AN IMPROVEMENT OF SOFTWARE VULNERABILITY AND CLASSIFICATION MODEL USING DEEP NEURAL NETWORK BASED ON BIG DATA

R.Kowshik, S.Sasmitha, N.Yuvaraj* And G.M.Sathyaseelan

Department of Information Technology, J K K Natraja College of Engineering and Technology, Komarapalayam-638183, Tamil Nadu, India.

ABSTRACT

Software vulnerabilities increase the risk of security breaches, potentially causing significant harm to systems. Automatic classification methods are essential for managing software vulnerabilities and enhancing system security. This project proposes an improved model named Automatic Vulnerability Classification Model (IGTF-DNN), which combines Information Gain based on Term Frequency (TF-IDF) and Deep Neural Network (DNN). By using the TF-IDF method to calculate the frequency and weight of words from vulnerability descriptions and IG to select features, an optimal set of feature words is obtained. The DNN then constructs an automatic classifier to effectively classify vulnerabilities. The model is tested using the National Vulnerability Database of the United States, showing better performance compared to the KNN model.

Keywords: Software Vulnerability, Deep Neural Network, Big Data, Information Gain, Term Frequency

1. INTRODUCTION

The rapid development of information technology has brought both convenience and significant security risks to various industries. Software vulnerabilities, which are defects in software or hardware, can be exploited by unauthorized individuals, posing severe risks to information systems. In 2017, Windows system vulnerabilities led to ransomware attacks on 100,000 organizations globally. The increasing number and variety of vulnerabilities highlight the importance of effective analysis and management to enhance system security and reduce attack risks. Traditional vulnerability classification methods have limitations, necessitating the development of automatic classification models using machine learning techniques.

The main objectives of the Vulnerability Classification are:

- Using TF-IDF to calculate the frequency and weight of each word from vulnerability descriptions.
- Using IG to select features to obtain an optimal set of feature words.
- Using a neural network model to construct an automatic vulnerability classifier for effective classification.

Specific objectives include:

- Extracting vulnerability information from records using program codes written in R.
- Collecting vulnerability data from 2015 to 2019 for training and testing.
- Processing data from multiple Excel worksheets and storing them as data frame objects.

2. LITERATURE SURVEY

2.1 Early Vulnerability Classification Methods

The RISOS method [1] is one of the earliest vulnerability classification approaches, targeting operating system vulnerabilities. This method categorized OS vulnerabilities into seven groups from the attack perspective but did not focus on how to exploit them.

2.2 Advanced Vulnerability Classification Methods

The PA vulnerability classification method [2] extended the focus to include application vulnerabilities. Andy Gray's method [3] introduced a ten-category system based on different analysis needs. However, as the complexity of vulnerabilities increased, traditional manual classification methods became less effective.

2.3 Machine Learning for Vulnerability Classification

Recent advancements in machine learning have shown promise in vulnerability classification. Shua et al. [4] applied the SVM classification method based on the LDA model, achieving good results in vulnerability grouping. Wijayasekara et al. [5] used the Naïve Bayes method for classifying textual information from error descriptions, illustrating its feasibility for this application. Gawron et al. [6] compared the Naïve Bayes algorithm with a simplified artificial neural network (ANN) algorithm, concluding that the ANN outperformed Naïve Bayes in vulnerability classification.

2.4 Limitations of Traditional Machine Learning Methods

Traditional machine learning algorithms often struggle with high-dimensional and sparse word vector spaces generated from vulnerability data. They also tend to overlook specific vulnerability information, resulting in lower classification accuracy. Recent studies have demonstrated the effectiveness of deep learning in various fields, including speech and image recognition, and natural language processing [7][8][9][10].

3. SYSTEM ANALYSIS

3.1 Existing System

The existing system uses the TF-IDF method to measure the importance of terms in documents and the Information Gain (IG) criterion to select features. However, it faces limitations in handling high-dimensional and sparse data, leading to reduced classification accuracy.

3.1.1 Drawbacks

- Fixed data set files limit the scope of input.
- Lack of backpropagation technique limits accuracy as hidden layer weights are not recalculated.
- Limited vulnerability categories are selected for classification.

3.2 Proposed System

The proposed system incorporates all existing methodologies with enhancements, including the use of deep neural networks (DNN) with backpropagation to improve accuracy. The figure represents feature selection and attack detection of IOT.



Fig 3.2.1 Dual Convolutional Neural Network Approach for Feature Selection and Attack Detection on Internet of Things Networks.

The system processes vulnerability details from 2014 to 2019 and classifies 27 categories of vulnerabilities.

3.2.1 Advantages

• Ability to process more data set files with new vulnerability types.

• Improved accuracy through the application of backpropagation technique.

• Increased number of vulnerability categories for classification.

3.3 Feasibility Study

The feasibility study assesses the economic, operational, and technical viability of the project.

3.3.1 Economic Feasibility

The system requires minimal hardware investment, making it cost-effective. The use of freely available technologies reduces overall costs.

3.3.2 Operational Feasibility

The system fits seamlessly into current organizational operations, solving existing problems efficiently.

3.3.3 Technical Feasibility

The system design ensures that technical requirements are met, enabling effective implementation and use of the TF-IDF and IG methods in R programming.

4. SYSTEM SPECIFICATION

4.1 Hardware Requirements

- **Processor**: Dual Core 2.1 GHz
- **RAM**: 2 GB
- Monitor: 17" Colour
- Hard disk: 500 GB
- Keyboard: Standard 102 keys
- Mouse: Optical mouse

4.2 Software Requirements

- Operating System: Windows 10 Pro
- Environment: R Studio 1.0
- Language: R 3.4.4

4.3 Software Description

4.3.1 Front End

The R programming language is used for statistical computing and graphics. It supports data manipulation, calculation, and graphical display, making it ideal for developing statistical software and data analysis applications.

5. SYSTEM IMPLEMENTATION

5.1 Modules

The system is divided into several modules, including Term Frequency-Inverse Document Frequency (TF-IDF), Information Gain (IG), Feature Words Extraction, and Optimizations using DNN.

5.1.1 Term Frequency-Inverse Document Frequency (TF-IDF)

TF-IDF measures the importance of words in a document relative to a corpus, helping to identify key terms.

5.1.2 Information Gain (IG)

IG selects features by measuring the importance of each feature in reducing uncertainty in class classification.

5.1.3 Feature Words Extraction

This module extracts feature words using TF-IDF and IG values to create a reduced, high-value feature set for classification.

5.1.4 Optimizations Using DNN

The DNN module includes one input layer, multiple hidden layers, and one output layer. It optimizes feature selection and classification accuracy through forward and backpropagation processes.

5.2 Input Design

Input design converts user-originated inputs into a computer-understandable format. It involves capturing, preparing, and ensuring the accuracy of data.

5.3 Output Design

Output design focuses on generating useful information for end-users. The outputs are designed to be attractive, convenient, and informative.

5.4 System Testing

System testing includes unit testing, integration testing, functional testing, system testing, and acceptance testing to ensure the system functions correctly.

5.4.1 Unit Testing

Unit testing involves testing individual components of the system to ensure they function correctly.

5.4.2 Integration Testing

Integration testing evaluates interactions between combined components.

5.4.3 Functional Testing

Functional testing verifies that the system functions as specified.

5.4.4 System Testing

System testing assesses the entire integrated system to ensure it meets all requirements.

5.4.5 Acceptance Testing

User acceptance testing (UAT) involves end-users testing the system to ensure it meets their needs.



Fig 5.4.5.1 Graph depicts the distribution of Vulnerability category

6. CONCLUSION AND FUTURE WORK

6.1 Conclusion

This project successfully applies deep neural networks to software vulnerability classification, demonstrating the effectiveness of the IGTF-DNN model in improving classification performance. The model outperforms traditional methods, providing a basis for future research using the benchmark vulnerability dataset.

6.2 Future Work

Future research will explore additional features and improve the model's performance. Potential areas include enhancing the neural network architecture and incorporating more comprehensive datasets.

7. APPENDIX

7.1 Source Code

The source code includes implementations of the sentiment analysis methods, data processing scripts, and the web application components.

7.2 Screenshots

Screenshots provide visual documentation of the system's user interface and key functionalities.

REFERENCES

- 1. Abbott, R. P., et al. (1976). Security Analysis and Enhancements of Computer Operating Systems. US Department of Commerce.
- Bisbey, I. R., & Hollingworth, D. (1978). Protection Analysis: Final Report. Univ. of Southern California.
- Gray, A. (2003). An historical perspective of software vulnerability management. Inf. Secur. Tech. Rep., 8(4), 34-44.
- 4. Shua, B., et al. (2013). Automatic classification for vulnerability based on machine learning. Proc. IEEE Int. Conf. Inf. Automat. (ICIA), 312-318.
- 5. Wijayasekara, D., et al. (2014). Vulnerability identification and classification via text mining bug databases. Proc. 40th Annu. Conf. IEEE Ind. Electron. Soc., 3612-3618.
- Na, S., et al. (2016). A study on the classification of common vulnerabilities and exposures using Naïve Bayes. Proc. Int. Conf. Broadband Wireless Comput. Commun. Appl., Springer, 657-662.
- Gawron, M., et al. (2017). Automatic vulnerability classification using machine learning. Proc. Int. Conf. Risks Secur. Internet Syst., Springer, 3-17.
- Deng, J., et al. (2009). ImageNet: A large-scale hierarchical image database. Proc. IEEE Conf. Comput. Vis. Pattern Recognit., 248-255.
- Russakovsky, O., et al. (2015). ImageNet large scale visual recognition challenge. Int. J. Comput. Vis., 115(3), 211-252.
- Xiong, W., et al. (2017). Toward human parity in conversational speech recognition. IEEE/ACM Trans. Audio Speech Lang. Process., 25(12), 2410-2423.
- 11. Krizhevsky, A., et al. (2012). ImageNet classification with deep convolutional neural networks. Proc. Adv. Neural Inf. Process. Syst., 1097-1105.
- 12. Silver, D., et al. (2016). Mastering the game of Go with deep neural networks and tree search. Nature, 529, 484-489.
- 13. Iyyer, M., et al. (2015). Deep unordered composition rivals syntactic methods for text classification. Proc. 53rd Annu. Meeting Assoc. Comput. Linguistics, 1681-1691.
- Jo, H., et al. (2016). Large-scale text classification with deep neural networks. Comput. Cognit., 23(5), 322-327.
- 15. Aziguli, W., et al. (2017). A robust text classifier based on denoising deep neural network in the analysis of big data. Sci. Program., 2017, 1-10.