

# Machine Learning Based Approach for the Automated Mapping of Discovered Vulnerabilities to Adversarial Tactics

Cameron Noakes

February 2022

## 1 Sentence Overview

The researchers converted each textual vulnerability description into a 512 fixed-dimensional numerical vector using Deep Averaging Network (DAN)-based USE (Universal Sentence Encoder) model. Exclusions of the characteristics of the vulnerabilities were not taken into account.

The defender needs to have a set of metrics that can be used, and automatically mapping vulnerabilities to potential attack tactics in the ATT&CK framework, this can allow security solutions to better respond to any vulnerability exploitations.

The researchers proposed a system that provides a multi label classification approach to automatically map vulnerabilities to the ATT&CK framework adversarial tactics. The proposed approach will help cyber defenders to prioritize their defense strategies.

Through the use of evaluating a set of machine learning algorithms (BinaryRelevance, LabelPowerset, ClassifierChains, MLKNN, BRKNN, RAKELd, NLSP, Neural Networks), it was concluded that ClassifierChains with Random-Forest are the most appropriate and is used to develop the proposed system.

## 2 Introduction

From the outline in the overview abstract, these security solutions should help simplify incident response (IR) and also support digital forensic investigations,

by providing an easier indication about the attacker mindset and the ATT&CK techniques behind it. This can be done numerous ways but the one the researchers chosen was to map exploited vulnerabilities to attackers' ATT&CK framework tactics.

This method can be used to detect security patterns and identify exploitable vulnerabilities. being able to accurately map information to possible attackers behavior and perspectives (in examples such as tactics and techniques), the defender will be able to know the possible upcoming outcomes of exploitation attacks.

An adequate defense strategy could be selected, or effective forensic analysis could be conducted based on the defender understanding attacker mindsets and perspectives as mentioned above.

The researchers converted each textual vulnerability description into a 512 fixed-dimensional numerical vector using Deep Averaging Network (DAN)-based USE (Universal Sentence Encoder) model. Exclusions of the characteristics of the vulnerabilities were not taken into account.

Not all of the existing CWEs are related to vulnerabilities in databases such as CVEs, and not all CAPEC IDs are related to the MITRE ATT&CK framework matrix, providing some incomplete datasets that will needed to be passed through data cleaning, more in-depth into this in the data preprocessing section.

The literature outlined how some related works researchers studied the use of two approaches (Doc2Vec and TF-IDF) to identify differences between the CAPEC textual descriptions and the CVE description text. It was proven that the TF-IDF is more accurate in regards to mapping any CAPEC-ID to CVE-ID.

## **3 Data Preprocessing**

### **3.1 Data Cleaning**

The datasets characterized by their incompleteness have to have some information removed due to them being incomplete, irrelevant, or duplicated data that could potentially negatively impact the performance of the machine learning algorithms. The researchers start this initial process by removing all the features

with data that may not be available (more than 80 percent) for all entries in the dataset. The obtained datasets include 27,471 samples where each one of them has 37 features.

Any null values in any incomplete datasets will be removed due to the issues that could arise when training and developing machine learning algorithms for the proposed solution. An obtained final dataset with 6,798 complete samples were concluded after removing all null values.

### **3.2 Data Encoding**

Due to the need for numerical values for the development and training of the Machine learning algorithms, the categorical textual descriptions and data needs to be encoded into numbers before applying to the chosen Machine Learning algorithms. Integer encoding was chosen by the researchers to provide the current problem with a solution. Integer encoding is where each categorical/textual value will be mapped to an integer.

Mathematical algorithmic equations are used to provide this solution using precision and recall variables to calculate things like Macro-F1 scores.

## **4 Conclusion**

In this literature, a proposition was made to take a Machine Learning-based solution to be able to then map the detected vulnerabilities and patterns of exploited vulnerabilities to attackers' tactics in the ATT&CK framework.

## **5 What problem does this solve?**

The given problem is that textual descriptions often hold a lot of data regarding the vulnerabilities and often are left out of models.

## **6 How do they solve it?**

The proposed system model uses machine learning algorithms and data encoding from textual descriptions to numerical to be able to accurately map vulnerability databases to ATT&CK tactics.

## **7 How do they assess it?**

There was no evaluation method or approach given for this literature piece.