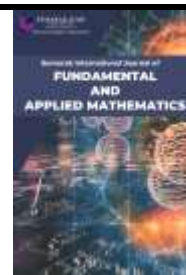




## Semarak International Journal of Fundamental and Applied Mathematics

Journal homepage:  
<https://semarakilmu.com.my/journals/index.php/sijfam/index>  
ISSN: 3030-5527



# Enhancing Asset Security in Malaysia: A Multivariate Regression and Time Series Analysis Approach

Nurul Azian Mohd Basharudin<sup>1</sup>, Hafiz Asyraf Ghani<sup>1</sup>, Mohammad Zaharudin Ahmad Darus<sup>2</sup>, Sahimel Azwal Sulaiman<sup>1,\*</sup>

<sup>1</sup> Centre for Mathematical Sciences, Universiti Malaysia Pahang Al-Sultan Abdullah, 26300 Gambang, Pahang, Malaysia

<sup>2</sup> Department Cyber Risk Intelligence, Cybersecurity Malaysia, 63000 Cyberjaya, Selangor, Malaysia

### ARTICLE INFO

#### Article history:

Received 18 December 2024

Received in revised form 17 January 2025

Accepted 24 February 2025

Available online 15 March 2025

#### Keywords:

Asset Security; multivariate regression; time series analysis Malaysia; risk management

### ABSTRACT

This research addresses the challenge of managing and securing external assets in an organization's digital infrastructure. As the attack surface grows due to factors like software updates, configuration changes, outdated security policies, and the addition of new assets, organizations become increasingly vulnerable to security threats. Without proper forecasting and analysis, these trends can lead to inefficiencies in resource allocation and expose critical assets to cyberattacks. The Autoregressive Integrated Moving Average model was employed to forecast changes in the external attack surface and prediction the quantity of total assets exposed over time. Next, multivariate linear regression was used to analyse the relationships between various factors. Influence diagrams was used to visualize the different factors, uncertainties, and decisions interact in the context of resource allocation and security planning. The results presented those certain factors, such as frequent software updates and the addition of new assets, significantly contributed to the expansion of the attack surface. Then, the strongest predictors of asset exposure were identified, which allowed for more targeted interventions. The influence diagrams provided a clear, visual representation of how these factors interact, aiding in the understanding of complex security scenarios. This research analysing critical relationships with multivariate linear regression, organizations can better allocate resources to mitigate risks.

## 1. Introduction

Today's IT ecosystems encompass numerous endpoints and assets spread across diverse locations and devices, including core networks, regional offices, subsidiaries, and third-party hosting providers. This dispersion presents challenges in monitoring and increases the risk of exploitation by adversaries. Failing external assets effectively can lead to unauthorized access, data breaches, service disruptions, malware infections, and reputational loss. Determining the external attack surface is crucial, especially for government entities aiming to bolster the security of digital assets facing the internet [1,2]. External Attack Surface Management (EASM) plays a vital role in identifying, assessing,

\* Corresponding author.

E-mail address: [sahimel@ump.edu.my](mailto:sahimel@ump.edu.my)

<https://doi.org/10.37934/sijfam.5.1.110>

and managing risks associated with internet-facing assets and systems. By deploying EASM, organizations can evaluate their external attack surface, which comprises potential entry points or weaknesses accessible from the outside world and exploitable by adversaries. This evaluation enables preemptive security measures to mitigate potential threats [3-5].

There are various approaches to analyzing external attack surfaces, including utilizing the STRIDE framework, Risk-Based Vulnerability Management (RBVM), and Attack Tree Analysis. According to [6], STRIDE represents several exploitable security threats and suggests a five-phase threat analysis process like threat trees. The framework for determining quantitative vulnerability measures is based on the attack tree model [7]. Meanwhile, RBVM is used to identify and address vulnerabilities based on the risk they pose to the organization. Remediation efforts are prioritized by combining threat intelligence, asset criticality, and vulnerability data.

Time series analysis is employed in this research to determine trends in exposed vulnerabilities and forecast future trends, enabling organizations to proactively reduce vulnerabilities and enhance security [8-10]. By leveraging Multivariate Linear Regression (MLR), correlations between various independent and dependent variables are uncovered. In [11-13], additionally, Influence Diagrams are constructed to visually represent decision-making processes and causal relationships within the security framework. Through these methodologies, this research aims to enhance risk assessment strategies and empower decision-makers to proactively manage cybersecurity challenges.

This research consists of four sections: introduction, methodology, results and discussion, and conclusion. A brief explanation of the research is provided in Section 1 (Introduction). The techniques employed in this study, including influence diagrams, ARIMA, and multivariate linear regression (MLR), are described in Section 2. The results of each method are discussed in Section 3, while the concise conclusion is presented in Section 4."

## 2. Methodology

This section discussed methodology on data requirements, Multivariate Linear Regression (MLR), Autoregressive Integrated Moving Average (ARIMA), Influence Diagram and Evaluation Method.

### 2.1 Data Requirements

This research employs MLR, ARIMA and Influence Diagrams to analyze and forecast total assets and total exposed assets. There are seven variables used in this research, two are dependent variables (Total Exposed, Total Asset) and the others include Low, Medium, High and Critical are categorized as independent variables. By leveraging MLR, correlations between various independent and dependent variables are uncovered, shedding light on factors influencing total exposed and total asset values. ARIMA is applied to forecast significant independent variables and then the ARIMA output is used in predicting dependent variables by using multivariate linear regression model. Additionally, Influence Diagrams are constructed to visually represent decision-making processes and causal relationships within the security framework.

Data for this study is sourced from the CrowdStrike application [14], which provides comprehensive insights into an organization's external assets categorized by risk levels. The data acquired for this project contains information related to risk assessment or exposure of certain assets. The dataset includes several key attributes, each with a specific explanation. The *Date* attribute indicates when the data was collected. *Total Exposed* represents assets that are visible or accessible to potential attackers, making them vulnerable due to their exposure on external networks or to unauthorized users. *Total Asset* refers to the total number of components within the system,

including servers, workstations, databases, applications, network devices, and any other essential infrastructure elements. *Low Risk Level* signifies assets with minimal vulnerabilities, posing less critical risks. *Medium Risk Level* indicates assets with moderate vulnerabilities that could impact organizational security if exploited. *High Risk Level* represents assets with significant vulnerabilities, which could have serious security implications if compromised. Finally, *Critical Risk Level* designates essential assets whose breach would immediately and severely impact project operations, data integrity, or confidentiality.

## 2.2 Multivariate Linear Regression (MLR)

In this research, MLR is used when there is more than one dependent variable and independent variables. It can identify the relationship between several independent variables and a dependent variable. The equation for MLR given as follows:

$$Y = \beta_0 + \beta_1x_1 + \beta_2x_2 + \dots + \beta_nx_n \quad (1)$$

where  $x_1, x_2, \dots, x_n$  is independent variables,  $Y$  is a dependent variable,  $\beta_0$  is intercept and  $\beta_1, \beta_2, \dots, \beta_n$  is coefficient. Since there are seven variables in this project, two of variables are dependent variables, namely Total Exposed and Total Asset. Meanwhile, all the variables except for the date are categorized as independent variables. These independent variables include low, high, medium, and critical categories.

## 2.3 Autoregressive Integrated Moving Average (ARIMA)

This project fully utilizes the ARIMA method due to its suitability for the characteristics of our dataset. ARIMA is equipped to handle both autocorrelation and noise in the data. ARIMA can capture correlations between observations and take irregular fluctuations into consideration, even with a smaller sample. Even though having more data is usually preferable, 31 rows may be enough to fit a simple ARIMA model. Given the dataset's size, ARIMA also is suitable for short-term forecasting. It can provide meaningful insights into short-term future values by leveraging its components effectively. The general ARIMA model is denoted as ARIMA ( $p, d, q$ ), where  $p$  is the number of lagged observations (autoregressive terms),  $d$  is the degree of differencing required to make the series stationary and  $q$  is the number of lagged forecast errors (moving average terms). The ARIMA model can be represented by the following formula:

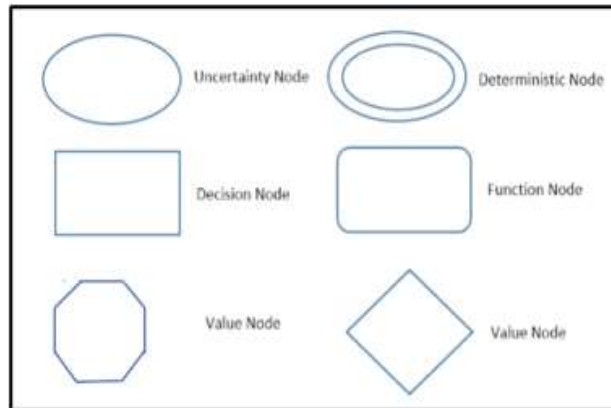
$$Y_t = \phi_1Y_{t-1} + \phi_2Y_{t-2} + \dots + \phi_pY_{t-p} + \varepsilon_t \quad (2)$$

where  $Y_t$  is the value at time  $t$ ,  $\phi$  are the parameter to be estimated,  $p$  is the number of lagged observations and  $\varepsilon_t$  is the error term.

## 2.4 Influence Diagram

An influence diagram is a simple visual representation of a decision-making process. The essential components are shown as nodes of various forms and colors, including decisions, uncertainties, and objectives. As arrows, it depicts the influences between them. However, nodes are not required to

be used to make a suitable influence diagram. It depends on the available variables. Figure 1 presented the list of nodes in the influence diagram.



**Fig. 1.** List of nodes in influence diagram

### 2.5 Evaluation Method

Mean Absolute Percentage Error (MAPE), Root Mean Square Error (RMSE) and Mean Absolute Error (MAE) are applied to performance of the model [15]. The statistic known as Mean Absolute Percentage Error (MAPE) is frequently used to assess a forecasting model's accuracy. It calculates the average percentage difference, represented as a percentage, between the anticipated values and the actual values. The equation for MAPE is as follows.

$$MAPE = \frac{1}{n} \sum_{t=1}^n \left| \frac{Y_t - \hat{Y}_t}{Y_t} \right| \times 100 \tag{3}$$

where  $n$  is the number of observations,  $Y_t$  is the actual value of the time series at time,  $t$  and  $\hat{Y}_t$  is the forecasted value at time. A typical statistic for assessing a forecasting or regression model's accuracy MAE. It calculates the average size of the discrepancies between the expected and actual values. The equation is given as follows.

$$MAE = \frac{1}{n} \sum_{t=1}^n |Y_t - \hat{Y}_t| \tag{4}$$

where  $n$  is the number of observations in the dataset,  $Y_t$  is the actual value of the time series at time and  $\hat{Y}_t$  is the forecasted value at time. MAE is frequently used to evaluate model performance and contrast various methodologies in a variety of fields, including regression analysis, machine learning, and time series forecasting. Better accuracy is indicated by a lower MAE since it suggests less differences between actual and anticipated values.

The accuracy and precision of a forecasting or regression model are frequently assessed using the RMSE statistic. It calculates the average square root of the disparities between the expected and actual values. The RMSE equation is as follows:

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2} \tag{5}$$

where  $n$  is the number of observations in the dataset,  $Y_i$  is the actual value of the time series at time  $i$  and  $\hat{Y}_i$  is the forecasted value at time  $i$ . The square of the differences gives greater errors more weight. The average squared differences can be translated and made easier to compare across different datasets or models by taking the square root of the average squared differences. Since smaller variances between projected and actual values are implied by a lower RMSE, it denotes more accuracy and precision.

### 3. Results

This section presents the results of our ARIMA and MLR models, with a primary focus on predicting the external attack surface. Additionally, MLR is used to assess feature correlations and employs an influence diagram to illustrate causal relationships, aiding in risk assessment.

#### 3.1 Multivariate Linear Regression

The main goals of employing multivariate linear regression are twofold. First, to identify the most effective model that includes only the significant variables. Second, to use this model to predict Total Exposed and Total Assets. Figure 2 illustrates the relationships between each independent variable and the dependent variables (Total Exposed and Total Asset).

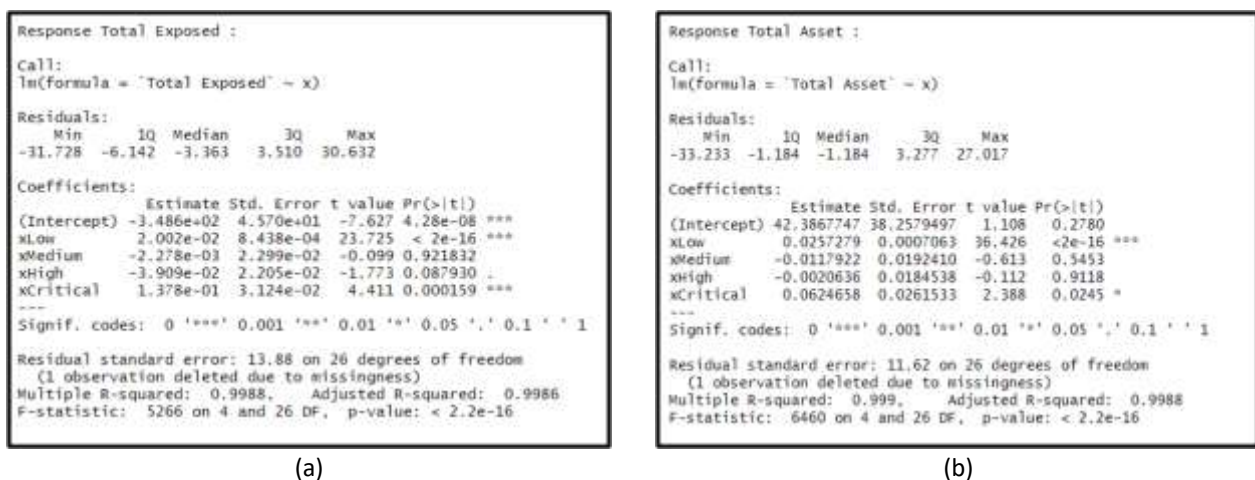


Fig. 2. The relationship for each independent variable to (a) Total Exposed (b) Total Asset

Based on Figure 2, the  $p$ -values for both dependent variables are 0.0000, which means that the model is significant where at least one of the independent variables is related to the dependent variable. The significant variables for Total Exposed and Total Assets are  $lq$  and critical. It is because their  $p$ -value is less than  $\alpha$  (0.05).

Figure 3(a) shows that three variables—Low, High, and Critical—fit the developed model, as their  $p$ -values are less than  $\alpha$  (0.05). The reduced model of MLR is given as follows.

$$y_1 = -348.6 + 0.02002x_1 - 0.0309x_3 + 0.1378x_4$$

$$y_2 = 42.39 + 0.0257x_1 - 0.0021x_3 + 0.0625x_4 \tag{6}$$

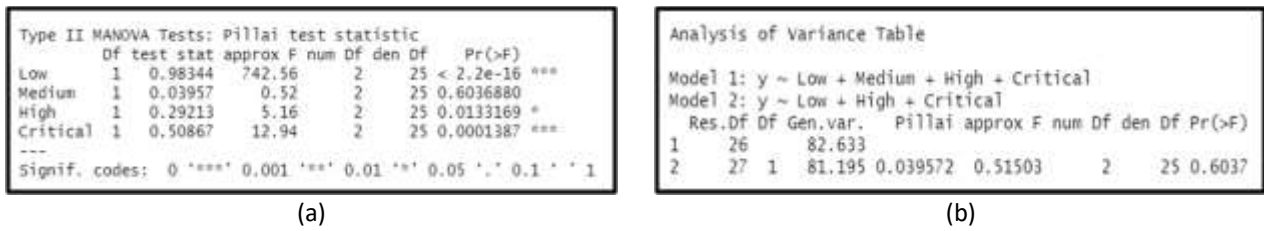


Fig. 3. The output for (a) MANOVA (b) ANOVA

Next, the difference between the full model and the reduced model is evaluated using ANOVA. Based on Figure 3(b), the full model (mvlr) does not offer a statistically significantly better fit than the reduced model (mvlr1) since its p-value exceeds  $\alpha$  (0.05). Therefore, the best multivariate model is the reduced model, which includes only the significant variables. Thus, the forecasting process will use five variables: Total Exposed, Total Asset, Low, High, and Critical. Here, Total Exposed and Total Asset are the dependent variables, while Low, High, and Critical serve as the independent variables.

Finally, the trained multivariate linear regression model is applied to predict Total Exposed and Total Asset. Multivariate linear regression enables analysis of the relationships between multiple independent variables (Low, High, and Critical) and multiple dependent variables (Total Exposed and Total Asset). By incorporating the forecasted values of Low, High, and Critical from Section 3.2 as input features, the model leverages these variables to make predictions about Total Exposed and Total Asset. Figure 4 presented the output of forecasted dependent variable.

	Date	Predicted_Total_Exposed	Predicted_Total_Asset
	<date>	<dbl>	<dbl>
1	2023-05-31	1332.429	2041.076
2	2023-06-01	1334.671	2044.516
3	2023-06-02	1336.289	2047.001
4	2023-06-03	1337.458	2048.795
5	2023-06-04	1338.302	2050.091
6	2023-06-05	1338.912	2051.027
7	2023-06-06	1339.352	2051.703
8	2023-06-07	1339.670	2052.191
9	2023-06-08	1339.900	2052.544
10	2023-06-09	1340.066	2052.798

Fig. 4. Forecast output for dependent variables

### 3.2 AutoRegressive Integrated Moving Average (ARIMA)

For ARIMA model, the dataset used for conducting SES consists of 31 rows, with 24 rows used for training and seven rows for testing. First, ARIMA models were used to fit the time series data for the Low, High, and Critical variables. Model fit refers to how well the selected ARIMA model aligns with the observed data, capturing the patterns, trends, and variability present in the dataset. Model fit is evaluated using metrics such as log likelihood, AIC, AICc, and BIC, as well as by comparing predicted values to actual data points. A well-fitted model provides reliable predictions and insights, while a poorly fitted model may produce inaccurate or misleading results. Table 1 presents the output of the fitting process.

**Table 1**  
 The Output of ARIMA model

Variables	Output	Explanation
Low	<pre>Series: data\$Low ARIMA(1,0,0) with non-zero mean  Coefficients:       ar1      mean     0.7222  53211.147 s.e.  0.1156  3035.905  sigma^2 = 27444328: log likelihood = -308.8 AIC=623.6  AICC=624.49  BIC=627.9</pre>	i) ARIMA model used is ARIMA(1,0,0) ii) The autoregressive coefficient (ar1) is 0.7222, indicating the strength and direction of the relationship between current and past values in the time series. iii) The standard errors ar1 is 0.1156.
High	<pre>Series: data\$High ARIMA(0,1,0)  sigma^2 = 22974020: log likelihood = -296.82 AIC=595.63  AICC=595.78  BIC=597.03</pre>	i) ARIMA model used is ARIMA(0,1,0) ii) An ARIMA model with no AR or MA terms (ARIMA(0,1,0)) indicates that the model does not consider past values or forecast errors when predicting future values.
Critical	<pre>Series: data\$Critical ARIMA(0,1,0)  sigma^2 = 6324554: log likelihood = -277.47 AIC=556.93  AICC=557.08  BIC=558.34</pre>	iii) Hence, there are no model coefficients or standard errors associated.

The next output presents the forecast for the *Low*, *High*, and *Critical* variables, based on the ARIMA models fitted to the training data. Table 2 displays the first 10 forecasted rows for each significant variable. From the table, the predicted values for high and critical produce a flat forecast, meanwhile the predicted value for low variable produce an increasing trend. Then this output is used as the as input variables in regression models (Eq. (6)). Therefore, by utilizing the information from these ARIMA models, we can forecast Total Exposed and Total Asset by using multivariate linear regression.

**Table 2**  
 Forecast output for significant independent variables

	Variable		
	Low	High	Critical
32	52704	30818	15509
33	52844	30818	15509
34	52946	30818	15509
35	53020	30818	15509
36	53073	30818	15509
37	53111	30818	15509
38	53139	30818	15509
39	53159	30818	15509
40	53173	30818	15509
41	53184	30818	15509

### 3.3 Influence Diagram

Next, the results on influence diagram. Figure 5 depicts an influence diagram created from a dataset containing two deterministic nodes (Total Assets and Exposed Assets) and five decision nodes (Low, Medium, High, Critical, and Risk Score). Deterministic nodes have fixed values that are known but may change over time due to factors such as system upgrades, asset additions or removals, network configuration changes, or updates to security controls. In contrast, decision nodes represent choices aimed at managing and reducing the attack surface. These decisions focus on minimizing vulnerabilities and decreasing the likelihood of successful attacks, involving actions such as risk mitigation, prioritization, change management, and resource allocation.

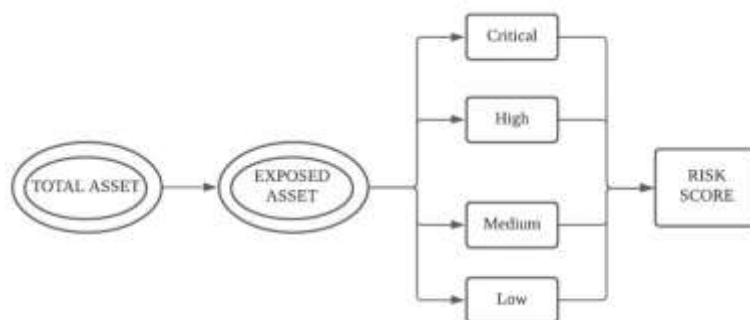


Fig. 5. Influence diagram

The *Risk Score* decision node helps prioritize and identify areas of the attack surface that require immediate attention and resource allocation. For *Low*-risk decisions, actions include regular monitoring, implementing baseline security procedures, and ensuring ongoing maintenance. *Medium* risk decisions may involve adding extra security measures, performing periodic vulnerability assessments, and improving surveillance and incident response capabilities. *High* risk decisions typically require allocating additional resources, enforcing stricter security controls, and performing rapid remediation to minimize risk. Finally, *Critical* risk decisions call for urgent security audits, deploying robust security controls, assigning dedicated personnel for continuous monitoring and response, and conducting thorough vulnerability assessments.

By employing this influence diagram, the impact of different decisions on the security of external assets can be systematically assessed. The deterministic nodes provide a stable basis for understanding the current state of assets, while the decision nodes guide the actions required to manage and reduce risks effectively. This structured approach facilitates a comprehensive analysis, enabling the identification of critical areas that require immediate attention and the formulation of targeted strategies to enhance security measures.

### 3.4 Evaluation Model

The model evaluations used in this project include MAPE, MAE, and RMSE for all risk levels, applied to a multivariate linear regression model. Table 3 shows that Total Asset recorded the lowest MAPE value, indicating that the forecasts have a relatively low error rate of 0.35% on average compared to the actual data. This suggests that the model's predictions for Total Asset are highly accurate, with minimal variance from the true values. In contrast, Total Exposed recorded a higher MAPE value than Total Asset, with a deviation of 12.99% from the true values. This result indicates a relatively greater inaccuracy and high error, suggesting that the model needs improvement to achieve more precise predictions.

Furthermore, Total Asset also recorded the lowest MAE value, at 9.071, meaning the projections deviate from the actual values by about 9.071 units on average. This level of deviation indicates moderate accuracy, suggesting that while the predictions are reasonably close, there is still room for improvement. However, Total Exposed recorded a much higher MAE value of 199.66, meaning the model's predictions deviate by around 199.66 units on average from the actual values, indicating a significant discrepancy and poor accuracy.

Additionally, Table 3 shows that the Low Risk Level recorded the lowest error, while the Medium Risk Level recorded the highest. The RMSE of 301.81 indicates that the average squared discrepancies between the model's forecasts and the actual values result in an error of 301.81 units, representing a significant deviation. This suggests that the model's performance needs improvement to reduce overall errors and produce more accurate predictions. For Total Exposed, the RMSE is 199.98, which



indicates an average deviation of approximately 199.98 units from the actual values. Meanwhile, Total Asset recorded an MAE of 7.11, indicating a much smaller average deviation and suggesting higher predictive accuracy.

**Table 3**  
The Model Evaluation

Evaluation	Total Exposed	Total Asset
MAPE	12.99%	0.35%
MAE	199.66	9.071
RMSE	199.98	7.11

#### 4. Conclusions

In this paper, three key objectives were pursued. First, using MLR, the goal was to uncover relationships between independent and dependent variables. This analysis revealed moderate to strong positive correlations and the absence of weak ones, resulting in the identification of the best model, validated through ANOVA testing, which included just three explanatory variables. Second, the project successfully achieved its objective of forecasting the dependent variable by using ARIMA results as inputs for the MLR model. The evaluation metrics, including RMSE, MAPE, and MAE, indicate varying levels of accuracy across different categories. While the MAPE values are relatively low—0.35% being the lowest for Total Asset—indicating a high level of accuracy, higher RMSE and MAE values for Total Exposed suggest limitations in the model's accuracy for predicting vulnerability levels within government organizations. Finally, the third objective was to represent deterministic relationships among variables through an influence diagram, with Total Asset and Total Exposed as deterministic nodes and Low, Medium, High, and Critical as decision nodes linked to the Risk Score. This understanding of deterministic nodes has the potential to enhance the efficient protection of external government assets.

#### Acknowledgement

This research was not funded by any grant.

#### References

- [1] Theisen, Christopher, Nuthan Munaiah, Mahran Al-Zyoud, Jeffrey C. Carver, Andrew Meneely, and Laurie Williams. "Attack surface definitions: A systematic literature review." *Information and Software Technology* 104 (2018): 94-103. <https://doi.org/10.1016/j.infsof.2018.07.008>
- [2] Bradbury, Matthew, Carsten Maple, Hu Yuan, Ugur Ilker Atmaca, and Sara Cannizzaro. "Identifying attack surfaces in the evolving space industry using reference architectures." In *2020 IEEE Aerospace Conference*, pp. 1-20. IEEE, 2020. <https://doi.org/10.1109/AERO47225.2020.9172785>
- [3] Zhang, Ying, Xin Guo, Ran Liu, and Haibo Zhang. "Research on Network Security Trend Prediction Based on Exponential Smoothing Algorithm." In *2020 IEEE 11th International Conference on Software Engineering and Service Science (ICSESS)*, pp. 507-510. IEEE, 2020. <https://doi.org/10.1109/ICSESS49938.2020.9237658>
- [4] Oruma, Samson O., Mary Sánchez-Gordón, Ricardo Colomo-Palacios, Vasileios Gkioulos, and Joakim K. Hansen. "A systematic review on social robots in public spaces: Threat landscape and attack surface." *Computers* 11, no. 12 (2022): 181. <https://doi.org/10.3390/computers11120181>
- [5] Straub, Jeremy. "Modeling attack, defense and threat trees and the cyber kill chain, att&ck and stride frameworks as blackboard architecture networks." In *2020 IEEE International conference on smart cloud (SmartCloud)*, pp. 148-153. IEEE, 2020. <https://doi.org/10.1109/SmartCloud49737.2020.00035>
- [6] Anand, Pooja, Yashwant Singh, Arvind Selwal, Pradeep Kumar Singh, Raluca Andreea Felseghi, and Maria Simona Raboaca. "Iovt: Internet of vulnerable things? threat architecture, attack surfaces, and vulnerabilities in internet of things and its applications towards smart grids." *Energies* 13, no. 18 (2020): 4813. <https://doi.org/10.3390/en13184813>

- [7] Ten, Chee-Wooi, Chen-Ching Liu, and Manimaran Govindarasu. "Vulnerability assessment of cybersecurity for SCADA systems using attack trees." In *2007 IEEE Power Engineering Society General Meeting*, pp. 1-8. IEEE, 2007. <https://doi.org/10.1109/PES.2007.385876>
- [8] Hyndman, R. J. *Forecasting: principles and practice*. OTexts, 2018.
- [9] Shumway, Robert H., David S. Stoffer, and David S. Stoffer. *Time series analysis and its applications*. Vol. 3. New York: springer, 2000. <https://doi.org/10.1007/978-1-4757-3261-0>
- [10] Nugus, Sue. *Financial planning using Excel: forecasting, planning and budgeting techniques*. Butterworth-Heinemann, 2009. <https://doi.org/10.1016/B978-1-85617-551-7.00015-X>
- [11] Ndilikilikesha, Pierre C. "Potential influence diagrams." *International Journal of Approximate Reasoning* 10, no. 3 (1994): 251-285. [https://doi.org/10.1016/0888-613X\(94\)90003-5](https://doi.org/10.1016/0888-613X(94)90003-5)
- [12] Pearl, Judea. "Influence diagrams—Historical and personal perspectives." *Decision Analysis* 2, no. 4 (2005): 232-234. <https://doi.org/10.1287/deca.1050.0055>
- [13] Rios Insua, David, Aitor Couce-Vieira, Jose A. Rubio, Wolter Pieters, Katsiaryna Labunets, and Daniel G. Rasines. "An adversarial risk analysis framework for cybersecurity." *Risk Analysis* 41, no. 1 (2021): 16-36. <https://doi.org/10.1111/risa.13331>
- [14] CrowdStrike: Stop breaches. Drive business. 2023.
- [15] Swanson, David A. "On the relationship among values of the same summary measure of error when used across multiple characteristics at the same point in time: an examination of MALPE and MAPE." *Review of Economics and Finance* 5, no. 1 (2015).