을수록 일관성이 높다. 토픽 그룹 범위를 2부터 15까지로 하여 계산한 결과는 그림 5, 그림 6이다. 최적의 토픽 그룹의 수는 3개로 확인하였다.

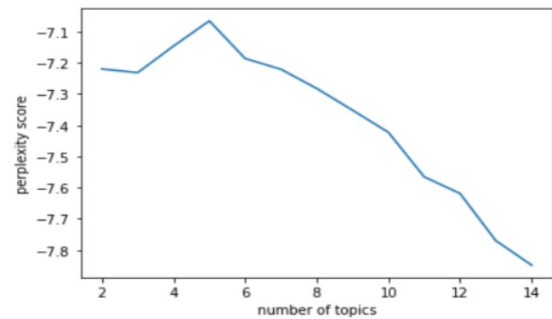


Fig. 5 Perplexity values

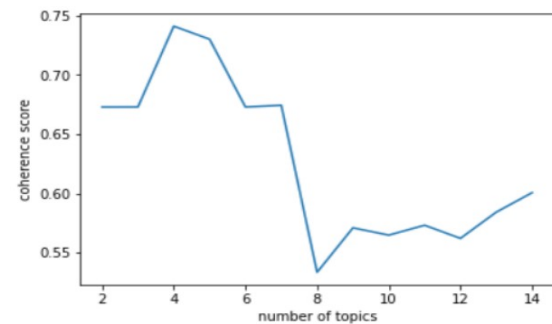


Fig. 6 Coherence values

LDA 모델 학습 결과를 시각적으로 표현하는 라이브러리인 pyLDAvis를 이용하여 최종 결과를 시각화한 결과는 그림7과 같다. 왼쪽의 원들은 각 토픽 그룹들을 의미하고 오른쪽의 막대그래프들은 토픽 그룹에 포함된

단어와 그 단어의 분포 수준을 보여준다.

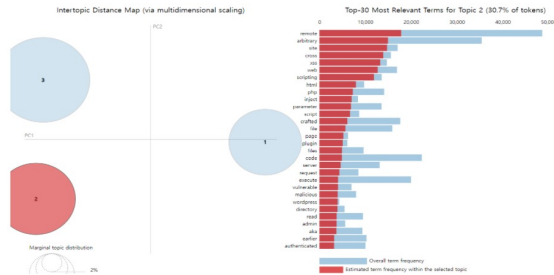


Fig. 7 Visualization by using pyLDAvis

3개의 각 토픽 그룹별 키워드를 리스트화하여 해당 그룹과 CVE 단어 목록 간의 유사도를 자카드 기법으로 측정한다. 각 CVE 별로 Topic 그룹들과의 유사도 값들을 모두 더해 최종값(L)으로 반영한다.

$$L = \sum_{n=1}^3 \mathcal{J}(C, T_n) \quad (3)$$

- C : CVE Word List
- T : Topic List according to LDA modeling

4.3. MITRE ATT&CK과 CVE 유사도 (3단계)

세 번째 식에서는 MITRE ATT&CK의 Tactic과 CVE 데이터 간의 유사도를 확인한다. MITRE ATT&CK 프레임워크 데이터의 동향을 식에 반영하고자 한다.

MITRE ATT&CK에서 제공하는 14개 Tactic과 222개의 Technique 키워드 데이터를 수식어, 부사 등을 제외한 단어별로 쪼개고, Tactic 별 중복 단어를 제외한 모든 단어를 포함하여 14개의 리스트를 생성한다. 14개 그룹별 키워드와 CVE 단어 리스트 간의 유사도를 자카드 기법으로 측정한다. 즉, 각 CVE 리스트별 각 Tactic 리스트와 비교하는 과정을 14번 거치게 된다. CVE 별 14개의 Tactic 그룹과의 유사도를 확인한 값들을 모두 누적하여 합산해 최종값(M)으로 반영한다.

$$M = \sum_{n=1}^{14} \mathcal{J}(C, K_n) \quad (4)$$

- C : CVE Word List
- K : Word List according to MITRE ATT&CK's Tactic and Technique

4.4. 최종 위험 산출식

최종 위험 산출식은 아래와 같다. 본 연구의 스코어링 산출식은 전체 취약점 동향을 반영하여 위험도를 평가하고자 하였다. 위의 3단계 과정에 따른 식 (2), (3), (4)인 W, L, M 각각에 가중치(α, β, γ)를 부여하고 모두 합산하여 최종 취약점 위험 스코어링 산출식(CTRS)을 도출하였다. W는 CVE와 Vendor, Attack Type과의 연관성, L은 전체 CVE의 토픽과의 연관성, M은 MITRE Tactic과의 연관성을 나타낸다. 각 식에 대한 영향도를 시간 흐름(동향)에 따른 기존 CVE 스코어와 비교하기 위해 α, β, γ 를 0부터 1까지로 설정한 결과, 취약점 위험도에 대해 가장 순위를 잘 표현한 비중은 0.6666:1:1 이었다. 이는 기존에 없는 플랫폼(또는 Vendor)을 대상으로 이뤄지는 공격이나 공격 유형에 대한 위험도를 가장 잘 나타내었으며, 결과적으로 W에는 1의 가중치를 부여하였고, 동향을 효과적으로 반영하고 있는 식 L과 M에는 1.5의 가중치를 부여하였다.

$$CTRS = W \times \alpha + L \times \beta + M \times \gamma \quad (5)$$

$$= \left(\frac{\sum_{n=1}^{16} V_n}{V_{count}} + \frac{\sum_{n=1}^{40} A_n}{A_{count}} \right) \times \alpha + \sum_{n=1}^3 \mathcal{J}(C, T_n) \times \beta + \sum_{n=1}^{14} \mathcal{J}(C, K_n) \times \gamma$$

4.5. 스코어링 활용 및 검증

본 논문에서의 최종 위험 산출식을 실제 공개 데이터에 적용해 그 활용도를 확인해보고자 한다. APT 공격 스코어링 기법을 제안한 기존 연구[14]에서는 제안한 방법을 최종적으로 APT 공격 사례 보고서에 적용하였으며, 이를 통해 여러 사이버 공격의 위험 수준과 시의성을 판단할 수 있는 가능성을 보여주었다.

4.5.1. 검증 데이터 및 방법

국내외 다양한 기관, 연구소 및 사이트 등에서 CVE와 위험도 높은 사이버 위협에 대한 취약점과 분석 정보들을 제공하고 있다. 2007년부터 15년간 보안과 관련된 정보를 제공하는 사이트인 'The Hacker News'의 기사 내용을 검증 데이터로 활용하였다. 본 연구에서는 위와 같이 공신력 있는 사이트에서 수집한 취약점 데이터를 토대로 최종 위험 산출 스코어링을 수행한다.

먼저, 취약점 정보 제공 사이트에서 수집한 요약 데이

터와 관련도 높은 취약점 리스트 선정을 위하여 여러 문장을 단어로 쪼개어 저장한다. 불용어(Stopwords) 목록을 불리와 제외할 단어로 활용한다. 문장들을 쪼갠 단어들 리스트 형태로 저장하고 중복 단어를 제거한다. 다음으로 리스트 데이터와 관련성 높은 CVE 번호를 확인하기 위하여 자카드 유사도 기법을 활용하였다. 리스트 목록과 CVE 데이터 간의 유사도 계산 결과, 값이 큰 상위 20여 개(전체 CVE 데이터의 0.001%)의 CVE 번호를 추출한다. 선정된 CVE 번호 20여 개에 대해 위협 산출식에 따른 위험도 점수를 확인하여 다시 위험도 점수 기준 내림차순으로 정리한다. 최종 확인된 CVE 번호 관련 내용이 해당 사이트 또는 기사에 존재하는 경우, 간단한 정보 데이터만으로 위협 산출식을 통해 위험도 높고 관련성 깊은 CVE를 확인할 수 있다.

4.5.2. 분석 결과

발표된 특정 취약점 위협 정보(Zoho ManageEngine Vulnerability)인 그림 8과 같은 데이터[15]에 대해 본 논문에서 산출한 취약점 위험도 스코어링 식을 적용하였다.

The U.S. Cybersecurity and Infrastructure Security Agency (CISA) on Thursday added a recently disclosed security flaw in Zoho ManageEngine to its Known Exploited Vulnerabilities (KEV) Catalog, citing evidence of active exploitation.

"Zoho ManageEngine PAM360, Password Manager Pro, and Access Manager Plus contain an unspecified vulnerability which allows for remote code execution," the agency said in a notice.

Fig. 8 Vulnerability information by 'Hacker News'

세 단계 과정을 거쳐, 표 3과 같은 결과를 확인하였다. 취약점 위협 정보에 대한 관련도와 위험도가 높은 상위 5개의 CVE와 해당 CVE의 CVSS 점수, 최종 취약점 위험도 점수를 보여준다.

Table. 3 Result of applying threat calculation scoring

No.	CVE Number	CVSS	1st	2nd	3rd	Final Score
1	CVE-2020-9496	6.1	0.302	0.421	0.581	180.5
2	CVE-2020-7199	9.8	0.342	0.285	0.448	144.15
3	CVE-2021-39179	8.8	0.381	0.195	0.443	133.8
4	CVE-2021-42099	9.8	0.398	0.216	0.261	111.35
5	CVE-2021-44515	9.8	0.435	0.145	0.231	99.9

가장 높은 점수를 보여주고 있는 CVE-2020-9496는 사이트 내에서 설명하고 있는 위협 정보와 가장 관련 있

는 취약점 번호이다. CVSS점수(심각도에 따라 0.0 ~ 10.0까지 분포)는 6.1점(Medium)이지만, 실제 발표된 정보 확인 결과, 위협 정보에 해당하는 최근 취약점 번호인 CVE-2022-35405는 이전에 발표된 XML RPC 역직렬화(Deserialization) 취약점인 CVE-2020-9496과 동일한 특정 클래스에 의해 XML 데이터가 역직렬화되어 원격 명령이 실행될 수 있는 취약점이다. 취약 대상, 공격 방식 등의 주요 포인트가 유사하고 실제 해당 취약점 패치 이후 유사한 공격 포인트를 통해 다시 활성화된 취약점으로 확인되었다. 그 외에 두 번째부터 다섯 번째 사이에 랭크된 취약점 4개도 모두 적용 위협 정보와 관련성이 높은 비슷한 유형의 취약점이며, 동일한 취약 대상 제품에 대한 취약점도 포함된 것을 확인할 수 있다.

이 외에도 여러 보안 취약점 관련 뉴스 데이터에 대해 본 논문의 위험도 스코어링 기법을 적용하고 관련 CVE 정보를 확인하였다. 4개의 예시 데이터[16-19] 적용 결과에서 확인한 유의미한 정보들, 즉, 관련도와 위험도 높은 상위 5개의 취약점 번호, 뉴스 데이터와의 유사도, 스코어링 시스템 점수를 표 4에 나타내었으며, 최종 점수 기준으로 내림차순 정리하였다.

첫 번째 뉴스 데이터는 스프링 자바 프레임워크 관련 취약점 정보이다. 실제 데이터 분석 결과, 과거에 있었던 CVE-2010-1622 취약점과 비슷하게 특정 객체에 접근할 수 있게 되어 발생한 취약점이었으며 관련 CVE 번호가 확인된다. 또한 스프링 프레임워크 DoS 취약점과 클라우드 자원 접근 취약점 CVE도 확인되었다.

두 번째 뉴스 데이터는 자바 기반의 팀 협업 소프트웨어인 Atlassian Confluence 관련 제로데이 취약점 정보이다. 해당 데이터 위험 산출 과정 적용 결과와 실제 데이터 분석 결과, 해당 취약점은 과거의 CVE-2021-26084와 유사한 취약점으로 확인되었다. 또한 자바 클래스의 메소드를 실행하는 오픈소스 표현식 언어인 OGNL을 악용하는 과거의 취약점과도 관련이 있다는 것을 확인할 수 있다.

세 번째 데이터 역시, 뉴스 데이터에 해당하는 취약점과 관련 있는 취약점 결과를 보여준다. 네트워크 접근 권한이 있는 악의적 공격자가 서버 측 템플릿 주입을 트리거 할 수 있는 방법의 취약점 정보들을 확인할 수 있다.

Table. 4 Result of applying threat calculation scoring

Vulnerability News Data			
1	Additional details of the flaw, dubbed SpringShell and Spring4Shell, have been withheld to prevent exploitation attempts and until a fix is in place by the framework maintainers, a subsidiary of VMware……		
	CVE Number	Similarity	Final score
	CVE-2010-1622	0.895	138.985
	CVE-2022-22963	0.923	137.488
	CVE-2022-22950	0.908	122.763
	CVE-2014-0112	0.765	102.564
2	Vulnerability News Data		
	Atlassian has warned of a critical unpatched remote code execution vulnerability impacting Confluence Server and Data Center products that it said is being actively exploited in the wild. Alternatively, it has recommended implementing a web application firewall (WAF) rule which blocks URLs containing "\${" to reduce the risk……		
	CVE Number	Similarity	Final score
	CVE-2022-26134	0.813	116.499
	CVE-2021-26084	0.711	109.666
	CVE-2021-31805	0.858	101.915
3	Vulnerability News Data		
	The security shortcoming relates to a remote code execution vulnerability that stems from server-side template injection in VMware Workspace ONE Access and Identity Manager. "A malicious actor with network access can trigger a server-side template injection that may result in remote code execution," the company noted in its advisory……		
	CVE Number	Similarity	Final score
	CVE-2022-22954	0.886	221.323
	CVE-2017-17215	0.753	175.324
	CVE-2020-3233	0.781	136.451
4	Vulnerability News Data		
	"A local, authenticated attacker could gain elevated local system or administrator privileges through a vulnerability in the Win32k.sys driver," the Windows maker said. "These types of vulnerabilities are a frequent attack vector for malicious cyber actors of all types and pose……		
	CVE Number	Similarity	Final score
	CVE-2022-21882	0.852	180.325
	CVE-2021-1732	0.795	152.680
	CVE-2019-1732	0.652	120.082
CVE-2022-21882	0.631	105.675	
CVE-2020-0668	0.598	99.683	

네 번째 데이터는 윈도우 취약점 관련 정보이다. Win32k.sys 드라이버의 취약점을 악용한 로컬 시스템 권한 상승 취약점으로 확인되었고, 높은 관련성과 위험도 결과로 확인된 CVE-2021-1732는 과거의 동일한 모듈의 권한 상승 취약 부분을 우회하여 발생한 취약점이다. 그 외 확인된 취약점들은 비슷한 유형의 커널 권한 상승 및 검증 미흡과 관련된 취약점이었다.

표 4에 제시한 뉴스 데이터들은 2022년 상반기에 이슈가 되었던 취약점과 관련 있는 내용으로 최신의 취약점이며, 알려진 공격 코드를 확인할 수 없었던 취약점에 대해서도 관련 취약점 정보를 확인할 수 있었다. 실제 새로운 취약점은 이전 취약점에서 파생되어 새로 발표되는 경우가 많기 때문에 취약점 관련 일부 데이터만으로 사전에 위험도를 확인하거나 관련 취약점 정보를 얻어 대응 및 인지에 활용할 수 있을 것으로 예상된다.

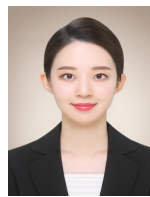
V. 결론

본 연구에서는 MITRE ATT&CK과 토픽 모델링을 활용한 동향을 반영하여 보안 CVE 취약점의 위험도를 스코어링 하는 방법을 제안하였다. 취약점 자체 정보뿐만 아니라, 표준 모델로 자리 잡아가고 있는 MITRE ATT&CK 프레임워크의 Tactic과 Technique 요소들을 이용하였고, CVE 전체 데이터 대상을 LDA 토픽 모델링 알고리즘을 이용해 전체 동향을 반영한 취약점의 위험도 기준을 제시하였다. 제안한 스코어링 방법을 통해 위협 데이터 일부만으로 관련 있고 위험도 높은 취약점을 확인할 수 있으며, 보안 취약점에 대한 간편한 분석과 관련 취약점을 미리 확인함으로써 신속한 대응의 방안으로 활용할 수 있다.

현재도 보안 취약점은 지속해서 발표되고 있고, 그에 따라 위협 산출식 기본 데이터 정보도 계속해서 누적되고 있다. 본 연구 개선을 위하여 해당 취약점 정보 변동 사항 및 추가되는 데이터의 전처리 과정에 대한 자동화 작업이 필요하며, 데이터가 쌓여 나갈수록 위험도 결과는 더 정교하고 세밀해져 취약점 위험 스코어링에 있어 더욱 활용도가 높아질 수 있을 것으로 기대된다.

References

- [1] CVE Details. Browse Vulnerabilities By Date [Internet]. Available: <https://www.cvedetails.com/browse-by-date.php>.
- [2] MITRE ATT&CK. Enterprise tactics [Internet]. Available: <https://attack.mitre.org/tactics/enterprise/>.
- [3] Towards Data Science. Topic Modeling in Python: Latent Dirichlet Allocation (LDA) [Internet]. Available: <https://towardsdatascience.com/end-to-end-topic-modeling-in-python-latent-dirichlet-allocation-lda-35ce4ed6b3e0>.
- [4] LearnDataSci. Jaccard Similarity [Internet]. Available: <https://www.learndatasci.com/glossary/jaccard-similarity/>.
- [5] H. Chen, J. Liu, R. Liu, N. Park, and V. S. Subrahmanian, “VEST: A System for Vulnerability Exploit Scoring & Timing,” in *Proceeding of the Twenty-Eighth International Joint Conference on Artificial Intelligence*, Macao, China, pp. 6503-6505, 2019. DOI: 10.24963/ijcai.2019/937.
- [6] M. C. Kim, S. Oh, H. Kang, J. Kim, and H. K. Kim, “Risk Scoring System for Software Vulnerability Using Public Vulnerability Information,” *Journal of the Korea Institute of Information Security & Cryptology*, vol. 28, no. 6, pp. 1449 - 1461, Dec. 2018. DOI: 10.13089/JKIISC.2018.28.6.1449.
- [7] H. H. Jin and H. K. Kim, “A Study on Web Vulnerability Risk Assessment Model Based on Attack Results: Focused on Cyber Kill Chain,” *Journal of the Korea Institute of Information Security & Cryptology*, vol. 31, no. 4, pp. 779 - 791, Aug. 2021. DOI: 10.13089/JKIISC.2021.31.4.779.
- [8] C. Park, “Vulnerability Risk Computation method for Attack Graph Generation,” in *Proceedings of the Korean Institute of Communication Sciences Conference*, Jeju, Korea, pp. 123-124, 2016. DOI: 10.17662/ksdim.2020.16.1.079.
- [9] H. -G. Chae, G. -H. Lee, and J. -Y. Lee, “Analysis of Domestic and Foreign Financial Security Research Activities and Trends through Topic Modeling Analysis,” *Journal of the Korea Industrial Information Systems Research*, vol. 26, no. 1, pp. 83 - 95, Feb. 2021. DOI: 10.9723/jksis.2021.26.1.083.
- [10] S. U. Lee and J. Lee, “Topic Modeling to Identify Cloud Security Trends using news Data Before and After the COVID-19 Pandemic,” *Convergence Security Journal*, vol. 22, no. 2, pp. 67 - 75, Jun. 2022. DOI: 10.33778/kcsa.2022.22.1.067.
- [11] H. S. Choi, W. S. Lee, and S. Y. Sohn, “Analyzing Research Trends in Personal Information Privacy Using Topic Modeling,” *Computers & Security*, vol. 67, pp. 244-253, Jun. 2017. DOI: 10.1016/j.cose.2017.03.007
- [12] The MITRE Corp. Download CVE List [Internet]. Available: <https://cve.mitre.org/data/downloads/index.html>.
- [13] Pythonspot. NLTK stop words [Internet]. Available: <https://pythonspot.com/nltk-stop-words/>.
- [14] S. Cho, Y. Park, K. Lee, C. Choi, C. Shin, and K. Lee, “An APT Attack Scoring Method Using MITRE ATT&CK,” *Journal of the Korea Institute of Information Security & Cryptology*, vol. 32, no. 4, pp. 673 - 689, Aug. 2022. DOI: 10.13089/JKIISC.2022.32.4.673.
- [15] The Hacker News. CISA Warns of Hackers Exploiting Recent Zoho ManageEngine Vulnerability [Internet]. Available: <https://thehackernews.com/2022/09/cisa-warns-of-hackers-exploiting-recent.html>.
- [16] The Hacker News. Unpatched Java Spring Framework 0-Day RCE Bug Threatens Enterprise Web Apps Security [Internet]. Available: <https://thehackernews.com/2022/03/unpatched-java-spring-framework-0-day.html>.
- [17] The Hacker News. Hackers Exploiting Unpatched Critical Atlassian Confluence Zero-Day Vulnerability [Internet]. Available: <https://thehackernews.com/2022/06/hackers-exploiting-unpatched-critical.html>.
- [18] The Hacker News. Critical VMware Workspace ONE Access Flaw Under Active Exploitation in the Wild [Internet]. Available: <https://thehackernews.com/2022/04/vmware-releases-patches-for-critical.html>.
- [19] The Hacker News. CISA Orders Federal Agencies to Patch Actively Exploited Windows Vulnerability [Internet]. Available: <https://thehackernews.com/2022/02/cisa-orders-federal-agencies-to-patch.html>.



임지혜(Jihye Lim)

2019년 8월 한국외국어대학교
컴퓨터전자시스템공학부 학사
2021년 3월 ~ 현재 중앙대학교 융합보안학과
정보보안 전공석사과정
※관심분야 : Cyber Threat Intelligence,
정보보안, 보안 취약점 대응



이재우(Jaewoo Lee)

2006년 2월 서울대학교 컴퓨터공학부 학사
2008년 2월 서울대학교 컴퓨터공학부 석사
2017년 8월 University of Pennsylvania, Ph.D in
Computer and Information Science
2018년 3월~현재 중앙대학교 산업보안학과
조교수
※관심분야 : Cyber Physical System Security