

UNIVERSITY OF OTAGO

DEPARTMENT OF COMPUTER SCIENCE

COSC PROJECT REPORT

An Exploration of the Common Vulnerability Scoring System

Author:

Jake NORTON (5695756)

Supervisor(s):

Dr. David EYERS

Dr. Veronica LIESAPUTRA

October 10, 2024



Abstract

The Common Vulnerability Scoring System (CVSS) is designed to produce scores for software vulnerabilities, serving as a crucial tool for triaging the increasing number of new vulnerabilities released each year. As manual scoring cannot keep pace with this growth, there is a need for automated prediction methods. While machine learning, particularly large language models (LLMs), have shown promise in predicting CVSS metrics given their descriptions. This paper aims to explore the CVE and CVSS processes to provide a foundation for future research. This paper explores the potential of using multiple data sources for cross-validation and increased training data volume, specifically comparing NVD with the MITRE database. The analysis reveals differences in how these databases rate Common Vulnerabilities and Exposures (CVEs), highlighting the challenges in establishing a consistent ground truth for CVSS scores. While the combination of multiple datasets did not yield the expected improvements in prediction accuracy, this investigation provides valuable insights into the complexities of CVSS scoring across different evaluation teams. The next focus of this paper is on enhancing the interpretability and analysis of CVE/CVSS data through clustering techniques and other dataset analysis methods. I employ Latent Dirichlet Allocation (LDA) and other clustering approaches to uncover patterns within the vulnerability descriptions. The findings demonstrate that certain types of vulnerabilities tend to cluster together, such as cross-site scripting, denial-of-service attacks.

1 Introduction

In 2023, 29,065 new vulnerabilities were recorded, as illustrated in Figure 1. This number continues to rise year on year, underscoring the growing challenge of managing and prioritizing software vulnerabilities. These vulnerabilities are documented using the Common Vulnerabilities and Exposure system (CVE [27]), with CVSS scores derived from these entries to indicate severity and impact. The National Vulnerability Database (NVD [33]) serves as the primary source for CVSS-enriched CVE data and is widely used in research (e.g., [10, 1, 5]). However, to explore the potential benefits of multiple data sources, we also investigated the MITRE database [30], which is the main repository for CVEs and includes a significant number of entries with CVSS scores. The initial comparison between NVD and MITRE databases employed a methodology inspired by the paper “Can the Common Vulnerability Scoring System be Trusted? A Bayesian Analysis” [23]. This approach aimed to assess the agreement between different data sources in scoring CVEs and explore the feasibility of establishing a ground truth. The analysis revealed differences in how

these databases rate CVEs (see Figure 2, 5, 4), highlighting the subjective nature of CVSS scoring and the challenges in achieving consistency across different evaluation teams. While the comparison of multiple datasets did not yield the anticipated improvements in CVSS prediction, it provided valuable insights into the variability of human-assigned scores. This variability suggests that automated prediction models may face similar challenges in areas where human evaluators disagree. Given these findings, the focus of the research shifted towards enhancing the interpretability and analysis of the CVE/CVSS data itself. I employed clustering techniques, primarily Latent Dirichlet Allocation (LDA [6]), along with other dataset analysis methods to uncover patterns within the vulnerability descriptions. This approach offers several advantages:

1. It provides a data-driven perspective that complements the LLM-based prediction methods.
2. The clustering techniques offer faster analysis times compared to complex LLM interpretability methods, with LDA processing the entire corpus in approximately one hour.
3. It reveals inherent groupings and relationships within the data that may not be immediately apparent through other analysis methods.

The clustering analysis has yielded interesting insights, such as the tendency for certain types of vulnerabilities to cluster together. For example the analysis of the integrity impact metric revealed meaningful patterns that align with expert knowledge in the field:

- Vulnerabilities classified as having no impact (**NONE**) on data integrity often clustered around denial of service attacks and system crashes. While these are serious issues, they typically do not directly affect data integrity.
- Low impact (**LOW**) on integrity was frequently associated with cross-site scripting vulnerabilities, particularly in WordPress plugins. This classification makes sense as such vulnerabilities can cause localized data integrity issues, often limited to individual user interactions rather than compromising the entire application.
- High impact (**HIGH**) integrity issues clustered around SQL injection and database-related vulnerabilities. This aligns with the severe nature of these attacks, which can directly manipulate or corrupt data at the database level.

The following report will try to answer these questions:

- Is there scoring consistency between MITRE and NVD?
- Can we improve the LLMs classifier performance when we augment the NVD dataset with data from the MITRE dataset?
- What are the important keywords for each category?

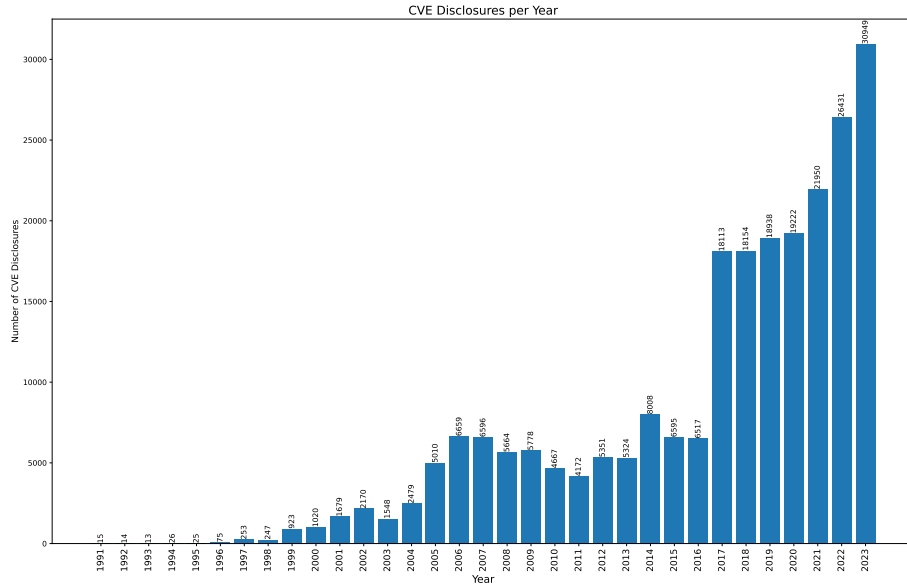


Figure 1: Number of new CVEs by year

2 Background

Vulnerabilities are stored in a consistent system called Common Vulnerabilities and Exposures (CVE [27]).

Here is an example CVE

- Unique Identifier: CVE-2024-38526
- Source: GitHub, Inc.

- Published: 06/25/2024
- Updated: 06/26/2024
- Description: pdoc provides API Documentation for Python Projects. Documentation generated with `pdoc -math` linked to JavaScript files from polyfill.io. The polyfill.io CDN has been sold and now serves malicious code. This issue has been fixed in pdoc 14.5.1.

Sourced from [NVD CVE-2024-38526 Detail \[32\]](#)

This has a unique identifier, which is given by one of the CVE numbering authorities (CNA [29]), such as GitHub, Google or any of these, [CVE list of partners \[9\]](#). The description is the most important part in our case. This should provide information about the vulnerability. What can be exploited (device / software component)? How is the product affected if the vulnerability is exploited? Ideally there would be a part of the description that relates to every metric, unfortunately these descriptions are not necessarily suited to machine learning as the people writing the descriptions are expecting a lot of intrinsic knowledge.

The Common Vulnerability Scoring System

CVSS scoring is a high level way to break up vulnerabilities into different categories. Organisations can use it to choose which vulnerability to focus on first. CVSS is broken up into three distinct sections: base, temporal and environmental scores.

For brevity I will only show the specifics of CVSS 3.1 [34] as this is by far the most commonly used version, even if it is not the most recent.

Base Score

- Attack Vector: Defines the avenues of attack that the vulnerability is open to. The more open a component is, the higher the score. This can have the values: **Network**, **Adjacent**, **Local** and **Physical**.
- Attack Complexity: Describes how complex the attack is to orchestrate. Encompasses questions like, what are the prerequisites? How much domain knowledge / background work is necessary? How much effort does the attacker need to invest to succeed? This can have the values: **Low** or **High**. **Low** gives a higher base score.

- Priviledges Required: The degree of privileges the user needs to complete the attack. Ranging from: **None**, **Low** (e.g. User level privilege), **High** (e.g. Administrator). The lower the privilege the higher the base score.
- User Interaction: Describes if the exploit requires another human user to make the attack possible, E.g. clicking a phishing link. This is either **None** or **Required**, the score is highest when no user interaction is required.
- Scope: Defines if the attack can leak into other security scopes. E.g. access to one machine gives the ability to elevate privileges on other parts of the system. This can take **Unchanged** or **Changed**, the score being highest when a scope change occurs.
- Confidentiality Impact: Detemines what is the impact on the information access / disclosure to the attacker. This can be: **High**, **Low** or **None** with **High** adding the most to the base score.
- Integrity Impact: Refers to the integrity of the information within the component. I.e. could the data have been modified by the attacker. This has: **High**, **Low** or **None**, as categories with **High** adding the most to the base score.
- Availability Impact: Refers to the impact of the attack on the availability of the component. E.g. the attacker taking the component off the network, denying the users access. This can be: **High**, **Low** and **None** with **High** adding the most to the base score.

This is a summarized version of the [3.1 specification document provided by the Forum of Incident Response and Security Teams \(FIRST\)](#) [34].

Temporal

- Exploit Code Maturity: The state of the attack itself, e.g. has this exploit been pulled off in the wild or is it currently academic.
- Remediation Level: Whether the exploit in question has a patch available.
- Report Confidence: The degree of confidence in the CVE report itself, the report may be in early stages where not all of the information is known.

This is a summarized version of the [3.1 specification document](#) provided by the [Forum of Incident Response and Security Teams \(FIRST\)](#) [34].

Temporal metrics would be useful in general for a CVSS score, however NVD do not store these temporal metrics. As far as I can tell there is no reason given for this specifically, though discourse ([Stack exchange post](#)) [4] around the subject suggests that this is due to a lack of verifiable reporting. From my perspective, both remediation level and report confidence feel like they could have scores attributed to them, however finding verifiable reports on the exploits seen in the wild is difficult. There are two relatively new organisations on this front, Cybersecurity & Infrastructure Security Agency (CISA, [public sector](#)) and inthewild.org ([private sector](#) [8]).

2.1 Data Options

I will be using NVD and MITRE as the sources of data. In 2016 when Johnson et al. did their paper on CVSS [23], they had access to five different databases. Unfortunately only two of these remain for modern data. There are others, but they are either in archival or proprietary status.

2.1.1 National Vulnerability Database

The National Vulnerability Database is the defacto standard dataset used for CVSS generation research [10, 1, 5]. This makes a lot of sense as it is built for the purpose with a team dedicated to enriching CVEs with CVSS scores. The dataset I am using was retrieved using the NVD API in March 2024 and contains ~100000 CVEs enriched with CVSS scores. This comes in a consistently formatted JSON dump.

2.1.2 MITRE Database

MITRE is the defacto database for the storage of CVEs themselves, their database contains ~40000 CVEs enriched with CVSS 3.1 scores. These are in a JSON dump retrieved in March 2024. The format for usage is a bit more cumbersome to use. The CVSS scores are only stored as CVSS vector strings (a simple text encoding [36]). These are not hard to parse, though they are stored slightly different between versions, as well as sometimes being inconsistent (~5000 had temporal metrics within the vector strings in the MITRE database).

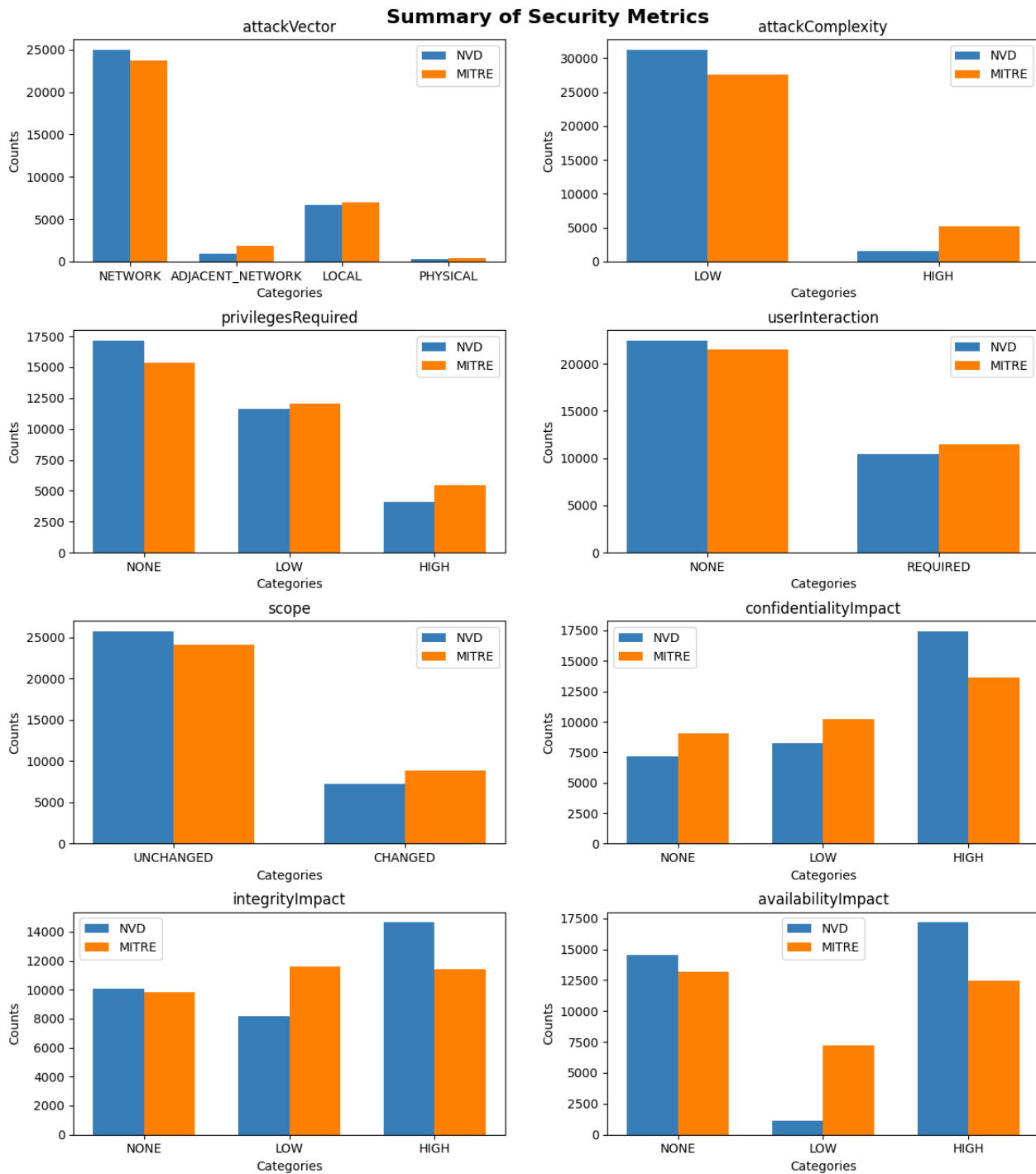


Figure 2: Comparison of CVSS ratings between MITRE and NVD

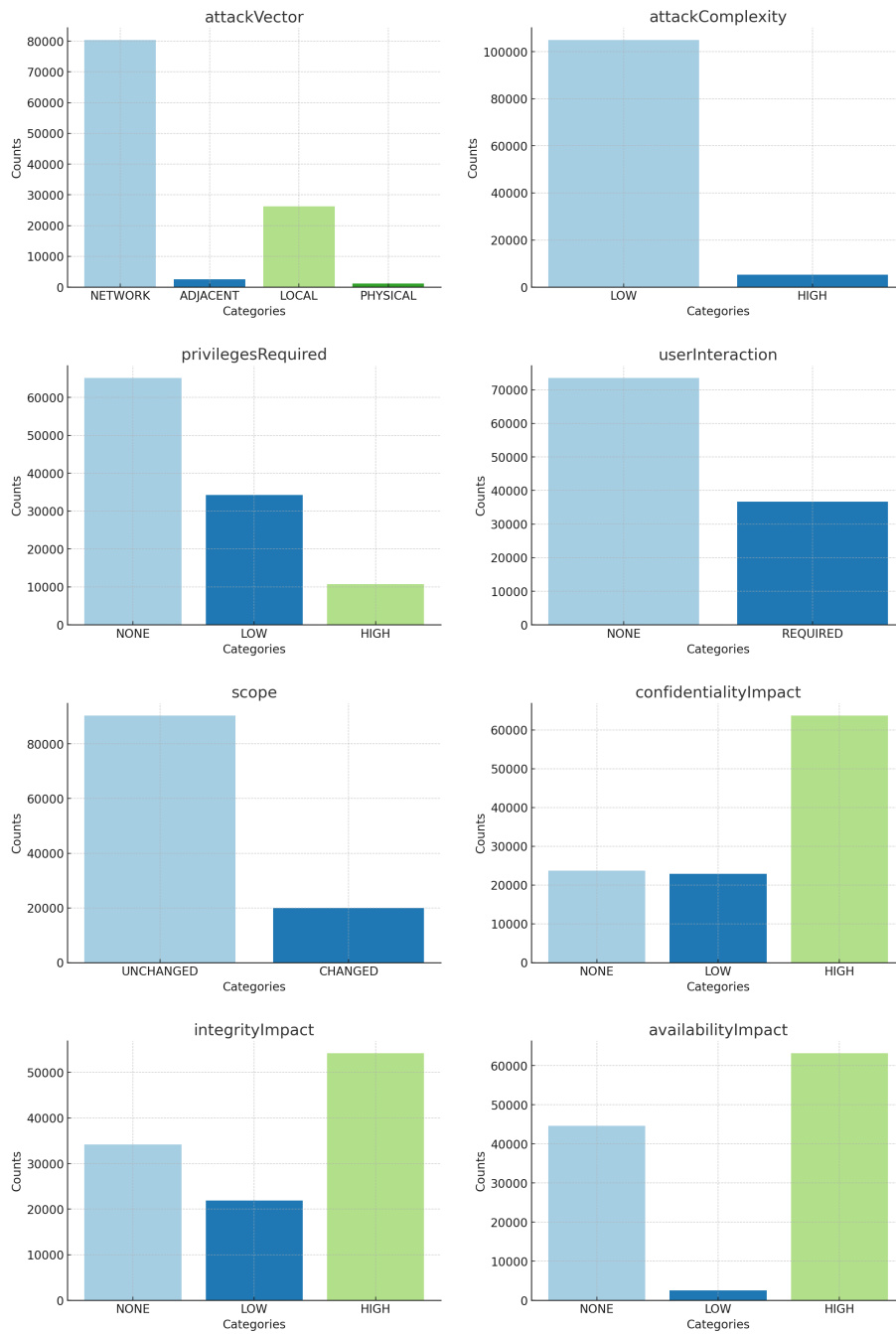


Figure 3: Distribution of Metric ratings for NVD

2.2 Evolution of CVSS and Its Identity Crisis

Possibly need some kind of link between these...

When CVSS 2.0 was released, it was promoted as a framework to help IT management prioritize and remediate vulnerabilities posing the greatest risk. The initial goal was to provide a comprehensive method for assessing risk, as indicated by its original documentation:

“Currently, IT management must identify and assess vulnerabilities across many disparate hardware and software platforms. They need to prioritize these vulnerabilities and remediate those that pose the greatest risk. But when there are so many to fix, with each being scored using different scales, how can IT managers convert this mountain of vulnerability data into actionable information? The Common Vulnerability Scoring System (CVSS) is an open framework that addresses this issue.” [18]

However, by the time CVSS 3.1 was released, the framework’s focus had shifted, partly due to complaints about CVSS being a poor judge of risk. The authors stated: *“CVSS measures severity, not risk.” [34]*

The identity crisis is a problem because the original stance, that it can be used as a primary prioritisation tool, has lured parts of the industry into doing just that. As mentioned by Henry Howland in [17], there are many large, mainly US based places mandating the sole use of CVSS base score for remediation, such as the Payment Card Industry Data Security Standard [2], the Department of Defense Joint Special Access Program implementation Guide [11] to name a few.

This change in stance has created confusion about the true purpose of CVSS.

2.2.1 CVSS Formula

How CVSS is computed under-the-hood is confusing at best. CVSS 3.1 is not explained to the same depth as version 2.0, but my understanding is that it followed a similar process. This is that process summarised from [CVSS version 2 FAQ](#): [14]

1. Divide the six metrics into two groups:
 - **Impact** (3 metrics)
 - **Exploitability** (3 metrics)
2. Create sub-vectors for each group:

- **Impact sub-vector:** 27 possible values (3^3)
 - **Exploitability sub-vector:** 26 possible values ($3^3 - 1$ for no impact)
3. Develop and apply a series of rules to order the sub-vectors. These rules are primarily based on the severity of components, e.g. vectors with more **Complete** components are rated higher.
 4. Assign scores to the ordered sub-vectors based on the derived rules and review by the Special Interest Group (SIG).
 5. Apply weightings to the sub-vectors:
 - **Impact:** 0.6
 - **Exploitability:** 0.4
 6. Develop a formula that approximates the scores derived from the ordered sub-vectors. Ensure the formula produces scores with ± 0.5 error from the originally defined vector score and does not exceed the maximum value of 10.
 7. Test and refine the formula through iterations, ensuring it aligns with desired values and corrects any issues, such as scores exceeding 10.

This process is inherently inaccurate, it is not a system designed to give precise scores. If we look at 6 above, the formula ([Appendix .2](#) shows the CVSS 3.1 base score formula for reference) which produces the score, does not match exactly the experts decision. There is a lot of rounding and approximation going on. This is designed to make a system which is easy to use and quick to complete by security professionals. There is a space for CVSS, however this along with the other mentioned reasons outlines the issues with using CVSS as a sole metric for prioritization. Perhaps an option is to triage the large swathes of new vulnerabilities coming in with an initial CVSS score, then move on to a deeper dive. This could be in the form of the extra CVSS metrics (Temporal score & Environmental score), or a look into other potential options like the Exploit Prediction Scoring System (EPSS [\[19\]](#)).

3 Is there scoring consistency between MITRE and NVD?

3.1 Related Work

The assessment and comparison of vulnerability scoring systems have been subjects of considerable research in recent years. This section outlines key studies that have contributed to our understanding of vulnerability scoring consistency, database reliability, and automated assessment techniques.

Comparative Analysis of Vulnerability Databases

The primary work most closely aligned with our Bayesian analysis between MITRE and NVD is the study by Johnson et al. titled "Can the Common Vulnerability Scoring System be Trusted? A Bayesian Analysis"[23]. Conducted in 2016, this research investigated the state of CVSS databases and their accuracy. Johnson et al. found that NVD was the most reliable database and concluded that the Common Vulnerability Scoring System could be trusted overall, as scorers rated CVEs consistently. Their methodology forms the foundation for the database comparison in Section 3.3 of our study, although the specific results may have evolved since their publication.

Dong et al.[13] conducted a comprehensive study on detecting inconsistencies in public security vulnerability reports. Their work, which analyzed vulnerability information from MITRE, NVD, and other sources, highlighted the prevalence of information inconsistencies and developed an automated system to detect such inconsistencies. This research underscores the importance of maintaining accurate and consistent vulnerability information across different databases and reports.

Automated Vulnerability Assessment Techniques

Jiang and Atif's work, "An Approach to Discover and Assess Vulnerability Severity Automatically in Cyber-Physical Systems" [22], presents an innovative pipeline for vulnerability assessment. Their approach utilizes multiple data sources and employs a majority voting system to mitigate the impact of incorrectly scored CVEs. While their methods differ from our focus on database comparison, their work underscores the value of leveraging multiple data sources in vulnerability assessment.

Consistency and Reliability of CVSS Scoring

Allodi and Massacci [3] conducted a critical examination of CVSS scoring in their paper "Comparing Vulnerability Severity and Exploits Using Case-Control Studies." Their research questions the correlation between CVSS scores and the actual exploitation of vulnerabilities in the wild, raising important considerations about the practical implications of vulnerability scoring systems.

Furthermore, Spring et al. [42] investigated the reliability of CVSS scoring across different human raters in their study "Towards Improving CVSS." They found significant variations in how different individuals interpret and apply CVSS scoring guidelines, highlighting the challenges in achieving consistent vulnerability assessments even within a standardized framework.

Summary

These studies collectively demonstrate the ongoing challenges and advancements in vulnerability scoring and database reliability. While Johnson et al. [23] provided a foundational analysis of CVSS trustworthiness, subsequent research has explored automated techniques, cross-database comparisons, and the practical implications of vulnerability scores. This study builds upon this body of work, focusing specifically on the consistency between MITRE and NVD scoring using updated data and methodologies.

3.2 Preliminary Data exploration

The scorers for both NVD and MITRE do rate CVEs reasonably similar, one pattern you can see as shown by Figure 2, is that NVD generally give the most common categorical output more ratings. They are less spread out across the full range of values. In addition, if we look at the **Attack Complexity** metric, there is a reasonably large difference in how they are rated, MITRE rate a lot more of the metrics with a **Low** score. This points to some of the difficulty with this kind of rating system, while in theory there is a true value for these metrics, it requires knowledge of the whole space around each of the vulnerabilities, this knowledge will always vary marker to marker. The model will not have direct access to this knowledge; however, it is hoped that it will be able to trace relationships between the different vulnerabilities and learn this intrinsically.

3.3 Hierarchical Bayesian Model

Analysis between the NVD [33] and MITRE [30] databases is conducted using a hierarchical Bayesian model. This model type is particularly suitable when the population is expected to share similarities in some respects while differing in others. In this context, while the databases share common knowledge about vulnerabilities, they differ in the experience of the individuals rating the metrics [23].

The model builds upon the original framework presented in Section 4.1 of [23]. It assumes the existence of a true value for each CVSS dimension, acknowledging that the database sample may deviate from this true value. These potential inaccuracies are represented using confusion matrices (see Equation 1).

A key distinction from the original model is the variability in the number of categorical choices for each CVSS metric. While the original model consistently used three variables for each CVSS metric, my adapted model accommodates between two to four categorical choices, depending on the specific CVSS dimension being evaluated.

3.3.1 Confusion Matrix

Below is an example of the confusion matrix for the CVSS dimension **Confidentiality Impact**:

$$\Pi_{ci} = \begin{bmatrix} \pi_{nn} & \pi_{nl} & \pi_{nh} \\ \pi_{ln} & \pi_{ll} & \pi_{lh} \\ \pi_{hn} & \pi_{hl} & \pi_{hh} \end{bmatrix} \quad (1)$$

where π_{nn} denotes the probability that the database correctly assigns the score **None** when the actual value is indeed **None**. Conversely, π_{nl} and π_{nh} represent instances where **None** was not the true value, indicating an incorrect score assignment by the database.

3.3.2 Prior Distributions

For categorical variables, we employ uninformative priors using Dirichlet distributions[16]. This approach provides a uniform prior over the probability space for all categorical options, minimizing the influence of prior beliefs on the results. The impact of these priors is negligible for categorical metrics with more options than the number of categories.

The confusion matrices also utilize Dirichlet distribution priors. However, we incorporate a slight initial belief to reflect that the individuals producing scores are

not acting entirely randomly and are likely to be correct more often than not. These priors remain weak given the thousands of observations in our dataset.

3.3.3 Parameter Estimation

Our parameter estimation follows the Bayesian approach, updating prior beliefs with observations to produce posterior beliefs. We employ Markov Chain Monte Carlo (MCMC [21]) methods, which allow for simulating data based on the previously created distribution by sampling values that the model believes are likely from the target distribution. This technique enables accurate sampling without requiring all data points.

For implementation, we utilized the PyMC library [35], which fulfills the same functions as JAGS in the original paper [23]. PyMC facilitates the modeling process and provides robust tools for Bayesian inference and MCMC sampling.

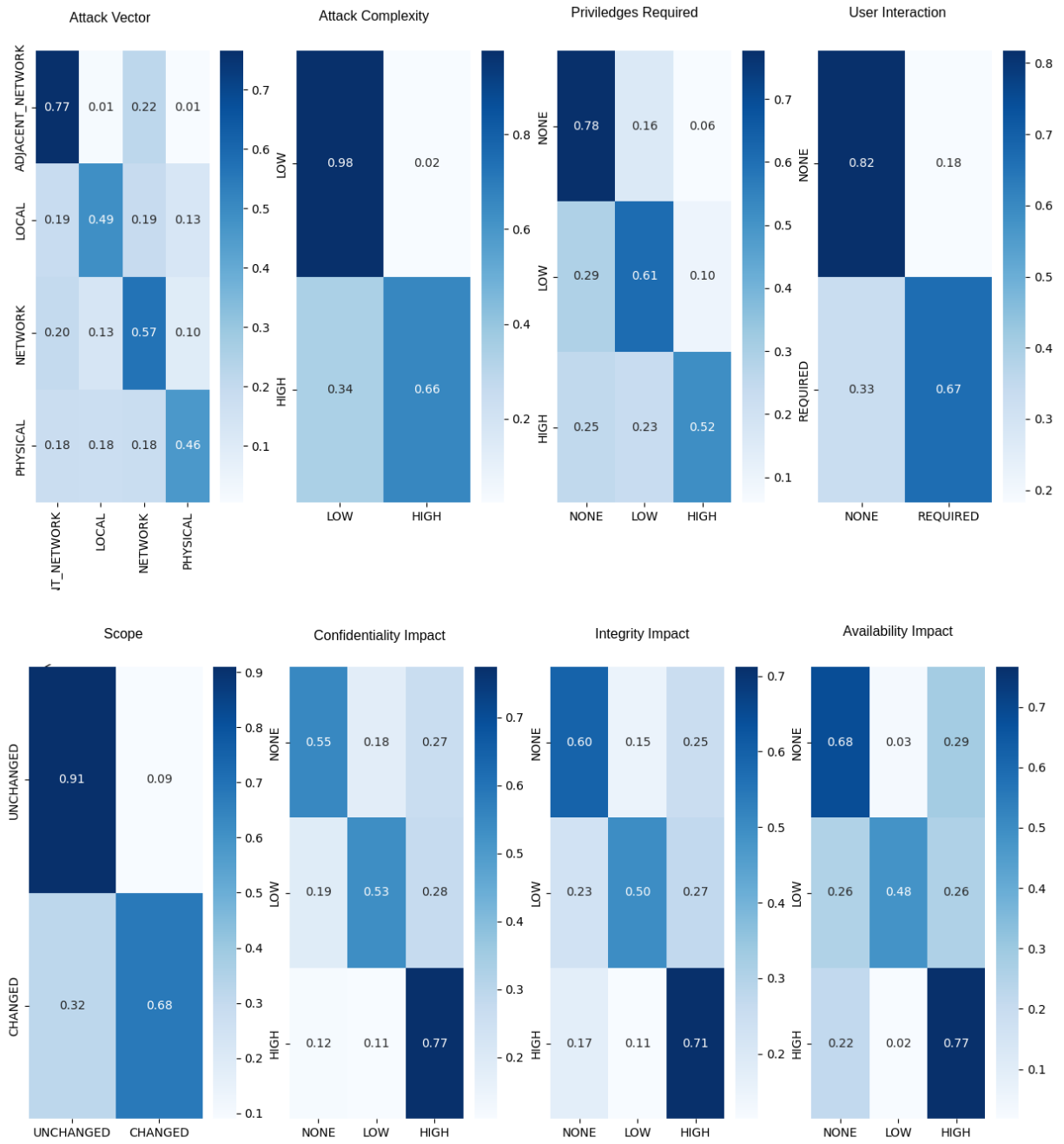
This methodology provides a comprehensive framework for analyzing and comparing vulnerability scoring across different databases, accounting for the inherent uncertainties and variations in expert assessments.

3.3.4 Bayesian Analysis Results

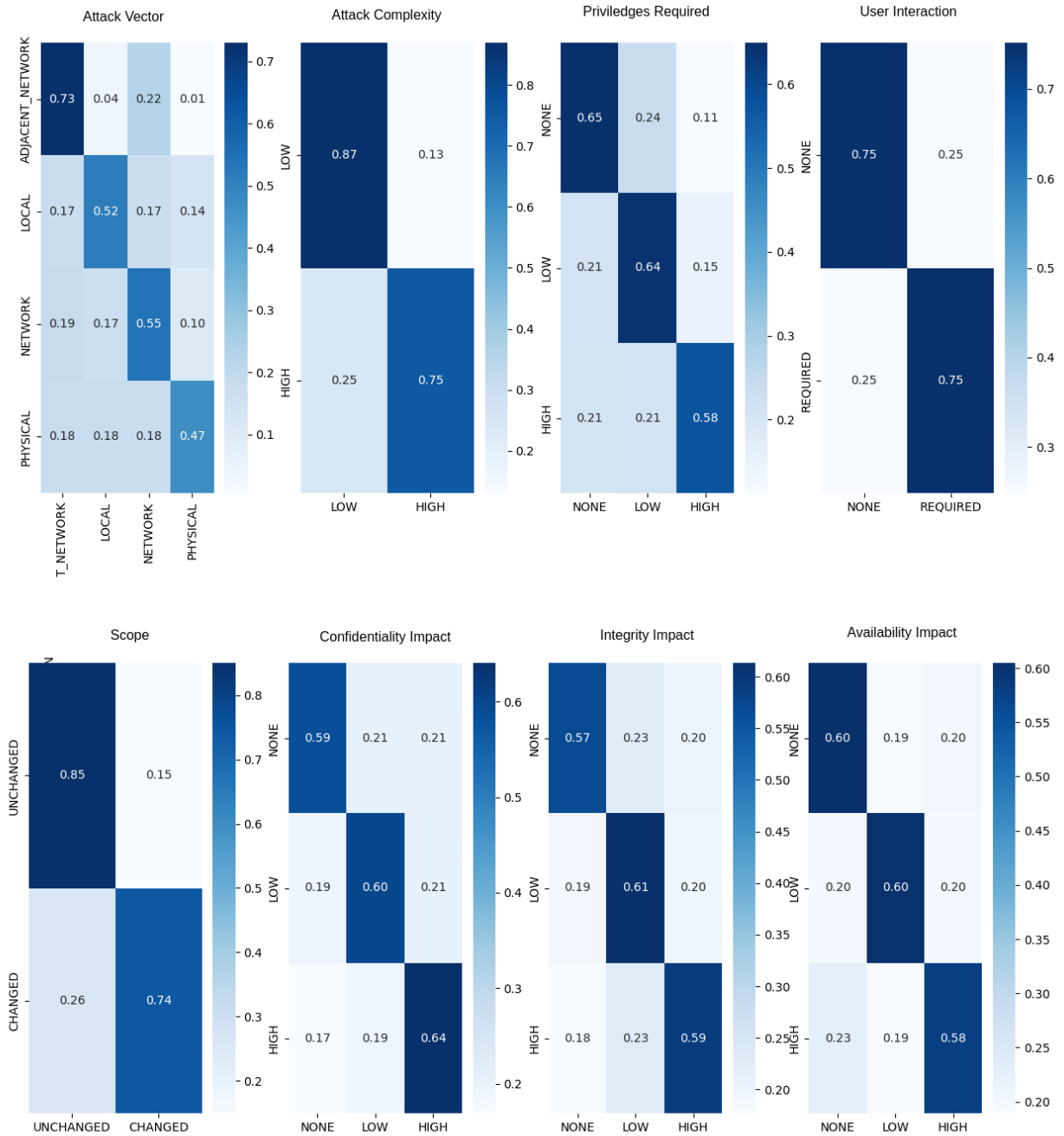
The results for the database analysis will be shown through the confusion matrices for the estimated accuracy of both NVD and MITRE. Unfortunately it is difficult to compare my results to the original paper [23] as they are on a totally different dataset. However, I will note that in general the estimated accuracy for both datasets is much lower than the scores from the original study. Across the board NVD often had $\sim 90\%$ accuracy for the highest accuracy field of any metric, with Confidentiality as a clear outlier as seen from Table 1.

Some general trends are that NVD (see Figure 4) have more extreme estimated accuracies. They do better for the higher frequency options, for the **Low** option on **Attack Complexity** for example, NVD have 98% estimated accuracy and 66% for **High** versus MITRE (see Figure 5), **Low** scored 87% and 75% for **High**. Instead of further analysing in this way, I will point out some of my worries around these results, not that they are wrong per se, but that they do not really tell us anything extra than that shown by in Figure 2, except perhaps it is helpful to see the results in a percentage estimated accuracy instead of a proportion. Unfortunately this is outside of my strengths, I did some cursory exploration into if doing this sort of analysis between two populations like this does make sense, the discourse here [15] suggests that such a thing can be done, though it also suggests the need for more informative

priors. This may apply in my case, and I intend on getting an expert opinion closer to home, however that will be after this report.



16
Figure 4: Confusion Matrices for NVD for CVSS version 3.1



17
Figure 5: Confusion Matrices for MITRE for CVSS version 3.1

Table 1: Confusion Matrices for NVD on CVSS version 2.0 from Johnson et al. [23]

Access vector				Access complexity			
	N	A	L		L	M	H
N	0.99	0	0	L	0.54	0.29	0.17
A	0.21	0.71	0.08	M	0.02	0.88	0.1
L	0.05	0.0	0.95	H	0.01	0.08	0.91
Authentication				Confidentiality			
	N	S	M		C	P	N
N	0.99	0.01	0	C	0.4	0.2	0.2
S	0.04	0.95	0.01	P	0.2	0.4	0.19
M	0.19	0.2	0.6	N	0.2	0.21	0.4
Integrity				Availability			
	C	P	N		C	P	N
C	0.91	0.08	0.01	C	0.92	0.07	0.01
P	0.05	0.93	0.02	P	0.07	0.91	0.02
N	0.02	0.07	0.92	N	0.02	0.11	0.87

4 Can we improve the LLMs classifier performance when we augment the NVD dataset with data from the MITRE dataset?

4.1 Related Work

Recent advancements in natural language processing and machine learning have led to significant progress in automating CVSS metric prediction from vulnerability descriptions. This area of research is particularly relevant to our work in Section 4.

Shahid and Debar [41] proposed CVSS-BERT, an approach leveraging state-of-the-art natural language processing techniques to determine the severity of a computer security vulnerability from its description. Their work, titled "CVSS-BERT: Explainable Natural Language Processing to Determine the Severity of a Computer Security Vulnerability from its Description," uses multiple BERT classifiers, one for each metric composing the CVSS vector. This approach is similar to our methodology in Section 4.

The authors of CVSS-BERT trained their model on CVE vulnerability data from

2018 to 2020, using BERT-small, a lighter version of BERT. They achieved high accuracy for predicting individual CVSS metrics, ranging from 83.79% to 96.07%. For 55.3% of the vulnerabilities in their test set, the predicted severity scores exactly matched the true severity scores. Their work also emphasized explainability, using gradient-based input saliency methods to determine the most relevant input words for each prediction.

Our work builds upon the foundations laid by CVSS-BERT, exploring the impact of augmenting the training data with additional sources (MITRE dataset) to potentially improve the model’s performance. We aim to investigate whether this augmentation can lead to more accurate and robust CVSS predictions.

Costa et al. [10] also contributed to this field with their paper "Predicting CVSS Metric via Description Interpretation." They tested various large language encoder-only models for generating CVSS from CVE descriptions, achieving state-of-the-art results with the DistilBERT model. Their work included text preprocessing techniques and explored the interpretability of the models using Shapley values, providing additional perspectives on enhancing CVSS prediction accuracy and explainability.

These studies collectively demonstrate the potential of transformer-based models in vulnerability assessment and highlight the challenges in accurately predicting CVSS metrics from textual descriptions. Our work in Section 4 builds upon these foundations, particularly the approach of CVSS-BERT, while exploring the impact of augmenting the training data with additional sources to potentially improve prediction accuracy and robustness.

Cody Airey –a classmate of mine– has been working on a similar problem. He has been reproducing results from Costa et al. [10], I am choosing to use Airey’s results due to them being on CVSS version 3.1 instead of a mix of versions 3.0 and 3.1 as Costa et al. used. My choice of model for CVSS prediction will very much bootstrap off his work. So far, a strong contender for state-of-the-art model for predicting CVSS metrics from CVE descriptions is the DistilBERT model [40]. This is a variant of BERT [12], a pre-trained deep bidirectional transformer large language model developed by Google. DistilBERT has advantages over other BERT models in terms of performance, but also on speed of training as well as size / memory footprint of the model.

4.1.1 Training

The model is trained separately for each of the eight CVSS metrics. To ensure robustness and reliability of our results, we trained each model on five different data splits. This approach allows us to calculate standard deviations, effectively

mitigating the impact of any potential "lucky" data split on our findings.

Our work diverges from Costa et al. and Airey’s approaches by utilizing a combined dataset from both NVD and MITRE, rather than relying solely on NVD data. We converted this merged data into a CSV format, containing vulnerability descriptions and their corresponding CVSS scores. This integration resulted in approximately 40,000 duplicate CVEs and a total of around 140,000 CVEs enriched with CVSS scores. Table 2 and Table 3 present the performance results of our DistilBERT model trained on this combined NVD and MITRE dataset. Contrary to our initial hypothesis, this data augmentation approach led to a decrease in performance across all metrics. However, we observed lower standard deviations for some metrics, indicating increased consistency in our model’s predictions. We noted anomalies in the balanced accuracy scores for certain metrics, which we attribute to the model not predicting all possible categorical options for these metrics. This phenomenon affects all balanced accuracy scores except for **Privileges Required** (PR) and **Confidentiality** (C), as evidenced in Table 2 and Table 3. The performance degradation can be attributed to the introduction of conflicting scores for overlapping CVEs from the two databases. While we anticipated that additional data would enhance the model’s generalization capabilities, our results indicate otherwise. This discrepancy between databases serves as a valuable indicator when analyzing generated CVE scores, suggesting that our model faces challenges similar to those encountered by human evaluators.

Our model’s lowest performance is observed in the **Availability Impact** category, a trend also noted in Airey’s work. This correlation between machine learning model difficulties and human scoring challenges provides insight into the inherent complexities of this particular metric. **Attack Complexity** presents another notable challenge, primarily due to data imbalance. The dataset is heavily skewed towards the **LOW** score for this metric, which impacts the model’s ability to accurately predict across all categories. These findings underscore the complexities involved in CVSS prediction and highlight areas for future improvement in both data curation and model development.

Metric	Model	AV	AC	PR	UI
Accuracy	NVD	91.28 \pm 0.26	95.64 \pm 0.68	82.77 \pm 0.24	93.86 \pm 0.19
	NVD and MITRE	72.81 \pm 0.32	92.62 \pm 0.15	81.18 \pm 0.18	66.35 \pm 0.24
F1	NVD	90.98 \pm 0.31	93.85 \pm 1.39	82.53 \pm 0.26	93.82 \pm 0.19
	NVD and MITRE	61.36 \pm 0.42	89.08 \pm 0.22	80.96 \pm 0.19	52.93 \pm 0.31
Bal Acc	NVD	67.88 \pm 2.11	55.82 \pm 7.23	75.98 \pm 0.47	92.46 \pm 0.21
	NVD and MITRE	25.00 \pm 0.00	50.00 \pm 0.00	75.18 \pm 0.31	50.00 \pm 0.00

Table 2: Comparison of the effects of the pre-trained models on the CVSS v3.1 dataset (Part 1).

Note: AV = Attack Vector, AC = Attack Complexity, PR = Privileges Required, UI = User Interaction

Bal Acc = Balanced Accuracy

Metric	Model	S	C	I	A
Accuracy	NVD	96.38 \pm 0.09	86.24 \pm 0.20	87.15 \pm 0.10	88.70 \pm 0.10
	NVD and MITRE	80.21 \pm 0.16	82.45 \pm 0.11	45.71 \pm 0.26	52.53 \pm 0.23
F1	NVD	96.30 \pm 0.10	86.09 \pm 0.21	87.11 \pm 0.10	88.04 \pm 0.11
	NVD and MITRE	71.40 \pm 0.22	82.34 \pm 0.12	28.68 \pm 0.28	36.18 \pm 0.27
Bal Acc	NVD	91.57 \pm 0.43	82.70 \pm 0.36	85.81 \pm 0.10	64.01 \pm 0.13
	NVD and MITRE	50.00 \pm 0.00	79.85 \pm 0.23	33.33 \pm 0.00	33.33 \pm 0.00

Table 3: Comparison of the effects of the pre-trained models on the CVSS v3.1 dataset (Part 2).

Note: S = Scope, C = Confidentiality Impact, I = Integrity Impact, A = Availability Impact

Bal Acc = Balanced Accuracy

Dont think this is the perfect question but leaving for now

5 In pursuit of keywords

This next section is the continuation of the exploratory study with a focus on what are the contributing factors for each metric / class for said metric. Clustering and analysis of the data makes sense in many ways. As the data is already grouped into Common Weakness Enumeration [28], it is likely that there are some sort of patterns

we can find. As a somewhat arbitrary place to start, choose K as 8 due to that being the number of CVSS metrics. For this I was hoping to see the data would naturally have some clusters based on the general theme of the vulnerability / description, e.g. a cluster based generally around SQL injection / databases in general.

5.1 K-Means Clustering of Vulnerability Descriptions

The process of clustering vulnerability descriptions using K-means consists of four main steps:

Data Preparation

The analysis begins by loading vulnerability descriptions from a dataset and using TF-IDF (Term Frequency-Inverse Document Frequency) vectorization to convert the text descriptions into numerical features [39]. This technique helps capture the importance of words within the context of the entire corpus often used in information retrieval.

Modelling

- Standard K-Means algorithm is utilized for clustering. This method aims to partition the descriptions into 8 clusters, each representing a group of similar vulnerabilities.
- The K-Means algorithm operates by iteratively assigning data points to the nearest cluster center and then updating the center based on the assigned points. The distance between data points is evaluated using the Euclidean distance between the TF-IDF vectors.

Evaluation

To assess the stability of the results, the clustering process is repeated multiple times with different random seeds. I did record a coherency score, but as this was just for exploratory purposes, the important part was if there was at least some indicator of this being a positive direction of inquiry.

Interpretation

Below are the example topics gained from the initial kmeans clustering:

- Cluster 0: vulnerability allows user service access attacker versions prior discovered issue
- Cluster 1: needed android id privileges lead possible execution interaction exploitation bounds
- Cluster 2: macos issue addressed improved fixed ios ipados able 15 13
- Cluster 3: code vulnerability attacker remote execution arbitrary execute file exploit user
- Cluster 4: site cross scripting xss plugin stored vulnerability wordpress forgery csrf
- Cluster 5:sql injection php parameter v1 vulnerability contain discovered admin vulnerable
- Cluster 6: manipulation identifier leads vdb vulnerability classified unknown remotely attack disclosed
- Cluster 7: cvss oracle mysql vulnerability attacks server access base unauthorized score

The most promising from this initial set is cluster four, highlighted in red above. Cross-site scripting (XSS) and wordpress plugins are a common trend within CVEs. Figure 6 shows the counts per year of CVEs published containing at least five of the words from that cluster. Five is arbitrary, this is more here just to show a potential insight in looking at these trends. In this case, from 2019 to 2023 there is a large increase in these types of vulnerabilities, mainly in WordPress plugins. You may notice the trend matches with the general trend of CVEs published (Figure 1), this is a positive result, in that, the clustering is still following the underlying distribution and has not failed to capture the overall trend.

Here are some example of the descriptions, for your viewing pleasure:

This was a good indicator that it is worth looking into clustering further.

CVE ID	Description
CVE-2023-24378	Auth. (contributor+) Stored Cross-Site Scripting (XSS) vulnerability in Codeat Glossary plugin \leq 2.1.27 versions.
CVE-2023-24396	Auth. (admin+) Stored Cross-Site Scripting (XSS) vulnerability in E4J s.R.L. VikBooking Hotel Booking Engine & PMS plugin \leq 1.5.11 versions.
CVE-2023-25062	Auth. (admin+) Stored Cross-Site Scripting (XSS) vulnerability in PINPOINT.WORLD Pinpoint Booking System plugin \leq 2.9.9.2.8 versions.

Table 4: CVE Descriptions for Various WordPress Plugins

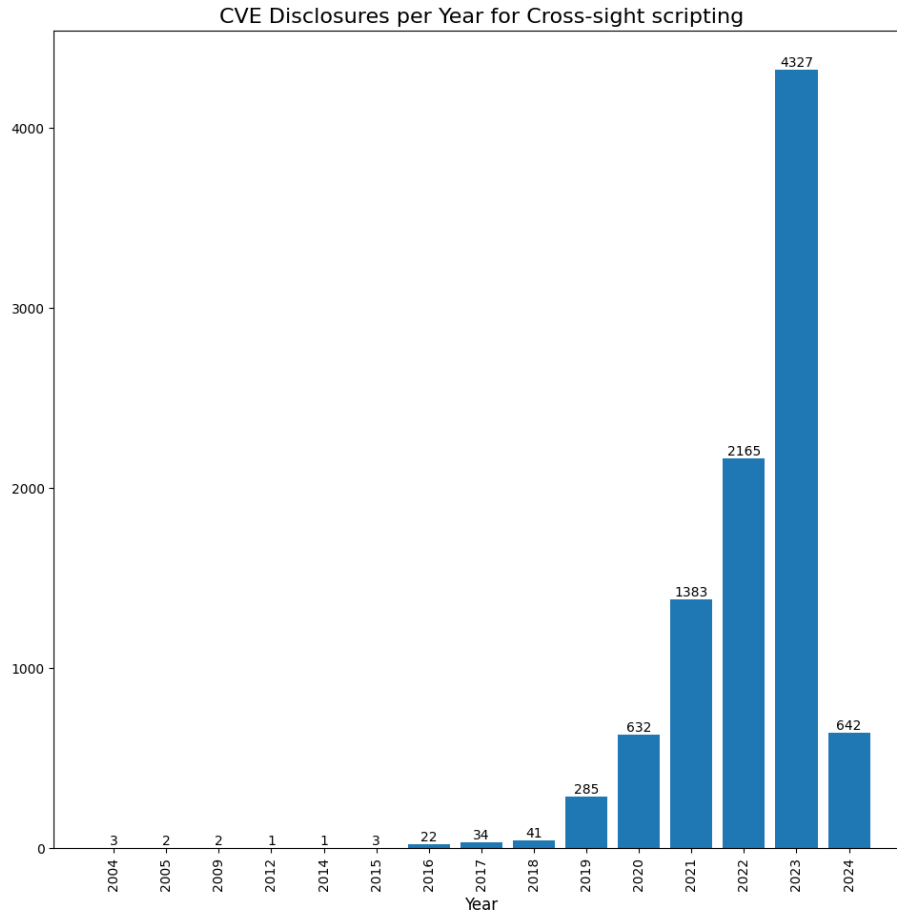


Figure 6: CVEs with descriptions containing at least five of the words from cluster four

5.2 Latent Dirichlet Allocation Methodology

The next approach to topic modeling utilizes Latent Dirichlet Allocation (LDA) [6] in conjunction with Word2Vec [26] embeddings to analyze Common Vulnerabilities and Exposures (CVE) descriptions. As an unsupervised learning method, LDA is well-suited for exploring large collections of text data without predefined categories [6, 44]. This is particularly valuable in the CVE context, where we aim to discover latent themes in vulnerability descriptions without prior knowledge of these themes. LDA’s soft clustering approach allows CVE descriptions to belong to multiple topics with varying probabilities [44], which is beneficial given the multi-faceted nature of many vulnerabilities. Furthermore, LDA produces human-interpretable topics [6], facilitating easier derivation of insights from the clustering results. The paper “Security trend analysis with CVE topic models, by Neuhaus & Zimmerman(2010)”[31] uses this technique to great effect on the NVD dataset (the analysis was for data ending in 2009) giving precedence to LDA being effective for this sort of task.

The methodology comprises several key steps:

Data Preprocessing

I begin by preprocessing the CVE descriptions:

- Custom stopwords are defined, combining standard English stopwords with domain-specific terms (e.g., “vulnerability”, “attacker”).
- Each description is tokenized and filtered, removing stopwords and short tokens (length ≤ 2).
- A dictionary is created from the processed texts, filtering out extremely rare and extremely common terms.

This data processing is necessary due to this type of topic discovery relying on word distributions, and therefore word frequencies. As a result stop words just muddy the water and put more interesting and relevant topic words further down.

Word Embedding

A Word2Vec [26] model is trained on the preprocessed texts:

- Vector size: 100
- Window size: 5

- Minimum word count: 2

This embedding captures semantic relationships between words in the CVE context.

The parameters chosen here just felt like reasonable defaults and have not been explored. There are also other options in using something like pretrained fasttext [7] as used in [43].

LDA Model Training

I have done many different iterations of the models with varying numbers of topics and hyperparameters. In grid search attempts I found that the different coherences metrics did not agree with each other. In addition, as we are searching for unknown topics in a unsupervised fashion, it is dubious how helpful these measures of fit are in any case. My results of such exploration will be found in the appendices. The hyperparameters below are broadly matching what I found to be the best, however the main metric came from looking at the actual clustering results, and the biggest impact came from the number of topics, with the final models being picked by their purity score (as described in Section 5.3) The key hyperparameters used in the LDA model training are as follows:

- Number of topics: 20, 40, 60, 80, 100
 - This parameter determines the number of topics the model attempts to discover in the corpus.
 - The broad range allows for exploration of both coarse-grained and fine-grained topic structures, which is crucial for identifying potential sub-categories within CVSS metrics or broader themes across different metrics.
- α : "symmetric"
 - Controls the prior distribution over topic mixtures for each document.
 - The "symmetric" setting assumes all topics are equally likely a priori for each document, providing an unbiased starting point for topic discovery in CVSS vulnerability descriptions.
- η : 0.1
 - This is the prior for word distributions over topics.

- The relatively small value of 0.1 can lead to more specific and distinct topics, potentially helping in distinguishing between different CVSS categories more clearly.
- Number of passes: 30
 - Represents the number of times the model cycles through the entire corpus during training.
 - 30 passes allow the model to refine its topic assignments multiple times, providing sufficient iterations to capture the underlying topic structure of CVSS data without overfitting.
- Number of iterations: 200
 - Sets the maximum number of iterations for each document during the inference process.
 - 200 iterations allow for thorough topic inference for each document, which is particularly important for CVSS vulnerability descriptions that might contain complex or technical language.

This set of hyperparameters collectively defines a thorough exploration of the topic space, allowing for the discovery of both broad and specific themes in the CVSS data. The wide range of topic numbers (20 to 100) was particularly impactful in identifying the optimal granularity for representing the vulnerability descriptions. The LDA model is implemented using Gensim’s LdaMulticore[37], allowing for parallel processing. In general I found Gensim very nice to work with, the model training was fast and the API made sense. The combination of Gensim’s efficient implementation and these carefully chosen hyperparameters facilitated a comprehensive exploration of the topic structure in the CVSS vulnerability descriptions.

Topic Assignment and Analysis

For each saved model:

- Topic assignments are generated for each document in the corpus.
- Results, including document index, description, assigned topic, and CVSS data, are saved in JSON format.
- The top 50 words for each topic are extracted and saved for interpretation.

Cluster-Class Association

For each class within each CVSS metric, we identify the most representative cluster:

- Calculate the proportion of documents from each class assigned to each topic cluster.
- The cluster with the highest count for a given class is considered its best representative.
- Calculate the purity score for all the different numbers of topics

With the best cluster from a given set discovered, we now need to decide which number of clusters provides the best overall fit. A rudimentary approach is to just use the purity score.

5.3 Purity as a Cluster Evaluation Metric

Purity is a simple external evaluation measure for cluster quality, particularly useful when predefined classes are available for the data [24].

Definition and Calculation

For a set of clusters $\Omega = \{\omega_1, \omega_2, \dots, \omega_K\}$ and a set of classes $C = \{c_1, c_2, \dots, c_J\}$, purity is defined as:

$$\text{purity}(\Omega, C) = \frac{1}{N} \sum_{k=1}^K \max_j |\omega_k \cap c_j| \quad (2)$$

where N is the total number of objects, and $|\omega_k \cap c_j|$ is the number of objects in cluster ω_k that belong to class c_j .

Interpretation

Purity ranges from 0 to 1, where 1 indicates perfect purity. A higher purity value generally suggests better clustering quality with respect to the ground truth classes [25].

In the context of topic modeling for CVSS metrics, high purity would indicate that each discovered topic strongly corresponds to a specific CVSS category.

Limitations

It's important to note that purity tends to increase as the number of clusters increases, with a trivial solution of perfect purity when each object is in its own cluster [38]. Therefore, it should be interpreted cautiously, especially when comparing models with different numbers of topics, however, as seen by Figures 7, 8, 9 this effect did not present itself too strongly. This effect incentivise picking the number of topics around the elbow of the graph (the point on a curve where the rate of change significantly shifts), as any minor increase in the purity score could just be based on the inherent increase in the purity score as the number of topics increase.

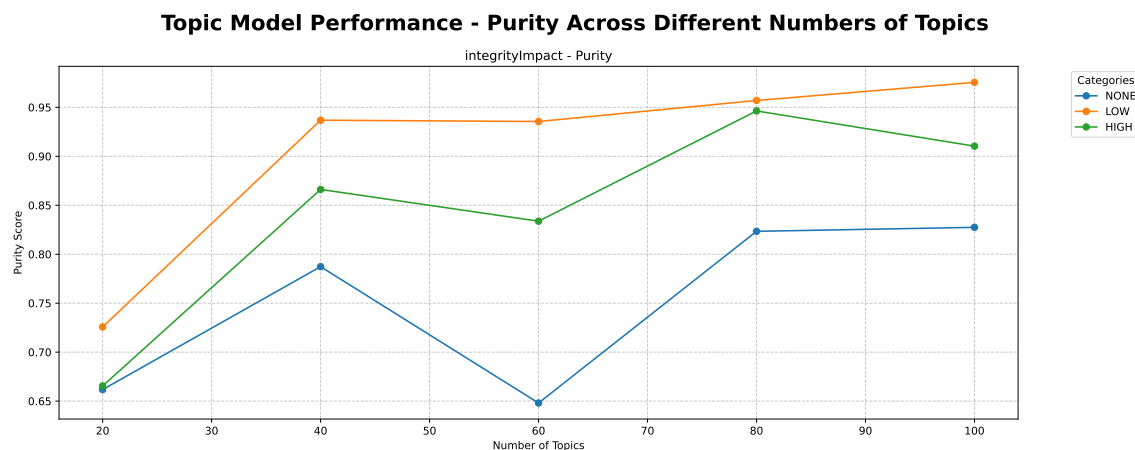


Figure 7: The purity score for when the topics have been assigned with a focus on Integrity Impact, split by the different possible classes

5.4 Evaluation and Insights

Elephant in the Room, Data Imbalance

As is obvious from Figure 3, the classes for each metric are highly imbalanced. This is just a reality of the data and so is something to contend with. This aspect made finding a well defined cluster for each topic difficult as well as making the graphing and interpreting the outcome of the clustering somewhat awkward.

As an example, if we take look at Figure 10, which shows the clustering for Privileges Required of the best cluster with a focus on the HIGH class, we can see that even though we are trying to find the cluster with the best representation of HIGH both of the other classes are more dominant. This is a result of the massive

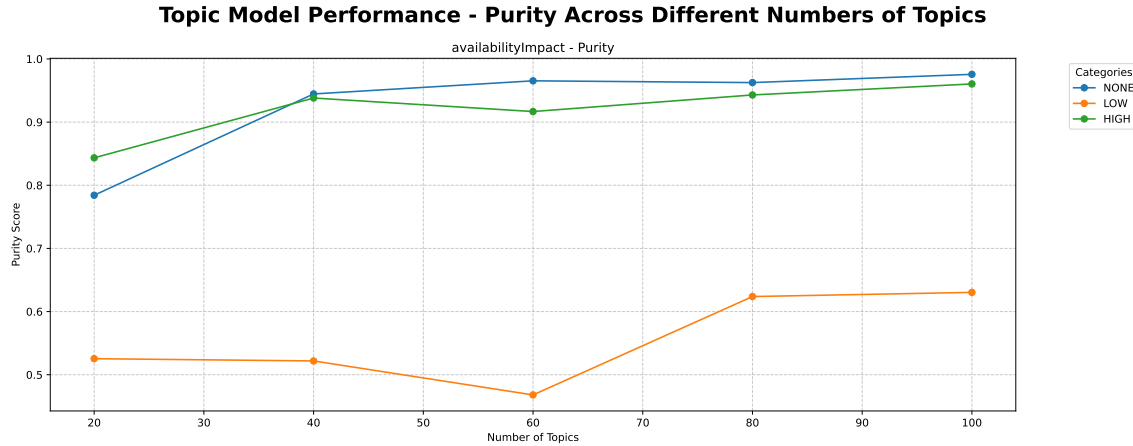


Figure 8: The purity score for when the topics have been assigned with a focus on Availability Impact, split by the different possible classes

class imbalance, the clustering cannot manage to find a class to represent such a little proportion of the data. As a result the analysis will be focused on the metrics which are the most naturally balanced as well as related to each other, these are the CIA triad, confidentiality, integrity, and availability impact.

5.5 Clustering Results

The clustering analysis using Latent Dirichlet Allocation (LDA) yielded insights into the patterns of vulnerabilities as they relate to the Common Vulnerability Scoring System (CVSS) metrics, particularly the integrity impact metric. While an ideal analysis would involve comparing our clusters with the assigned Common Weakness Enumeration (CWE), the time constraints and the complex nature of both our topic-based clusters and the CWE descriptions made automated cross-referencing unfeasible for the time being.

Methodology and Presentation

The primary results of the analysis are presented in graphical form. Each graph shows the distribution of documents across the three classes of integrity impact (NONE, LOW, HIGH) for the topic cluster that best represents each target class. This approach allows us to visualize how well this clustering method separates vulnerabilities based on their integrity impact.

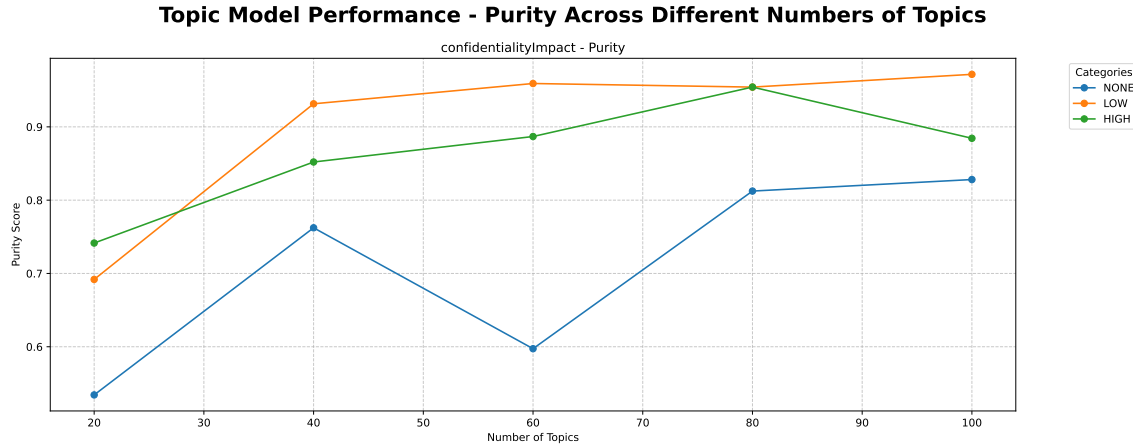


Figure 9: The purity score for when the topics have been assigned with a focus on **Confidentiality Impact**, split by the different possible classes

Fix the position of the myriad graphs that are in weird places in the document

5.5.1 Integrity Impact Analysis

The **Integrity Impact** metric in CVSS refers to the degree to which a vulnerability, if exploited, could affect the trustworthiness and veracity of data. A high integrity impact implies that data could be significantly corrupted or altered, potentially leading to serious consequences for the affected system or its users.

NONE Category

Figure 11 shows the distribution for the cluster best representing the **NONE** category of integrity impact.

In this cluster, we observe a predominance of vulnerabilities related to denial of service attacks and system crashes. While these issues are serious and can affect system availability, they typically do not directly compromise data integrity. This aligns with the expectations for vulnerabilities classified as having no integrity impact.

Key findings:

- The cluster shows a clear majority of **NONE**-rated vulnerabilities.

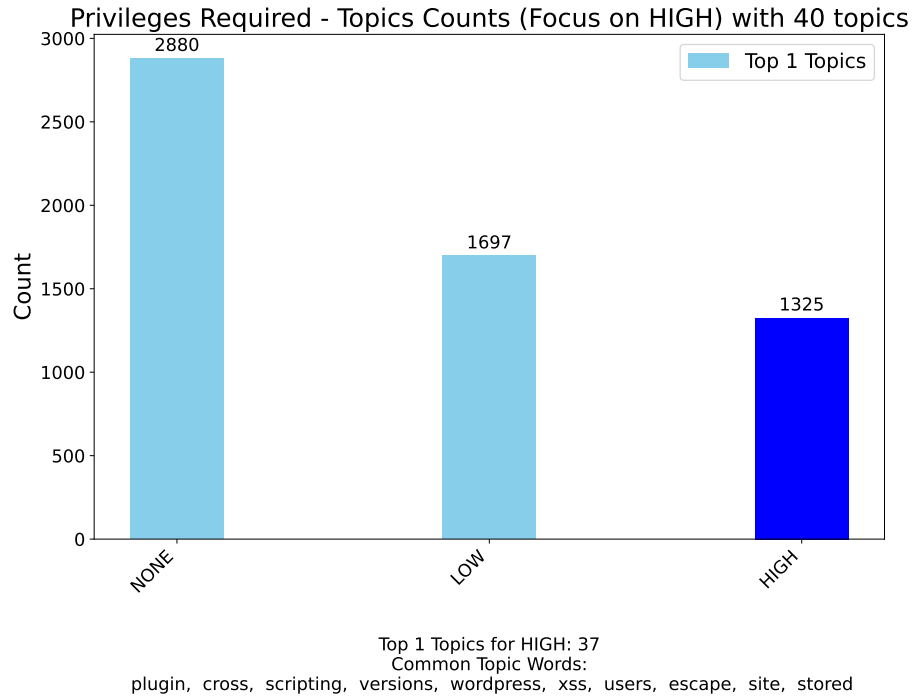


Figure 10: The counts of the documents within the best topic in relation to Privileges Required with target class HIGH

- Common terms in this cluster likely include "denial of service", "crash", and "dos".
- This result supports the effectiveness of our clustering in identifying vulnerabilities with no integrity impact.

LOW Category

Figure 12 represents the distribution for the cluster best representing the LOW category of integrity impact.

This cluster predominantly features cross-site scripting (XSS) vulnerabilities and other generally web browser based attacks. These types of vulnerabilities can indeed cause data integrity issues, but they are usually limited in scope, typically affecting individual user interactions rather than compromising the entire application's data integrity.

Key findings:

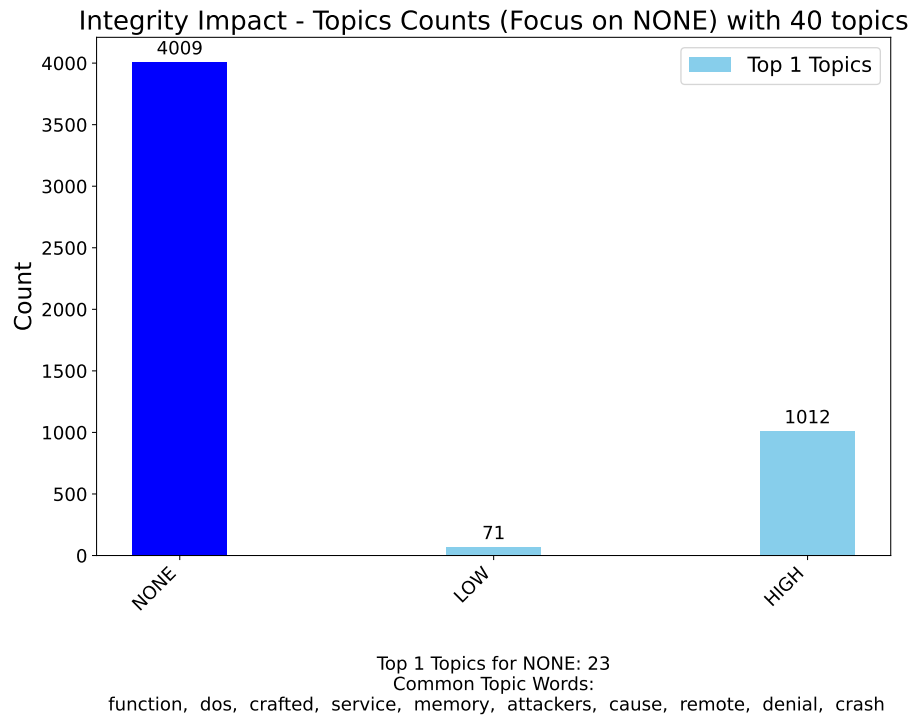


Figure 11: The counts of the documents within the best topic in relation to Integrity Impact with target class NONE

- The cluster shows a majority of LOW-rated vulnerabilities, with some overlap into other categories.
- Common terms likely include "cross-site scripting", "XSS", "javascript", and "html".
- This cluster aligns with our earlier k-means clustering results, providing some similarities between the clustering

HIGH Category

Needs a bit of work, when switching to purity this changed a little what was the best topic...

Figure 13 shows the distribution for the cluster best representing the HIGH category of integrity impact.

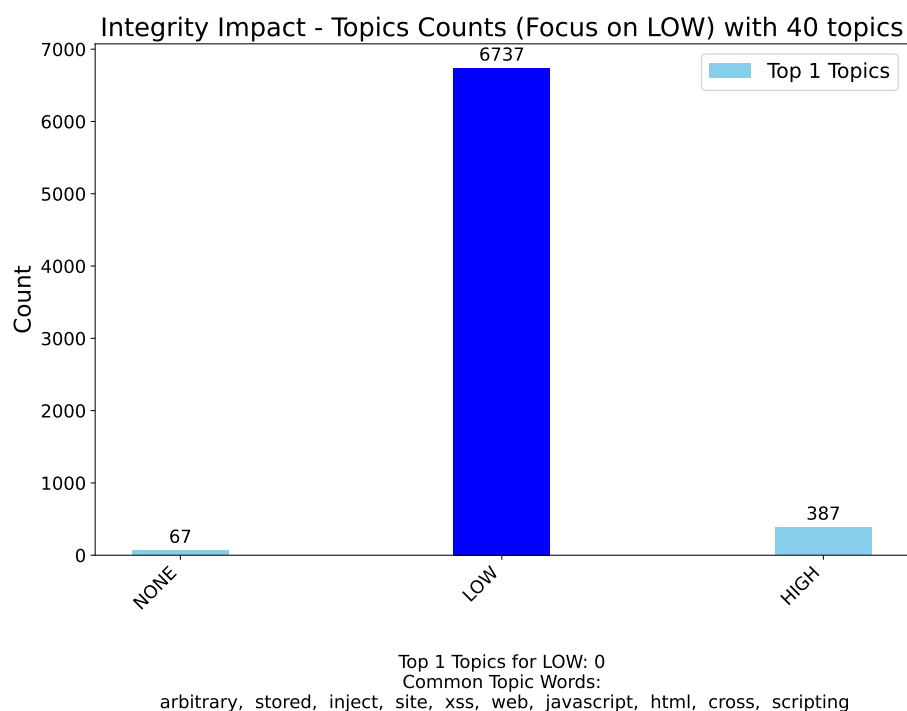


Figure 12: The counts of the documents within the best topic in relation to Integrity Impact with target class LOW

In this cluster, we see a focus on remote code execution vulnerabilities. These types of attacks have the potential to severely compromise data integrity by allowing unauthorized modification or corruption of server contents.

Key findings:

- The cluster shows a strong representation of HIGH-rated vulnerabilities.
- Common terms likely include "remote", "execution"
- The prevalence of these severe vulnerabilities in this cluster aligns with cybersecurity expert expectations for high-impact integrity issues.

I was also planning on adding confidentiality and availability impact after here which I have not done yet, still a good idea? So that I can talk about them together a bit more

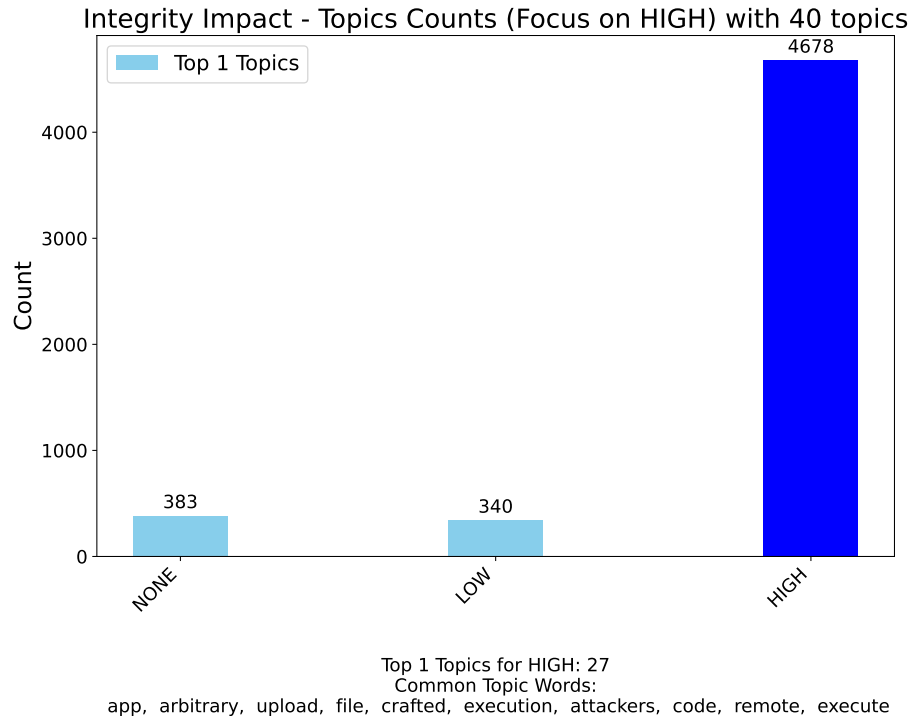


Figure 13: The counts of the documents within the best topic in relation to integrityImpact with target class HIGH

5.6 Cluster Merging and Refinement

In an attempt to enhance class representation and potentially improve the overall performance of the topic model, I explored the concept of merging related clusters. This process involved several steps:

- Identifying clusters with semantic similarity or overlapping representation of CVSS classes.
- Merging these clusters by combining their topic-word distributions and reassigning documents.
- Recalculating class representation metrics for the merged clusters.
- Comparing the performance of the merged model to the original in terms of:
 - Overall topic coherence

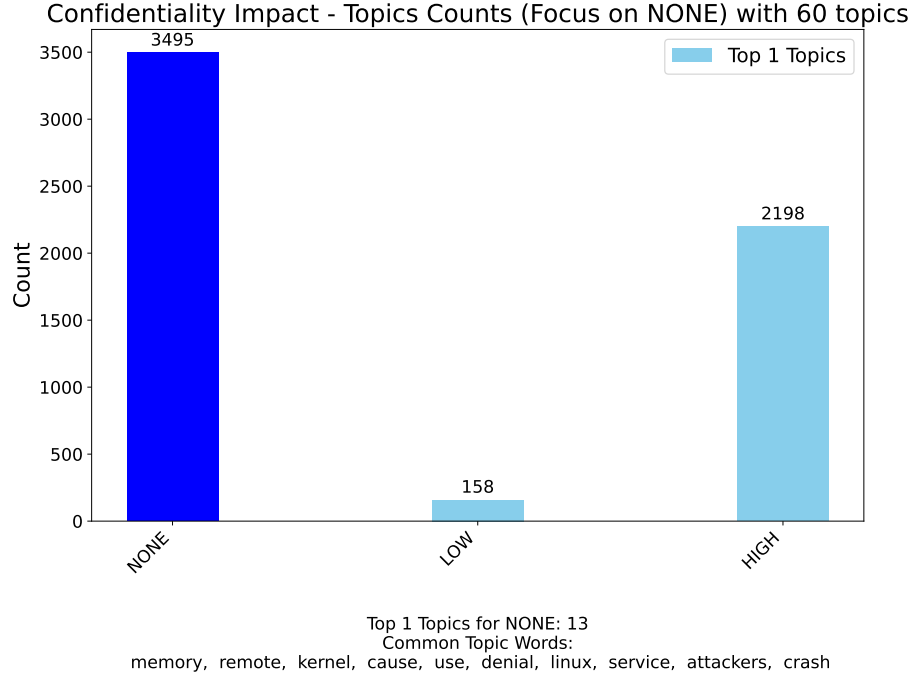


Figure 14: The counts of the documents within the best topic in relation to integrity-Impact with target class NONE, for the number of clusters (60) with the lowest purity

- Clarity of class representation
- Interpretability of resulting topics

5.6.1 Methodology

The approach to cluster merging was based on a simple similarity metric between topic-word distributions. Clusters were merged if their similarity exceeded a predefined threshold. While this method is straightforward to implement, it is admittedly crude and may not capture all the nuances of semantic similarity between topics.

Other approaches exist, for instance, Neuhaus & Zimmerman (2010) [31] merged clusters based on the similarity of the topic words found by their model. This other approach potentially offers a more nuanced way of identifying truly related topics.

5.6.2 Results

To illustrate the effects of cluster merging, we present the following figures:

need to add these in!

Contrary to our initial expectations, the process of merging clusters did not lead to improved results. Instead, I observed the following effects:

The merging of clusters resulted in several negative impacts on the model’s quality and utility. The process led to decreased cluster purity, with a higher mix of different CVSS classes within each cluster. This reduction in homogeneity made the larger, merged clusters more difficult to interpret, diluting their defining characteristics and semantic meanings. Consequently, the model lost granularity, as specific types of vulnerabilities or attack vectors were grouped into broader, less nuanced categories. Overall, these changes contributed to a decrease in topic coherence, indicating that the merged clusters were less internally consistent than their more granular predecessors.

5.7 Interpretation and Future Directions

These results suggest that our initial clustering approach was already capturing meaningful distinctions between different types of vulnerabilities and their associated CVSS ratings. The process of merging clusters, at least with the current methodology, appears to obscure these distinctions rather than enhance them.

However, this does not necessarily mean that cluster merging is inherently unhelpful. Instead, it suggests that more sophisticated merging strategies may be needed to preserve the valuable information captured in the original clusters while still allowing for the combination of truly related topics.

Possible future work on this topic includes:

- Implementing the cluster merging approach of Neuhaus & Zimmerman (2010), which focuses on the similarity of topic words rather than overall distribution similarity.
- Exploring hierarchical topic modeling approaches, which might allow for the representation of both fine-grained and broader topic categories simultaneously.
- Investigating dynamic topic modeling techniques to capture how vulnerability types and their CVSS ratings evolve over time, potentially revealing patterns in how initially distinct vulnerability categories merge or diverge over time.

5.8 Overall Observations

1. Our clustering approach successfully differentiated between vulnerabilities with varying levels of integrity impact, as evidenced by the distinct profiles of each

cluster.

2. The results align well with expert knowledge in the field of cybersecurity, suggesting that our unsupervised learning approach can capture meaningful patterns in vulnerability data.
3. While our current analysis provides a static view of the vulnerability landscape, the temporal aspect of the data presents an opportunity for future research. Analyzing how these clusters and trends have evolved over time could provide additional insights into the changing nature of cybersecurity threats.
4. The overlap between categories in some clusters (particularly visible in the LOW category) highlights the complexity of vulnerability classification and the potential for ambiguity in CVSS scoring.

5.9 Future Work

As I am a Masters student, I will be continuing the work that I have done here for COSC480 into next year. My future work will focus on several key areas:

- **Dataset Refinement:** The CVE dataset contains numerous poorly written descriptions that hinder machine learning model performance. We plan to develop robust methods for identifying and filtering out low-quality entries, potentially using the insights gained from our clustering analysis.
- **Advanced Clustering Techniques:** Building on our current work, we aim to explore more sophisticated clustering approaches. This includes implementing the cluster merging method proposed by Neuhaus & Zimmerman (2010) [31], which focuses on the similarity of topic words.
- **Temporal Analysis:** We intend to investigate how vulnerability patterns and their CVSS ratings evolve over time. This could involve applying dynamic topic modeling techniques to capture temporal trends in the data.

Through these efforts, I aim to develop a more robust, data-driven approach to vulnerability analysis and CVSS scoring that can adapt to the evolving landscape of cybersecurity threats.

6 Conclusion

The Common Vulnerability Scoring System (CVSS) plays a vital role in prioritizing and managing the ever-increasing number of software vulnerabilities. This research has highlighted both the strengths and limitations of current CVSS implementation and scoring practices.

Key findings from our study include:

- **Database Variability:** We observed variability in CVSS scores between different databases, such as NVD and MITRE. This inconsistency underscores the subjective nature of vulnerability scoring and the challenges in establishing a reliable ground truth.
- **Clustering Insights:** The application of Latent Dirichlet Allocation (LDA) for clustering CVE descriptions revealed meaningful patterns in vulnerability types and their associated CVSS ratings. This demonstrates the potential of unsupervised learning techniques in capturing latent structures within vulnerability data.
- **Cluster Merging Limitations:** Our experiments with cluster merging, while not yielding immediate improvements, provided valuable insights into the robustness of our initial clustering and the complexity of semantic relationships between vulnerability types.
- **Data Quality Importance:** The analysis highlighted the critical role of data quality in both manual and automated CVSS scoring processes. Poorly written or inconsistent CVE descriptions pose significant challenges for accurate vulnerability assessment.

While CVSS remains a crucial tool for vulnerability management, our research suggests that its current implementation has limitations that need to be addressed. The integration of machine learning models, particularly those focused on natural language processing and topic modeling, offers promising avenues for automating and enhancing the accuracy of CVSS scoring.

By combining advanced machine learning techniques with domain expertise in cybersecurity, we can work towards a more consistent, accurate, and interpretable system for assessing and prioritizing software vulnerabilities. This improved system would not only enhance the efficiency of vulnerability management but also contribute to a more secure digital ecosystem overall.

References

- [1] M Ugur Aksu et al. “Automated generation of attack graphs using NVD”. In: *Proceedings of the Eighth ACM Conference on Data and Application Security and Privacy*. 2018, pp. 135–142.
- [2] Ellen Richey et al. *Payment card industry (PCI) data security standard*. https://www.pcisecuritystandards.org/document_library?category=pcidss&document=dss4aag. [Online; accessed July-2024]. 2018.
- [3] Luca Allodi and Fabio Massacci. “Comparing Vulnerability Severity and Exploits Using Case-Control Studies”. In: *ACM Transactions on Information and System Security* 17 (Aug. 2014), pp. 1–20. DOI: [10.1145/2630069](https://doi.org/10.1145/2630069).
- [4] Anonymous. *CVSS v3 and v3.1 missing temporal metrics exploit code maturity and remediation*. <https://security.stackexchange.com/questions/270257/cvss-v3-and-v3-1-missing-temporal-metrics-exploit-code-maturity-and-remediation>. [Online; accessed July-2024]. 2024.
- [5] Hodaya Binyamini et al. “A framework for modeling cyber attack techniques from security vulnerability descriptions”. In: *Proceedings of the 27th ACM SIGKDD conference on knowledge discovery & data mining*. 2021, pp. 2574–2583.
- [6] David M. Blei, Andrew Y. Ng, and Michael I. Jordan. “Latent dirichlet allocation”. In: *J. Mach. Learn. Res.* 3.null (Mar. 2003), pp. 993–1022. ISSN: 1532-4435.
- [7] Piotr Bojanowski et al. *Enriching Word Vectors with Subword Information*. 2017. arXiv: [1607.04606](https://arxiv.org/abs/1607.04606) [cs.CL]. URL: <https://arxiv.org/abs/1607.04606>.
- [8] CISA. *Known Exploited Vulnerabilities Catalog*. <https://www.cisa.gov/known-exploited-vulnerabilities-catalog>. [Online; accessed June-2024]. 2024.
- [9] The MITRE Corporation. *List of Partners*. <https://www.cve.org/PartnerInformation/ListofPartners>. [Online; accessed June-2024]. 2024.
- [10] Joana Cabral Costa et al. “Predicting CVSS Metric via Description Interpretation”. In: *IEEE Access* 10 (2022), pp. 59125–59134. DOI: [10.1109/ACCESS.2022.3179692](https://doi.org/10.1109/ACCESS.2022.3179692).

- [11] Kenneth Brown David Beën. *Department of Defense (DOD) Joint Special Access Program (SAP) Implementation Guide (JSIG)*. [https://www.dcsa.mil/portals/91/documents/ctp/nao/JSIG_2016April11_Final_\(53Rev4\).pdf](https://www.dcsa.mil/portals/91/documents/ctp/nao/JSIG_2016April11_Final_(53Rev4).pdf). [Online; accessed July-2024]. 2016.
- [12] Jacob Devlin et al. *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. 2019. arXiv: [1810.04805 \[cs.CL\]](https://arxiv.org/abs/1810.04805). URL: <https://arxiv.org/abs/1810.04805>.
- [13] Ying Dong et al. “Towards the Detection of Inconsistencies in Public Security Vulnerability Reports”. In: *28th USENIX Security Symposium (USENIX Security 19)*. Santa Clara, CA: USENIX Association, Aug. 2019, pp. 869–885. ISBN: 978-1-939133-06-9. URL: <https://www.usenix.org/conference/usenixsecurity19/presentation/dong>.
- [14] FIRST. *CVSS Frequently Asked Questions*. <https://www.first.org/cvss/v2/faq#Explanation-of-CVSS-v2-formula-and-metric-valued-development>. [Online; accessed July-2024]. 2007.
- [15] Andrew Gelman. *Hierarchical modeling when you have only 2 groups: I still think it’s a good idea, you just need an informative prior on the group-level variation*. <https://statmodeling.stat.columbia.edu/2015/12/08/hierarchical-modeling-when-you-have-only-2-groups-i-still-think-its-a-good-idea-you-just-need-an-informative-prior-on-the-group-level-variation/>. [Online; accessed July-2024]. 2024.
- [16] Rameshwar D. Gupta and Donald St. P. Richards. “The History of the Dirichlet and Liouville Distributions”. In: *International Statistical Review / Revue Internationale de Statistique* 69.3 (2001), pp. 433–446. ISSN: 03067734, 17515823. URL: <http://www.jstor.org/stable/1403455> (visited on 10/09/2024).
- [17] Henry Howland. “CVSS: Ubiquitous and Broken”. In: *Digital Threats* 4.1 (Feb. 2022). DOI: [10.1145/3491263](https://doi.org/10.1145/3491263). URL: <https://doi.org/10.1145/3491263>.
- [18] Forum of Incident Response and Security Teams (FIRST). *A Complete Guide to the Common Vulnerability Scoring System*. <https://www.first.org/cvss/v2/guide>. [Online; accessed June-2024]. 2024.
- [19] Forum of Incident Response and Security Teams (FIRST). *The EPSS Model*. <https://www.first.org/epss/model>. [Online; accessed June-2024]. 2024.
- [20] Forum of Incident Response and Security Teams (FIRST). *The EPSS User Guide*. <https://www.first.org/epss/user-guide>. [Online; accessed June-2024]. 2024.

- [21] Mark Jerrum and Alistair Sinclair. “The Markov chain Monte Carlo method: an approach to approximate counting and integration”. In: *Approximation Algorithms for NP-Hard Problems*. USA: PWS Publishing Co., 1996, pp. 482–520. ISBN: 0534949681.
- [22] Yuning Jiang and Yacine Atif. “An Approach to Discover and Assess Vulnerability Severity Automatically in Cyber-Physical Systems”. In: *13th International Conference on Security of Information and Networks*. SIN 2020: 13th International Conference on Security of Information and Networks. Merkez Turkey: ACM, Nov. 4, 2020, pp. 1–8. ISBN: 978-1-4503-8751-4. DOI: [10.1145/3433174.3433612](https://doi.org/10.1145/3433174.3433612). URL: <https://dl.acm.org/doi/10.1145/3433174.3433612> (visited on 02/28/2024).
- [23] Pontus Johnson et al. “Can the Common Vulnerability Scoring System be Trusted? A Bayesian Analysis”. In: *IEEE Transactions on Dependable and Secure Computing* 15.6 (2018), pp. 1002–1015. DOI: [10.1109/TDSC.2016.2644614](https://doi.org/10.1109/TDSC.2016.2644614).
- [24] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. *Introduction to Information Retrieval*. Cambridge University Press, 2008.
- [25] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. *Introduction to Information Retrieval*. Cambridge University Press, 2008.
- [26] Tomas Mikolov et al. *Efficient Estimation of Word Representations in Vector Space*. 2013. arXiv: [1301.3781](https://arxiv.org/abs/1301.3781) [cs.CL]. URL: <https://arxiv.org/abs/1301.3781>.
- [27] MITRE. *Common Vulnerabilities and Exposures — CVE® The Standard for Information Security Vulnerability Names*. <https://cve.mitre.org/docs/cve-intro-handout.pdf>. [Online; accessed July-2024]. 2024.
- [28] MITRE. *Common Vulnerabilities and Exposures — CVE® The Standard for Information Security Vulnerability Names*. https://cwe.mitre.org/about/new_to_cwe.html. [Online; accessed September-2024]. 2024.
- [29] MITRE. *CVE Numbering Authorities (CNAs)*. <https://www.cve.org/ProgramOrganization/CNAs>. [Online; accessed May-2024]. 2024.
- [30] MITRE. *MITRE landing page*. <https://cve.mitre.org/>. [Online; accessed February-2024]. 2024.
- [31] Stephan Neuhaus and Thomas Zimmermann. “Security Trend Analysis with CVE Topic Models”. In: *2010 IEEE 21st International Symposium on Software Reliability Engineering*. 2010, pp. 111–120. DOI: [10.1109/ISSRE.2010.53](https://doi.org/10.1109/ISSRE.2010.53).

- [32] NVD. *CVE-2024-38526 Detail*. <https://nvd.nist.gov/vuln/detail/CVE-2024-38526>. [Online; accessed June-2024]. 2024.
- [33] NVD. *NVD landing page*. <https://nvd.nist.gov/>. [Online; accessed February-2024]. 2024.
- [34] Sasha Romanosky Peter Mell Karen Scarfone. *Common Vulnerability Scoring System v3.1: Specification Document*. <https://www.first.org/cvss/v3.1/specification-document>. [Online; accessed February-2024]. 2024.
- [35] PyMC. *PyMC landing page*. <https://www.pymc.io/welcome.html>. [Online; accessed May-2024]. 2024.
- [36] Qualys. *CVSS Vector Strings*. https://qualysguard.qualys.com/qwebhelp/fo_portal/setup/cvss_vector_strings.htm. [Online; accessed July-2024]. 2024.
- [37] Radim Rehurek. *Gensim Landing Page*. <https://radimrehurek.com/gensim/>. [Online; accessed September-2024]. 2024.
- [38] Andrew Rosenberg and Julia Hirschberg. “V-Measure: A Conditional Entropy-Based External Cluster Evaluation Measure.” In: Jan. 2007, pp. 410–420.
- [39] G. Salton and M.J. McGill. *Introduction to Modern Information Retrieval*. International student edition. TF-IDF concept discussed in Chapter 1, pages 63–71. McGraw-Hill, 1983. Chap. 1, pp. 63–71. ISBN: 9780070544840. URL: <https://books.google.co.nz/books?id=7f5TAAAMAAJ>.
- [40] Victor Sanh et al. *DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter*. 2020. arXiv: [1910.01108](https://arxiv.org/abs/1910.01108) [cs.CL]. URL: <https://arxiv.org/abs/1910.01108>.
- [41] Mustafizur R. Shahid and Hervé Debar. “CVSS-BERT: Explainable Natural Language Processing to Determine the Severity of a Computer Security Vulnerability from its Description”. In: *CoRR* abs/2111.08510 (2021). arXiv: [2111.08510](https://arxiv.org/abs/2111.08510). URL: <https://arxiv.org/abs/2111.08510>.
- [42] Jonathan Spring et al. “Towards improving CVSS”. In: *SEI, CMU, Tech. Rep* (2018).
- [43] Mark-Oliver Stehr and Minyoung Kim. *Vulnerability Clustering and other Machine Learning Applications of Semantic Vulnerability Embeddings*. 2023. arXiv: [2310.05935](https://arxiv.org/abs/2310.05935) [cs.CR]. URL: <https://arxiv.org/abs/2310.05935>.
- [44] Mark Steyvers and T. Griffiths. “Probabilistic topic models. Handb. Latent Semant”. In: *Anal* 427 (Jan. 2007), pp. 424–440.

Appendix .1 CVSS figures from other versions

Below is a collection of confusion matrices which are results from the Bayesian analysis of CVSS versions 2.0 & 3.0 for both the NVD and MITRE databases. Version 2.0 for MITRE is especially rough as there was very little data.

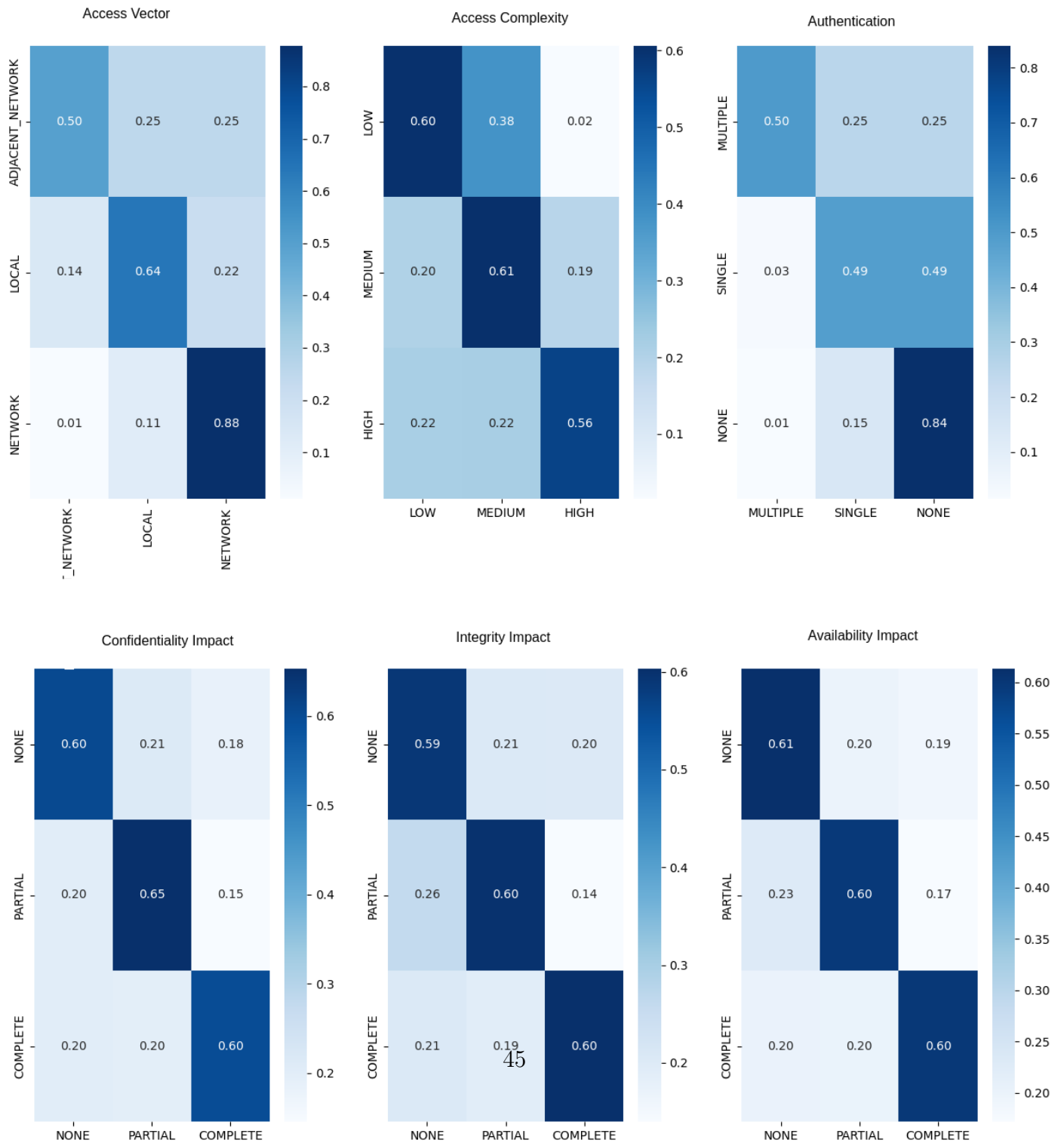


Figure 15: Confusion matrix of estimated accuracy for CVSS metrics for version 2.0 for NVD

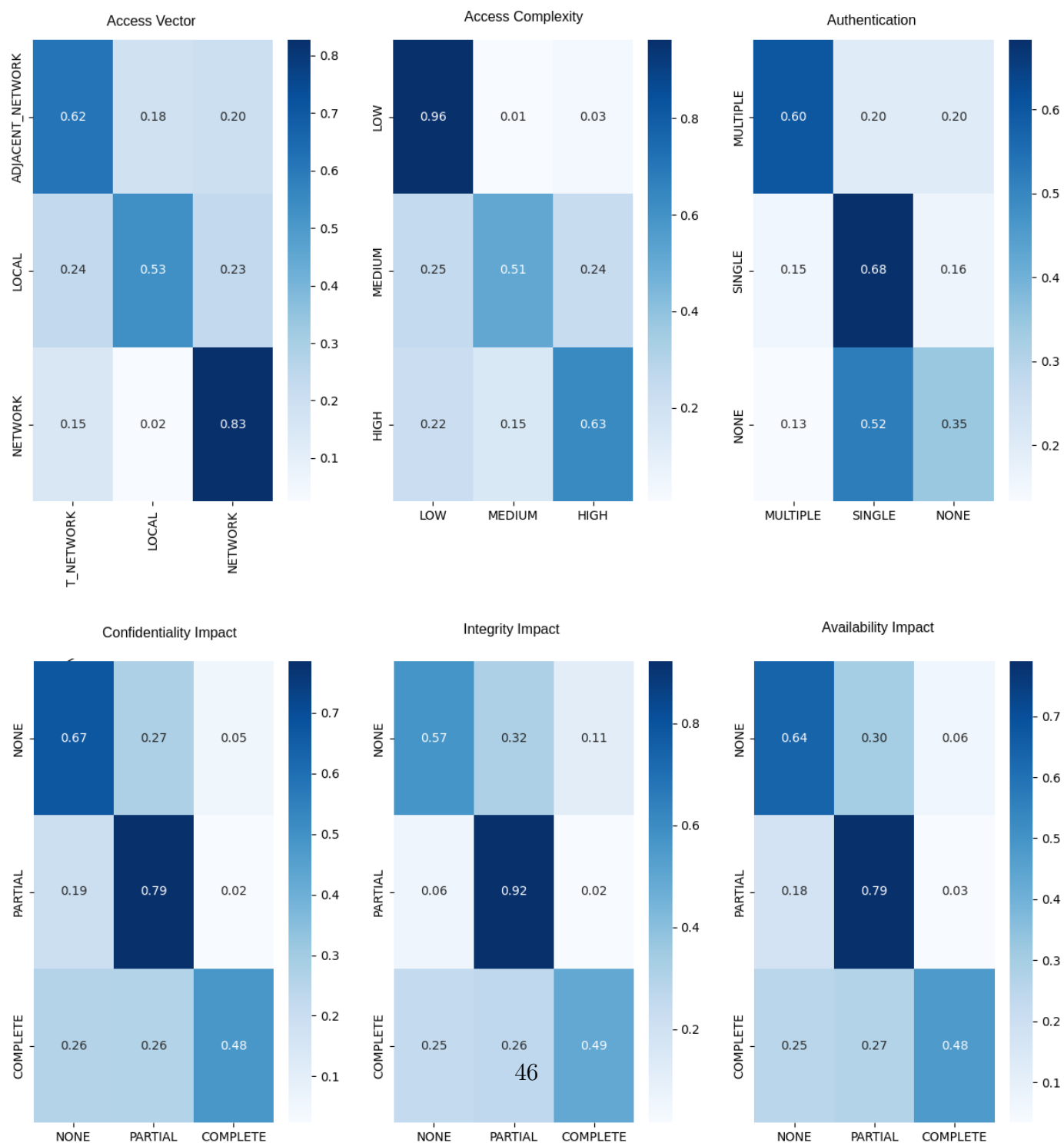


Figure 16: Confusion matrix of estimated accuracy for CVSS metrics for version 2.0 for MITRE

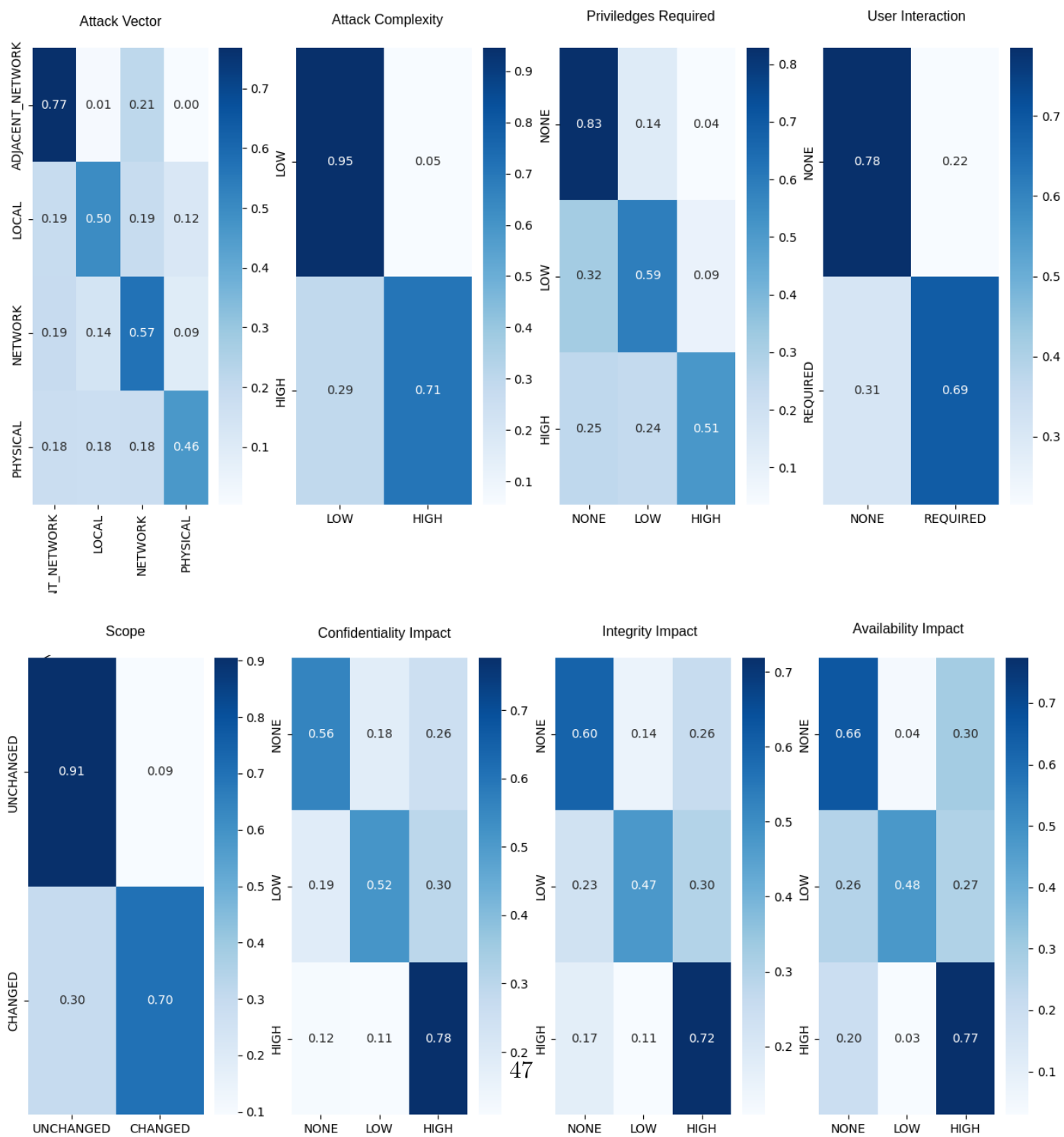


Figure 17: Confusion matrix of estimated accuracy for CVSS metrics for version 3.0 for NVD

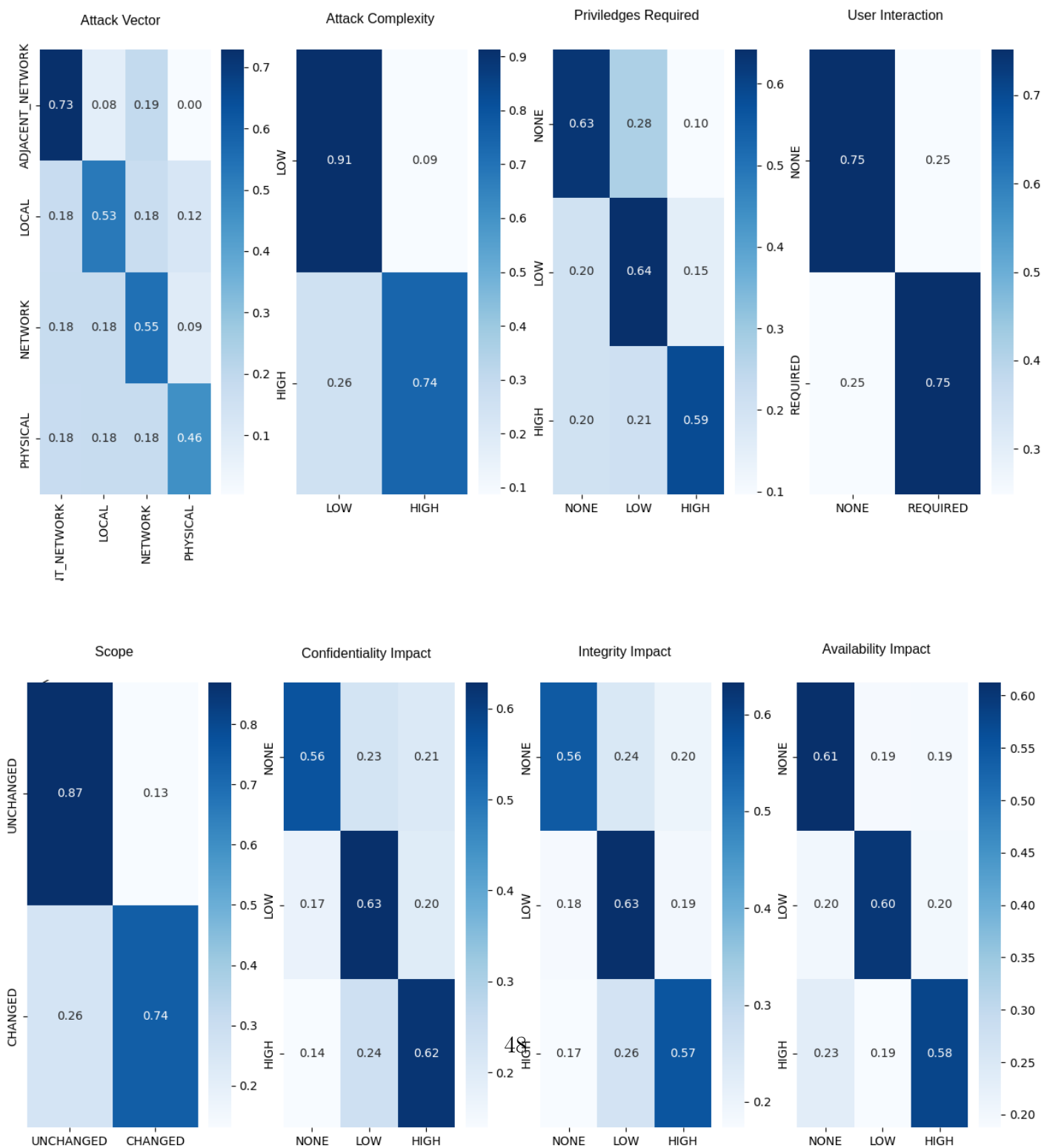


Figure 18: Confusion matrix of estimated accuracy for CVSS metrics for version 3.0 for MITRE

Appendix .2 CVSS 3.1 Base Score formula

Below is a formulaic representation of the CVSS 3.1 base score formula.

$$ISS = 1 - ((1 - Confidentiality) \times (1 - Integrity) \times (1 - Availability)) \quad (3)$$

$$Impact = \begin{cases} 7.52 \times (ISS - 0.029) - 3.25 \times (ISS - 0.02)^{15} & \text{if } Scope \text{ is } Changed \\ 6.42 \times ISS & \text{if } Scope \text{ is } Unchanged \end{cases} \quad (4)$$

$$Exploitability = 8.22 - AttackVector \times AttackComplexity \times PrivilegesRequired \times UserInteraction \quad (5)$$

$$BaseScore = \begin{cases} 0 & Impact \leq 0 \\ Roundup(Minimum(1.08 \times (Impact + Exploitability), 10)) & \text{if } Scope \text{ is } Changed \\ Roundup(Minimum((Impact + Exploitability), 0)) & \text{if } Scope \text{ is } Unchanged \end{cases} \quad (6)$$

Appendix .3 Exploit Prediction Scoring System diagrams for reference

Below are two graphs from FIRST, the creator of EPSS. These show some of the results they have found for this system, especially when comparing EPSS to CVSS in terms of efficiency of effort if you use EPSS to guide remediation of vulnerabilities.

Appendix A Aims and Objectives

Original

Aims The primary aim of this research is to develop sophisticated predictive models capable of accurately determining the severity levels of security threats based on the CVSS. This will involve a comprehensive review and comparison of current datasets, with a focus on leveraging natural language descriptions provided in security vulnerability reports. The project intends to utilize advanced transformer-based models to achieve this goal, contributing to the field of cybersecurity by enhancing the precision of threat severity assessments.

Objectives

- Conduct a comprehensive literature review to understand the current landscape of CVSS score prediction and the methodologies employed in existing models.
- Replicate successful methodologies to verify the accuracy of CVSS score databases, with a particular focus on alignment with recent CVSS standards and datasets.

Comparing Metrics: CVSS 7+ vs EPSS 10%+

Pulling EPSS and CVSS scores from October 1st, 2023 and measuring predictive performance at arbitrary thresholds against exploitation activity October 1-30, 2023. Data is limited to CVEs with CVSS 3.x scores published in NVD as of Oct 1, 2023.

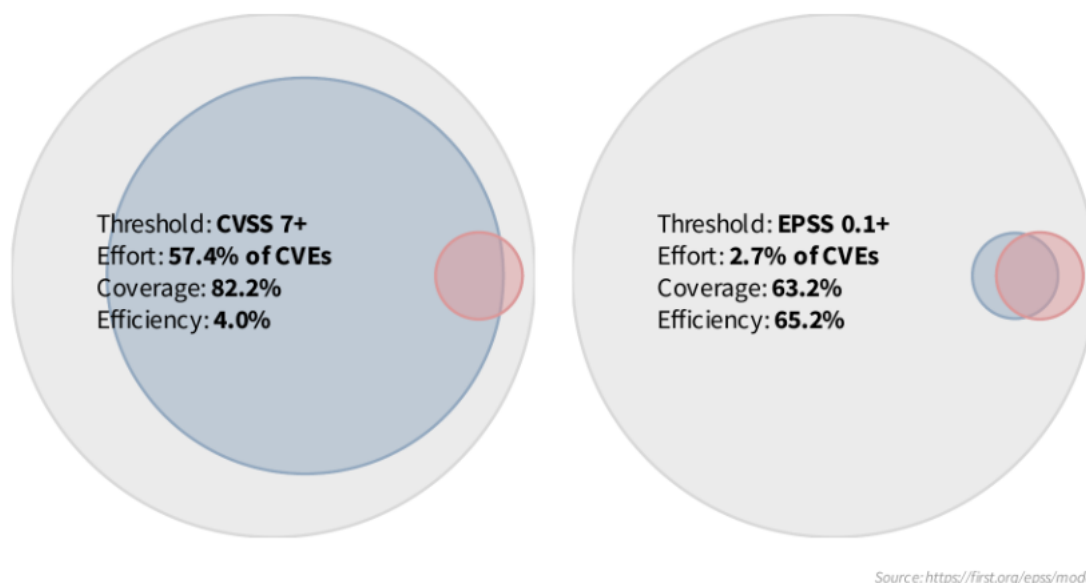


Figure 19: Comparing Metrics: CVSS 7+ vs. EPSS 10%+ sourced from [19]

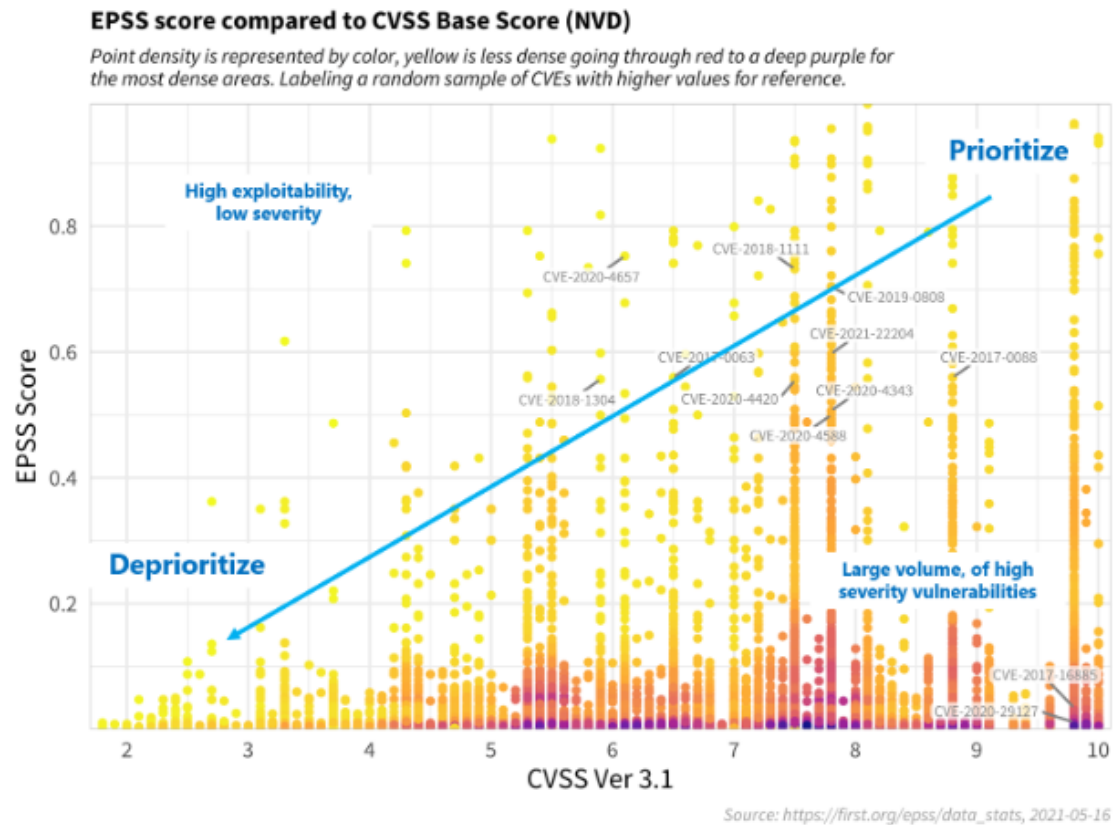


Figure 20: EPSS score compared to CVSS Base Score (NVD) sourced from [20]

- Explore opportunities for enhancing existing methodologies, including the investigation of data amalgamation from multiple databases to ascertain improvements in model performance.
- Experiment with various model architectures to identify the most effective approach in terms of predictive accuracy, specifically focusing on metrics such as the F1 score and balanced accuracy.

Timeline

- March: Initiate the project with a literature review, system environment setup, and resource gathering.
- March-April: Replicate existing methodologies to validate findings and ensure alignment with current standards.
- May-June: Generate preliminary results and compile an interim report detailing findings and methodologies.
- July-August: Conduct experiments with various data source combinations and model architectures to identify optimal configurations.
- September-October: Finalize experimental work, analyze results, and prepare the comprehensive final report.

Revised

Aims The primary aim of this research is to develop sophisticated predictive models capable of accurately determining the severity levels of security threats based on the CVSS. This will involve a comprehensive review and comparison of current datasets, with a focus on leveraging natural language descriptions provided in security vulnerability reports. The project intends to utilize advanced transformer-based models to achieve this goal, contributing to the field of cybersecurity by enhancing the precision of threat severity assessments.

Objectives

- Conduct a comprehensive literature review to understand the current landscape of CVSS score prediction and the methodologies employed in existing models.
- Replicate successful methodologies to verify the accuracy of CVSS score databases, with a particular focus on alignment with recent CVSS standards and datasets.
- Explore opportunities for enhancing existing methodologies, including the investigation of data amalgamation from multiple databases to ascertain improvements in model performance.

- Look into data cleaning and clustering, to improve the efficacy of the models, as well as a look into interpretability through data analysis.