# Phishing Detection for Secure Operations of UAVs

Jan Bohacik

Department of Informatics
Faculty of Management Science and Informatics, University of Zilina
Zilina, Slovakia
Jan.Bohacik@uniza.sk

*Abstract*—**Unmanned Aerial Vehicles (UAV) or drones are used in various domains more and more, including military operations, monitoring, rescue of victims, and transport. Often, UAV resources are developed as web services so that they can be accessed anywhere on the Internet through the World Wide Web. However, this makes them vulnerable to phishing activities of criminals who may try to access these resources and other sensitive information. Therefore, the development of a phishing detection tool based on data mining is presented in this paper. It consists of a browser extension monitoring visited webpages and a backend communicating with the browser extension for the purposes of executing some specific tasks. The browser extension is implemented in JavaScript and the ReactJS framework, and it contains an implementation of classifications with a Bayesian network, decision tree, nearest neighbor classifier and neural network. The backend uses PHP, Python scripts and the Apache HTTP Server. In addition, a browser extension is implemented so that data about webpages can be collected and this data is used for the creation of data mining models. Experimental validation with 10-fold cross-validation and through the browsing of real-world websites show promising results in phishing detection.**

*Keywords—UAV; security; phishing; data mining*

## I. INTRODUCTION

Flying vehicles which do not have a human pilot on board are called Unmanned Aerial Vehicles (UAV) or drones [2]. Traditionally, they were used for tasks in military operations whose capabilities were beyond the capabilities of many. But recent advances in technology and reductions in costs have allowed the spread of UAVs to other areas as well. And so, in addition to military operations as shown in [7], they are currently used in environmental and other monitoring [14], rescue of victims [12], and transport [4]. UAVs use waves on specific radio frequencies for the communication with the ground station so that they receive or send commands. This leads to a situation when the user is restricted to a location, which is not desirable for applications placed anywhere by any user. Therefore, concepts for the inclusion of UAVs into the cloud infrastructure have been developed recently [8]. This inclusion makes the UAV and its resources available to the client ubiquitously. This client could be a person who uses a browser on the Internet or another UAV. In general, cloud infrastructure has servers which are very powerful computers providing services in the cloud. The UAV is included into the cloud as a backend server which provides its services and resources. However, this is a common client-server system which may be vulnerable to phishing techniques in the same way as Internet banking for example.

Phishing is defined as the fraudulent attempt to obtain sensitive data or information by passing oneself off as a trustworthy entity in some digital communication [11]. This data or information may include usernames, passwords, or other sensitive details. It is often accompanied with the redirection of the user to a fake webpage that matches the look and feel of the legitimate one. Fake webpages are easy to make nowadays due to the existence of available user-friendly tools [5]. This is alarming for the users of UAVs who access them and their resources through webpages in a browser, and some anti-phishing approaches are required. Several anti-phishing approaches can be found and each is based on something different [13]. The oldest ones are blacklist-based approaches which check if webpages are in lists of forbidden URL links. However, new fake webpages can appear easily and so the functionality of them is limited these days. Heuristic-based approaches look at the URLs of webpages and apply heuristics for checking their features. Content-based approaches look for some terms in webpages or images in screenshots of webpages. KDD-based approaches where KDD is Knowledge Discovery in Data use attributes extracted from webpages for the creation and use of data mining classification models. There are also hybrid-based approaches which combine several of the approaches. The most complex and potentially most effective approaches are KDD-based ones and their merges in hybrid-based approaches. In addition, KDD-based approaches are actively researched at this moment [3][6][9][10]. Among them, the most used methods are classification methods. These methods use data about instances described by describing attributes and classified into the class attribute for the creation of a data mining model which is used for the classification of new instances. In this paper, a Bayesian network, a decision tree, a nearest neighbor classifier and a neural network are utilized in the development of a browser extension for monitoring visited webpages, of a browser extension for data collection and of a backed communicating with these extensions for the purposes of executing some specific tasks related to describing attributes.

The paper is organized in a way which is explained in this paragraph. The developed browser extension for monitoring visited webpages, browser extension for data collection and backed communicating with these extensions are presented in Section II. In Section III, the webpage data collected with the browser extension for data collection is described in detail. The creation of data mining models with the collected webpage data and carried experimental evaluation are analyzed in Section IV. Section V contains summarized conclusions of results.

## II. BROWSER EXTENSIONS AND THE BACKEND

There are two developed browser extensions which have been implemented for the purposes of phishing detection. In general, a browser extension is a module for an Internet browser and it is used for the customization of the browser. One of the developed extensions is for further advancement and for data mining specialists who can collect data about visited webpages with it and assign if particular webpages are phishing or legitimate. In addition, data mining specialists use this data collected about webpages for the development of data mining classification models. These models are utilized in the other developed browser extension which is used for monitoring visited webpages. A screenshot of this extension displayed in a browser is shown in Fig. 1. Currently, four types of created classification data mining models can be imported into the extension: a) a Bayesian network marked as Naive Bayes network in Fig. 1; b) a decision tree marked as j48 decision tree; c) a nearest neighbor classifier marked as NNge; and d) a neural network marked as MLP. During the importation, the models are transformed automatically from files containing representations taken from Weka into scripts in JavaScript, where Weka [15] is an open-source data mining tool. Both of the browser extensions are written in JavaScript and the frontend uses the ReactJS framework for the creation of their user interface.

Another important part of the implementation is the backend. It communicates with the browser extensions through the HTTP GET method and runs on a Linux server with PHP, Python scripts and the open-source Apache HTTP Server. For a successful execution of the backend, packages apache, sqlite, mysql, php7.*, and python have to be installed. Dependencies for Python scripts are in files 'requirements.txt' and they can be installed with command 'pip install -r requirements.txt'. Data about webpages are stored with the open-source MySQL relational database management system. Some tasks related to the data collection about webpages which cannot be done in the browser extensions are performed in the backend as well. In addition, the backed allows to collect data about thousands of webpages without the necessity of visiting each webpage in an Internet browser manually. This functionality uses some webpages which are available on the Internet, contain links to other webpages and allow determination if the webpages on these links are likely to be phishing or legitimate. For the creation of models in Weka, data about webpages has to be selected from the database in the backend. Then, the data is exported to the CSV format and loaded into Weka. The models from Weka are processed by the backend into files which are stored in the backend. When the extensions are started, these files are downloaded by them if updates are detected.



Figure 1. A screenshot of the developed browser extension for monitoring visited webpages.

### III. COLLECTED WEBPAGE DATA

A group of websites collected with the browser extension for data collection introduced in Section II is described here. This group contains data about 21755 webpages described by 29 attributes whose values were saved in the backend and value legitimate or phishing was assigned to each website on the basis of an inspection. The attributes had been inspired by [1]. Suppose a set $W$ which represents all collected webpages is defined. In other words, $W$ is a set of known instances. The cardinality of $W$ is 21755. Suppose each webpage $w \in W$ is described by 29 attributes $A_k \in A = \{A_1; \dots; A_k; \dots; A_{29}\}$. Symbol $A_k(w) = a$ is employed for any numerical attribute $A_k$ whose value is $a$ for webpage $w$. Symbol $A_k = P$ for any numerical attribute $A_k$ where $P$ is a set of numbers represents possible numerical values of $A_k$. Symbol $A_k(w) = a_{k,l}$ is employed for any categorical attribute $A_k$ whose value is $a_{k,l}$ for webpage $w$. Possible categorical values $a_{k,1}, \dots, a_{k,l}, \dots a_{k,l_k}$ for attribute $A_k$ are represented by symbol $A_k = \{a_{k,1}; \dots; a_{k,l}; \dots; a_{k,l_k}\}$. The assigned value legitimate or phishing for each webpage $w \in W$ is represented by $D(w)$. Symbol $d_1$ is employed for a categorical value meaning legitimate and $d_2$ is for phishing. Symbol $D = \{d_1; d_2\}$ means $d_1, d_2$ are possible values for $D$. A description of the collected data is in Table I.

TABLE I.      DESCRIPTION OF COLLECTED DATA

| Particular Attribute | Type | Possible Values | Units |
|---|---|---|---|
| *BaseURLLengthWithParams* ($A_1$) | Numerical | 4, 5, 6, … | count |
| *AtSymbolInURL* ($A_2$) | Categorical | *absent* ($a_{2,1}$) | N/A |
| | | *present* ($a_{2,2}$) | |
| *DashInDomain* ($A_3$) | Categorical | *absent* ($a_{3,1}$) | N/A |
| | | *present* ($a_{3,2}$) | |
| *DoubleSlashInURL* ($A_4$) | Categorical | *absent* ($a_{4,1}$) | N/A |
| | | *present* ($a_{4,2}$) | |
| *DNSRecord* ($A_5$) | Categorical | *absent* ($a_{5,1}$) | N/A |
| | | *present* ($a_{5,2}$) | |
| *PercentageOfExternalPictures* ($A_6$) | Numerical | [0; 100] | % |
| *HostInURL* ($A_7$) | Categorical | *absent* ($a_{7,1}$) | N/A |
| | | *present* ($a_{7,2}$) | |
| *HTTPSInDomain* ($A_8$) | Categorical | *absent* ($a_{8,1}$) | N/A |
| | | *present* ($a_{8,2}$) | |
| *IFrame* ($A_9$) | Categorical | *absent* ($a_{9,1}$) | N/A |
| | | *present* ($a_{9,2}$) | |
| *NonStandardPort* ($A_{10}$) | Categorical | *absent* ($a_{10,1}$) | N/A |
| | | *present* ($a_{10,2}$) | |
| *TextFieldInPopUpWindow* ($A_{11}$) | Categorical | *absent* ($a_{11,1}$) | N/A |
| | | *present* ($a_{11,2}$) | |
| *ShorteningService* ($A_{12}$) | Categorical | *absent* ($a_{12,1}$) | N/A |
| | | *present* ($a_{12,2}$) | |
| *SubmittingToMail* ($A_{13}$) | Categorical | *absent* ($a_{13,1}$) | N/A |
| | | *present* ($a_{13,2}$) | |
| *HTTPSProtocol* ($A_{14}$) | Categorical | *absent* ($a_{14,1}$) | N/A |
| | | *present* ($a_{14,2}$) | |
| *ServerFormHandler* ($A_{15}$) | Categorical | *absent* ($a_{15,1}$) | N/A |
| | | *present* ($a_{15,2}$) | |
| *IndexationByGoogle* ($A_{16}$) | Categorical | *absent* ($a_{16,1}$) | N/A |
| | | *present* ($a_{16,2}$) | |
| *WebsiteRanking* ($A_{17}$) | Numerical | [0; 10] | Open Page Rank |
| *AllowOnMouseOver* ($A_{18}$) | Categorical | *absent* ($a_{18,1}$) | N/A |
| | | *present* ($a_{18,2}$) | |
| *AllowRightClick* ($A_{19}$) | Categorical | *absent* ($a_{19,1}$) | N/A |
| | | *present* ($a_{19,2}$) | |
| *IPAddress* ($A_{20}$) | Categorical | *absent* ($a_{20,1}$) | N/A |
| | | *present* ($a_{20,2}$) | |
| *NumberOfLinksPointingToPage* ($A_{21}$) | Numerical | 0, 1, 2, … | count |
| *NumberOfSubDomains* ($A_{22}$) | Numerical | 0, 1, 2, … | count |
| *NumberOfRedirections* ($A_{23}$) | Numerical | 0, 1, 2, … | count |
| *DomainAge* ($A_{24}$) | Numerical | 0, 1, 2, … | day |
| *PercentageOfExternalURLsInATags* ($A_{25}$) | Numerical | [0; 100] | % |
| *PercentageOfExternalOrNoURLsInATags* ($A_{26}$) | Numerical | [0; 100] | % |
| *PercentageOfExternalURLsInMetaLinkScriptTags* ($A_{27}$) | Numerical | [0; 100] | % |
| *TimeToDomainExpiration* ($A_{28}$) | Numerical | 0, 1, 2, … | day |
| *WebTraffic* ($A_{29}$) | Numerical | 0, 1, 2, … | Alexa Rank |
| *Class* ($D$) | Categorial | *legitimate* ($d_1$) | N/A |
| | | *phishing* ($d_2$) | |

Attribute *BaseURLLengthWithParams* ($A_1$) is related to the URL and represents the sum of the number of characters in the domain and the number of characters in the parameters after ?. $A_1 = \{4, 5, 6, \dots\}$. Attribute $A_2 = AtSymbolInURL = \{a_{2,1}; a_{2,2}\}$ = {*absent*; *present*} gives information if the URL address of a webpage $w$ contains the @ symbol. If it does not contain this

symbol, $A_2(w) = absent$. Otherwise, $A_2(w) = present$. The other attributes can be read from Table I in a similar way. In addition to the description in Table I, the collected data has been analyzed and the summary of this analysis is in Table II.

TABLE II. WEBPAGE DATA ANALYSIS

| Attribute | Possible Values | Frequency | Median | Mode |
|---|---|---|---|---|
| $A_1$ | 4, 5, 6, … | N/A | 22 | 21 |
| $A_2$ | absent ($a_{2,1}$) | 21555 | N/A | absent |
| | present ($a_{2,2}$) | 200 | | |
| $A_3$ | absent ($a_{3,1}$) | 18708 | N/A | absent |
| | present ($a_{3,2}$) | 3047 | | |
| $A_4$ | absent ($a_{4,1}$) | 21701 | N/A | absent |
| | present ($a_{4,2}$) | 54 | | |
| $A_5$ | absent ($a_{5,1}$) | 4393 | N/A | present |
| | present ($a_{5,2}$) | 17362 | | |
| $A_6$ | [0; 100] | N/A | 0 | 0 |
| $A_7$ | absent ($a_{7,1}$) | 21729 | N/A | absent |
| | present ($a_{7,2}$) | 26 | | |
| $A_8$ | absent ($a_{8,1}$) | 9 | N/A | present |
| | present ($a_{8,2}$) | 21746 | | |
| $A_9$ | absent ($a_{9,1}$) | 17894 | N/A | absent |
| | present ($a_{9,2}$) | 3861 | | |
| $A_{10}$ | absent ($a_{10,1}$) | 21752 | N/A | absent |
| | present ($a_{10,2}$) | 3 | | |
| $A_{11}$ | absent ($a_{11,1}$) | 20780 | N/A | absent |
| | present ($a_{11,2}$) | 975 | | |
| $A_{12}$ | absent ($a_{12,1}$) | 19532 | N/A | absent |
| | present ($a_{12,2}$) | 2223 | | |
| $A_{13}$ | absent ($a_{13,1}$) | 19004 | N/A | absent |
| | present ($a_{13,2}$) | 2751 | | |
| $A_{14}$ | absent ($a_{14,1}$) | 7491 | N/A | present |
| | present ($a_{14,2}$) | 14264 | | |
| $A_{15}$ | absent ($a_{15,1}$) | 17797 | N/A | absent |
| | present ($a_{15,2}$) | 3958 | | |
| $A_{16}$ | absent ($a_{16,1}$) | 8225 | N/A | present |
| | present ($a_{16,2}$) | 13530 | | |
| $A_{17}$ | [0; 10] | N/A | 0.76 | 0 |
| $A_{18}$ | absent ($a_{18,1}$) | 15774 | N/A | absent |
| | present ($a_{18,2}$) | 5981 | | |
| $A_{19}$ | absent ($a_{19,1}$) | 21582 | N/A | absent |
| | present ($a_{19,2}$) | 173 | | |
| $A_{20}$ | absent ($a_{20,1}$) | 21690 | N/A | absent |
| | present ($a_{20,2}$) | 65 | | |
| $A_{21}$ | 0, 1, 2, … | N/A | 0 | 0 |
| $A_{22}$ | 0, 1, 2, … | N/A | 0 | 0 |
| $A_{23}$ | 0, 1, 2, … | N/A | 0 | 0 |
| $A_{24}$ | 0, 1, 2, … | N/A | 0 | 0 |
| $A_{25}$ | [0; 100] | N/A | 54 | 0 |
| $A_{26}$ | [0; 100] | N/A | 58 | 100 |
| $A_{27}$ | [0; 100] | N/A | 92 | 100 |
| $A_{28}$ | 0, 1, 2, … | N/A | 0 | 0 |
| $A_{29}$ | 0, 1, 2, … | N/A | 0 | 0 |
| $D$ | legitimate ($d_1$) | 17662 | N/A | legitimate |
| | phishing ($d_2$) | 4093 | | |

## IV. EXPERIMENTAL EVALUTATION

The data mining models used in the browser extension for monitoring visited webpages from Section II and created with the collected webpage data from Section III are evaluated here. The models are created with Weka, which is an open-source data mining tool that is widely used for research and industrial applications [15]. An example of a data mining model produced by Weka is shown in Fig. 2. Normally, there are a model created on the basis of all data and some statistics related to its validation. Due to a low percentage of phishing webpages within the total number of webpages in the real world, 10-fold cross-validation is employed. This is similar to the medical field where 10-fold cross-validation is typical and where only a small percentage of patients within the whole universe has a particular medical issue. 10-fold cross-validation partitions the webpage data into ten subsets with equal numbers of cases with legitimate and phishing webpages. Nine subsets are used for the creation of a model and the model is then validated on the remaining subset, which is possible thanks to the availability of the values for $D$. This is repeated ten times for the purposes of using each subset for validation exactly once. The statistics from Weka is recomputed into sensitivity, specificity, accuracy and criterion. TP Rate for phishing in Weka is equal to sensitivity which expresses the ability to correctly identify phishing webpages. FP Rate for phishing in Weka is used for the computation of specificity as $1 - $ FP Rate. Specificity expresses the ability to correctly identify legitimate webpages. Ideally, sensitivity equals one and specificity equals one. When sensitivity is too low, the browser extension for monitoring visited webpages does not show any warnings for many dangerous phishing webpages. This would make the use of the extension meaningless. When specificity is too low, the extension shows many warnings even if the visited webpages are legitimate. This would make the use of the extension bothersome, especially when the high number of existing legitimate webpages is considered. Therefore, specificity should be very close to one. Both sensitivity and specificity are important and so the criterion defined as the average of sensitivity and specificity is used. Correctly classified instances in Weka correspond to accuracy, which is not very determinative in unbalanced situations.
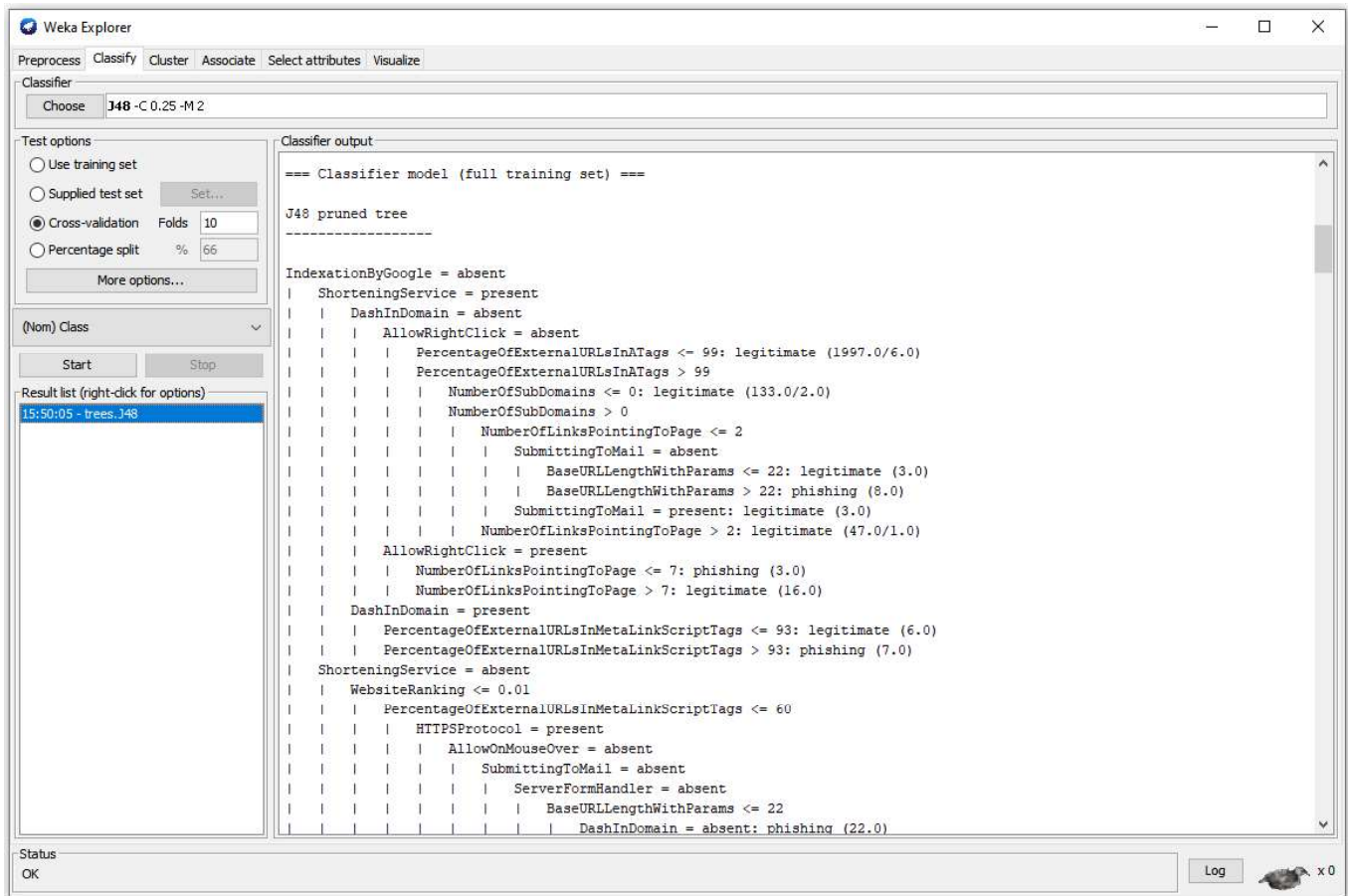
Figure 2.   A screenshot of the creation of a data mining model in Weka.

The results of the evaluation are summarized in Table III where the rows correspond to the data mining models supported in the browser extension for monitoring visited webpages and the columns are related to the particular recomputed measures. BayesNet is a Bayesian network model implemented in Weka as a Java class with the same name and displayed in Fig. 1 as Naive Bayes network. J48 is a decision tree model from Weka which is shown in Fig. 1 as j48 decision tree. MultilayerPerception is a neural network model from Weka whose name is MLP in Fig. 1. NNge is a nearest neighbor classifier. The best value of the criterion was achieved by NNge with 0.97350. Its specificity 0.99292 is also very close to one, which means that the use of the extension should not be bothersome. Very similar values were achieved by J48 as well. BayesNet has the worst values.

TABLE III.        RESULTS IN 10-FOLD CROSS-VALIDATION

| Model | Measure | | | |
|-------|---------|---------|---------|---------|
| | *Sensitivity* | *Specificity* | *Accuracy* | *Criterion* |
| BayesNet | 0.84266 | 0.95006 | 0.92986 | 0.89636 |
| J48 | 0.93599 | 0.99156 | 0.98111 | 0.96378 |
| MultilayerPerceptron | 0.88712 | 0.98800 | 0.96902 | 0.93757 |
| NNge | 0.95407 | 0.99292 | 0.98561 | 0.97350 |

## V.    CONCLUSIONS

Webpage data was collected for the development of data mining models classifying webpages into legitimate or phishing. The data has 21755 instances corresponding to webpages described by 29 attributes. In addition, two browser extensions and a backend were developed for the purposes of phishing webpages detection and data collection. Four data mining models such as a Bayesian network, decision tree, nearest neighbor classifier and neural network were created with the collected data and included for monitoring visited webpages. 10-fold cross-validation showed promising results with the best combination of sensitivity and specificity equaling to 0.95407 and 0.99292. Future work may include further development of models for an even more accurate classification of webpages.

REFERENCES

[1] N. Abdelhamid, A. Ayesh, F. Thabtah, "Phishing detection based associative classification data mining," Expert Systems with Applications, vol. 41, no. 13, pp. 5948-5959, 2014.

[2] M. Alwateer, S. W. Loke, N. Fernando, "Enabling drone services: Drone crowdsourcing and drone scripting," IEEE Access, vol. 7, art. no. 110035, 2019.

[3] J. Bohacik, I. Skula, M. Zabovsky, "Data mining-based phishing detection," in Federated Conference on Computer Science and Information Systems, 2020, pp. 27-30.

[4] D. Cvitanic, "Drone applications in transportation," in International Conference on Smart and Sustainable Technologies, 2020, art. no. 20133277.

[5] C. L. Evans, "Clone Zone is an easy tool for building fake websites," available at https://www.vice.com/en/article/3dkyyb/clone-zone-is-an-easy-tool-for-building-fake-websites, 2015.

[6] N. N. Gana; S. M. Abdulhamid, "Machine learning classification algorithms for phishing detection: A comparative appraisal and analysis," in International Conference of the IEEE Nigeria Computer Chapter,

[7] A. Y. Husodo, G. Jati, N. Alfiany, W. Jatmiko, "Intruder drone localization based on 2D image and area expansion principle for supporting military defence system," in IEEE International Conference on Communication, Networks and Satellite, 2019, pp. 35-40.

[8] S. Mahmoud, N. Mohamed, J. Al-Jaroodi, "Integrating UAVs into the cloud using the concept of the Web of Things," Journal of Robotics, vol. 2015, art. no. 631420.

[9] S. Paliath, M. A. Qbeitah, M. Aldwairi, "PhishOut: Effective phishing detection using selected features," in International Conference on Telecommunications, 2020, art. no. 20135977.

[10] S. Shukla, P. Sharma, "Detection of phishing URL using Bayesian optimized SVM classifier," in International Conference on Electronics, Communication and Aerospace Technology, 2020, pp. 1385-1389.

[11] A. J. Van der Merwe, M. Loock, M. Dabrowski, M, "Characteristics and responsibilities involved in a phishing attack," in Winter International Symposium on Information and Communication Technologies, 2005, pp. 249-254.

[12] S. Nair, G. Rodrigues, C. Dsouza, S. Bellary, V. Gonsalves, "Designing of beach rescue drone using GPS And Zigbee technologies," in International Conference on Communication and Electronics Systems, 2019, pp. 1154-1158.

[13] S. Patil, S. Dhage, "A methodical overview on phishing detection along with an organized way to construct an anti-phishing framework," in International Conference on Advanced Computing & Communication Systems, 2019, pp. 588-593

[14] R. Thomazella, J. E. Castanho, F.R.L. Dotto, O.P. Rodrigues Junior; G. H. Rosa; A. N. Marana; J. P. Papa, "Environmental monitoring using drone images and convolutional neural networks," in IEEE International Geoscience and Remote Sensing Symposium, 2018, pp. 8941-8944.

[15] I. H. Witten, E. Frank, M. A. Hall, C. J. Pal, Data Mining: Practical Machine Learning Tools and Techniques (4th edition). USA: Morgan Kaufman, 2016.