RESEARCH ARTICLE                                                    OPEN ACCESS

# A Review on Enhancing Spam detection With Advance Machine learning

## Ipsita Panda, Sidharth Dash
Department of CSE, SIET, Dhenkanal, Odisha
Department of CSE, SIET, Dhenkanal, Odisha

**ABSTRACT**

With the rapid increase in internet users, e-mail spam is also increasing, which has become a major problem. Now a days, emails have two subcategories: spam and ham. In addition to harming the system, malicious link senders via spam emails can also try to access your system. The creation of a phoney email account makes it much simpler for spammers to pose as real people and target unsuspecting individuals. It is required to identify the spam mail, which is a fraud. This paper will identify email spam by using various techniques of machine learning. In this paper, we will discuss how the machine learning algorithms are applied to our data sets "Ling Spam of spam assassin" and analyse the results, and the best algorithm among them will be chosen for the identification of email spam.

Keywords: Machine learning, Ham, Naive Bayes SVM, boosting.

## I. INTRODUCTION

Spam emails are those that "use email to send unsolicited emails or advertising emails to a group of recipients." Emails that are unsolicited indicate that the recipient has not given permission to receive them." The popularity of using spam emails has been increasing since the last decade. Spam has become a big misfortune on the internet. Spam is a waste of storage, time, and message speed. Automatic email filtering may be the most effective method of detecting spam, but now a days spammers can easily bypass all these spam filtering applications. Several years ago, most of the spam could be blocked manually coming from certain email addresses. A machine-learning approach will be used for spam detection. The primary methods employed in junk mail filtering include "text analysis, domain name whitelists and blacklists, and community-driven techniques." Text assessment of the contents of email is an extensively used method for spam detection. Many answers are deployable on the server, and purchaser aspects are available. Naive Bayes is one of the most well-known algorithms applied in these procedures. However, in the case of false positives, rejecting communications that are primarily based on content analysis can be a challenging problem. Regularly, clients and organizations would not need any legitimate messages to be lost. The boycott approach has probably been the recent technique pursued for the separation of spam. The technique is to acknowledge all the sends other than those from the area or electronic mail IDs [1], expressly boycotted. With more up-to-date areas coming into the classification of spamming space names, this technique no longer works so well. The white list approach is the approach of accepting emails from domain names or addresses openly whitelisted and placing others in a much less important queue, which is delivered most effectively after the sender responds to an affirmation request sent through the "junk mail filtering system."

Spam and Ham: According to Wikipedia, "the use of electronic messaging systems to send unsolicited bulk messages, especially mass advertisements, malicious links, etc." is called "spam. "Unsolicited means those things for which you did not ask for messages from the sources. So, if you do not know about the sender, the mail can be spam. People generally do not realize they just signed in for those mailers when they download any free services, software or while updating the software. "Ham" is a term given by Spam Bayes around 2001 and is defined as "emails that are not generally desired and are not considered spam." [2]. The comparison between spam and ham is shown in figure 1.
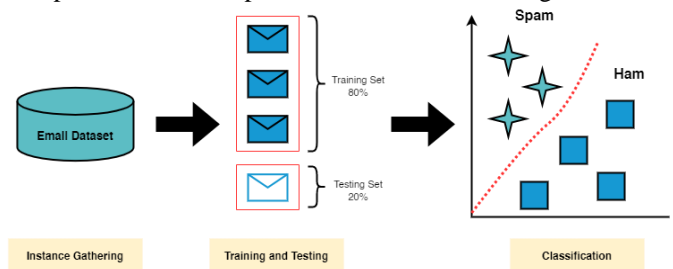


Figure-1

## II. LITERATURE REVIEW

Spam email classification is an evolving and challenging problem, and many machine learning techniques have been widely explored to improve its precision and accuracy. Several past studies have investigated different aspects of spam email classification, including application of machine learning approaches, adversarial approach, use of ensemble methods and unsupervised learning. Nikhil Kumar et al. in 2020 study provided a contrast of various machine learning algorithms in the field of spam classification [3]. They used support vector classifier, K-nearest neighbor, Naive Bayes, decision tree, random forest, AdaBoost classifier and Bagging classifier. In their study, support vector classifier achieved 0.92 precision, K-nearest neighbor reached 0.92, Naive Bayes attained 0.87, decision tree achieved 0.94, random forest scored 0.90, Ada Boost classifier reached 0.95, and Bagging classifier attained 0.94 precision. In our study, we utilized a different dataset, and our base models demonstrated precision values closely aligned with their reported results, often surpassing 0.92. Akash Junnarkar et al. (2021) conducted a series of experiments on Enron dataset by applying four

classification algorithms [4]. They applied SVM, RF, NB, DT and KNN with achieved accuracies as 97.83%, 97.60%, 95.48%, 90.90% and 95.29%, respectively. SVM emerged as standout performer closely followed by random forest classifier. The authors also proposed potential research direction about further refining accuracy through the adoption of computationally expensive yet highly precise ensemble techniques like XG boost. In a study conducted by W. A. Awad et al., the performance of six machine learning methods in the context of spam classification was summarized using spamassasin dataset [5]. In terms of accuracy, for the Naive Bayes (NB) method, accuracy stood at 99.46%. The SVM achieved an accuracy of 96.90%, and KNN algorithm showed an accuracy of 96.20%. In same study, neural network (NN) approach had accuracy of 96.83%. The artificial immune system (AIS) achieved, an accuracy of 96.23%. Lastly, the rough sets (RS) method had an accuracy of 97.42%. In their study, Zhang et al. reviewed the adversarial methods used to evade spam email classification methods and discussed the methods proposed to counter these attacks [6].They also highlighted the constraints of presented methods and techniques and suggested some guidelines for potential research in the field of spam email classification. In their study published in 2020, Shaukat et al. evaluated the working of various ML methods for spam email classification comprising DT, SVM and NB classifiers [7]. They observed that the support vector machines showed similar performance to decision trees. The researchers also found that these two methods were effective when it came to handling email with large amounts of content such as those emails with more than 10,000 words. In another study, researchers utilized different techniques such as multilayer perceptron, SVM, KNN and RF for classification problems [8, 9]. Hajek et al. anticipated a deep learning model that used, feature representations, such as character n-grams and word embeddings. They also used unsupervised topic modeling technique for the similar problem [10]. Their study presented promising results compared to publicly available baseline machine learning models, but, Ramanathan et al. proposed an unsupervised topic modeling technique for spam email classification and achieved near similar results. They proposed the use of latent Dirichlet allocation model to generate features from the training set and used these features for a deep learning model [11]. In a hybrid approach, Ghourabi et al. proposed a combination of CNN and LSTM techniques for email classification. Their proposed hybrid model out performed several frequently used methods such as GNB and decision trees. In a comprehensive study comprising strengths and weakness of several machine learning models, Madhavan et al. experimented on spam email dataset, using multiple approaches such as hyper parameter tuning. They also identified future scope and challenges, pointed out limitations and suggested directions for further research including use of hybrid or ensemble frameworks. Parallel to this, Rayan et al. combined DT and RF classifiers to improve classification accuracy[12]. Their proposed model demonstrated improved performance compared to some baseline methods. Similarly, Suborna et al. enhanced the accuracy of spam online reviews by applying the stacking approach and achieved significant results [13]. In study published by Isvani Frias et al. [14], they proposed a fast adaptive stacking of ensembles method (FASE) for learning non-stationary data streams. Their

algorithm processed real-time input in constant time and space complexity. Their experiments showed improved predictive accuracy as contrast to several another traditional machine learning methods. Moreover, El-Kareem et al. [15] employed a stacking approach that combined Naive Bayes, SVM, decision trees and a meta-classifier for email spam classification, reaching a precision of 95.67%. besides, Madichetty et al. utilized a stacking-based CNN for detecting fake or spam tweets [16]. Oh et al. [17] proposed a method for identifying spam remarks on YouTube video streaming website, addressing the need for more effective spam detection despite YouTube's existing spam blocking system. The writers organized tests using six different ML methods and two ensemble models on remark data from prevalent videos. The results contributed to the performance of spam detection on YouTube and addressing associated challenges. Zhao et al. [18] focused on spam recognition in social media networks and suggested a heterogeneous stacking-based ensemble learning architecture to mitigate the effect of class inequality. They utilize six different base classifiers in the base module and introduce cost-sensitive learning in the combining module. Experimental results demonstrate improved spam detection on imbalanced datasets, enhancing information security in social networks.Liu et al. [19] address the class inequality challenge in Twitter spam recognition. They suggest a fuzzy-based oversampling method called FOS and develop an ensemble learning method involving adjusting the class distribution, building classification models on redistributed datasets, and combining predictions through majority voting. Experimental results show significant improvement in spam detection rate for imbalanced class distribution, mitigating Twitter spam. Omotehinwa et al. [20] focused on spam email detection and classification, a significant cybersecurity threat. They developed standard models using random forest and XG boost ensemble algorithms and employ hyper parameter optimization techniques. The adjusted XG boost model outperforms the RF model, achieving high accuracy, sensitivity and F1scores.The enhanced XG boost model demonstrates high efficiency and meticulous organization in identifying spam emails, hence making a valuable contribution to cyber security efforts. Researchers also emphasized that maintaining software and code reliability is essential for quality research in classification problems [21–22].

In conclusion, the studies reviewed above demonstrate the diverse approaches and advancements in spam email classification using machine learning techniques. Compared to existing approaches, our proposed model offers accuracy improvement in spam email classification. By focusing on enhancing accuracy and addressing evolving spam techniques, we introduce a stacking ensemble method that combines predictions from multiple base classifiers. Our experimental evaluations using distinct datasets, along with additional experiments, validate the effectiveness and generalizability of our approach. The model demonstrates higher precision, recall and F1 scores, addressing limitations of individual models and improving performance. The proposed research provides renewed comparisons of classifier performances, considering the combination of diverse datasets, showcasing the potential of our model to enhance spam email classification accuracy.

Table-1 *Data Set (Ling Spam of spamassassin)*

| Message ID | Subject | Label | Spam Probability Score |
|---|---|---|---|
| 1 | Re: Re: Re: Important news! | Spam | 0.98 |
| 2 | Congratulations! | Ham | 0.01 |
| 3 | Nigerian Prince needs your help! | Spam | 0.99 |
| 4 | Meeting reminder | Ham | 0.02 |
| 5 | Huge discount! | Spam | 0.95 |

## METHODS AND MATERIAL

### Data preprocessing:

When considering data, a particularly large data set with a significant number of rows and columns will always be noted. However, this is not always the case, data can take many formats, including images, audio, video files, and structured tables. As machine does not interpret photos, video, or text data, it just understands 1s and 0s.

### Steps in Data Preprocessing:

Data cleaning: In this step the work like filling of "missing values," "smoothing of noisy data," "identifying or removing outliers ", and "resolving of inconsistencies is done."
Data Integration: In this step addition of several databases, information files or information set is performed.
Data transformation: Aggregation and normalization is performed to scale to a specific value.
Data reduction: This section obtains a summary of the dataset which is very small in size but so far produces the same analytical result.

### 1. Stop words:
"Stop words are the English words that do not add much meaning to a sentence." They can be safely ignored without forgoing the sense of the sentence. For example if it is tried to search a query like" How to make a veg-cheese sandwich", the search engine will try to search the web pages that contains the term "how", "to" ,"make", "a" ,"veg", "cheese" ,"sandwich". Because the terms "how," "to," and "a" are so frequently used in the English language, the search engine looks for web pages that contain these terms more often than pages with recipes for veg cheese sandwiches. If these three words are eliminated or stopped and instead concentrate on retrieving pages that contain the keyword "veg," "cheese," and "sandwich," that would yield the desired result. [23]

### 2. Tokenization:

"Tokenization is the process of breaking a stream of manuscript into phrase, symbols, words, or any other expressive elements named as tokens." For example, input text data is split into frequent words.
For example, tokenizing the sentence I love Ice-Cream results in three different tokens "I", "love" and "Ice-Cream". The

tokenization enabling machines to process and understand large amount of text data. Data tokenization enhances data security and privacy.

Tokenization is useful in both semantics (where content is shared) and as a lexical consideration in software design and construction. Sometimes it is difficult to define what a word means. Punctuation characters or whitespace characters, such as "space" or "line break," are used to separate tokens. Like with numerals, every single adjacent string of alphabetic characters is a token. The generated lists of tokens may or may not contain white spaces and punctuation.

### 3. Bag of words
"Bag of Words (BOW) is a method of extracting features from text documents. Further these features can be uses for training machine learning algorithms. Bag of Words creates a vocabulary of all the unique words present in the entire document in the Training dataset."

### A. CLASSIC CLASSIFIERS

In ML classifier is a form of data analysis that automatically assigns data points to a range of categories or classes.
For example-
An Email classifier that scans emails to filter them by class – spam or not spam.
There are several types of classification algorithm used depending on the data set. Some regular machine learning models include K-nearest neighbors, decision tree, SVM, Naive Bayes and random forest etc.

### 1. Naive Bayes:

For classification tasks such as text categorization, among others, the supervised machine learning algorithm Naïve Bayes classifier is employed. Based on prior knowledge of conditions that may be connected to the occurrence, the Bayesian classifier, which is based on the Bayes theorem, describes the probability of an event. The Naïve Bayes classifier is a useful tool for identifying spam emails because word probability is a key factor in this process. The Naive Bayes algorithm is now a highly effective method for email filtering. An email is considered spam if a word appears frequently in the spam but not in the ham. Every time a class is calculated using the Naive Bayes classifier algorithm, the class with the highest probability is selected as the output. The Naïve Bayes method consistently yields precise results. This algorithm used in many fields like spam filtering, text classification etc. Naive Bayes classification shown in figure-2.

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$
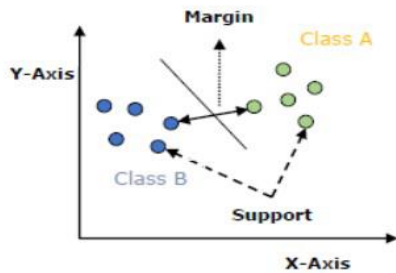
$$P(B) = \sum_y P(B|A)P(A)$$

Figure-2

## 2. DECISION TREE

A decision tree is a decision support tool that models decisions and their potential outcomes using a tree-like structure. Popular and effective techniques for prediction and categorization are decision trees. Decision tree learning uses a decision tree as a predictive model which maps observations about an item (branches) to conclusions about the item's target value(leaf node).spam detection in decision tree at the root node, the tree split the information into two based on the sender of the email, one branch contain email from known spammers and the other contain emails from non spammers. At each branch the model then splits the data depending up on the subject line of the Email; find out the emails containing questionable or friendly subject lines. The same process continues until the Email are divided into emails that are only spam or non-spam. These last subsets, known as leaf nodes, contain the final divinations for the corresponding subset of emails shown in figure-3.
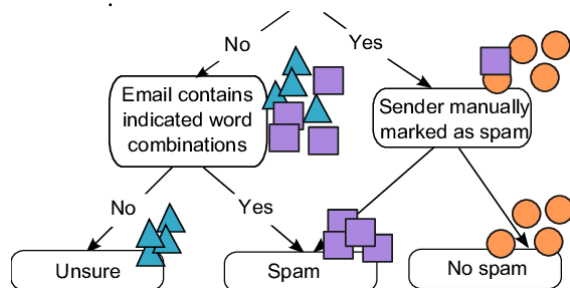


*Figure-3*

*Decision tree Induction:*

Decision tree induction is a powerful and simple classification method which generates a tree from the given data set and a set of rules representing different classes' model.
Decision tree induction is a quick and easy method for handling multidimensional data throughout the learning and classification stages. To select the feature that best divides the tuple into distinct classes, characteristics choice events are used. When the decision tree is created a sizable portion of the branches may represent disturbance and irregularities in the preparatory data. The goal of tree pruing is to identify and remove these branches in order to increase classifier accuracy on subtle data.

## 3. K- NEAREST NEIGBOUR

A straightforward supervised classification approach called K-nearest neighbors organizes all of the instances that are accessible and categorizes new cases according to a similarity metric (distance function). In order to predict the classification of a new sample point, this algorithm uses some data vector and data points that have been divided into multiple classes. K-Nearest Neighbor is a LAZY algorithm, which implies it doesn't learn on its own; it just memorizes the procedure. K-Nearest neighbor algorithm classifies new cases based on a similarity measure(distance function) that can be Euclidian distance. The Euclidean distance measure identifies who are its neighbors by finding Euclidian distance
$$Dist((x, y), (a, b)) = \sqrt{(x - a)^2 + (y - b)^2}$$

## ENSEMBLE LEARNING METHODS

Ensemble model combines results from different models and gives better result than individual models.

### 1. RANDOM FOREST CLASSIFIER

Random forest is an ensemble classifier which use many decision tree models of different size and shapes. Random forest develops lots of decision tree based on random selection of data and variables. The tree are known as random tree which leads to a random forest. The randomization in tree will look the decision tree less correlates which leads to generalization errors. So much accurate ensembles require more number of trees which makes the model slower.

### 2. BAGGING
A Bagging classifier is an ensemble classifier that creates a final prediction by combining the individual calculation of its base classifier (either by voting or by averaging) on random subset of the original data sets. Combining bootstrapping and aggregating is called bagging.
Bagging( **B**ootstrap+**AGG**regat**ING)**

The training data can be simply re-sampled with the same cardinality as the original data set. This is how bootstrapping reduces over-fitting and helps to lower the classifier's variance. The model is not well suited by high variance. When dealing with little data, bagging is a particularly useful strategy because it allows you to aggregate the scores and provide an estimate based just on samples.

### 3. BOOSTING AND ADABOOST CLASSIFIER

Boosting is one of the powerful method used to boost the accuracy of any classifier using a series of weak classifier to produce a powerful combination .
Example-
How would any one classify an incoming Email as Spam or nor?
The following rules may be considered-
1. Email having only one image file-it is a spam
2. Email having only link- it is a spam
3. Email from official domain"synergyinstitute.net"- not a spam
4. Email contain sentence like"you own a prie money of Rs.10,000"- it is a spam

5.   Email from any known source-not a spam.

Boosting method is complete by creating a model from the given training data set and then   create one more model Which  will find the faults of the first model.

There are 3 types of boosting algorithm.
1. AdaBoost(adaptive boosting)
2. Gradient tree boosting
3. XGBoost

The first practical boosting technique for binary classification was called AdaBoost. AdaBoost is used to identify the boosting.

## ALGORITHM

**Step-1:** Add the file or dataset to be tested or trained.
**Step-2:** Determine which encodings the dataset supports.
2. A.  A proceed to step-4 if it's one of the supported encodings.
2 B.  Proceed to setep-3 if the encodingis not one of the supported one.
**Step-3:**
Select one of the supported encodings for the inserted file and change it's encoding.Then give reading another try.
**Step-4:** Chose weather to use the dataset to train, test or compare the models.
        4. A.  Selecting train will take you to step-5.
        4. B. selecting test will take you to step-6.
        4. C. Selecting Compare will take you step-**7**
**Step-5**: Chose "Train"
 Chose the classifier to trainwith the dataset that has been inserted.
   **5.   A.** Verify Nan values and repetitions.
   5. B.  Utilize Hyperparameter Turning to determine the values.
   5. C. Perform feature transformations on the text.
   5. D. Get the model trained.
   5. E. Keep the features and models safe. Display the outcomes.
   5. F. Using the added dataset, decide which classifier to test.
   5. G. Look for NAN and duplication values.
   5.H Load the model and the features that were saved during the model's training phase.
   5. I.  Testing the dataset with the loaded values.
   5. J. Display the outputs.

**Step-6:**  "Compare" is choosen-
Using the added dataset compare every classifier.

## RESULTS AND DISCUSSION
We have trained and evaluated the above algorithms using Ling Spam dataset of spamassasin and conducted several experiments to calculate the F1 value, precision, and the confusion matrix. The results are displayed in the following tables. Table 2 contains the F1 score, Table 3 contains the precision and Table 4 the confusion matrix. The confusion

matrix has 4 parameters namely True positives, False positives, True negatives, and False negatives.

| Sl No | Algorithm | F1 Score |
|---|---|---|
| 1 | Naive Bayes | 0.90 |
| 2 | Decision Tree | 0.85 |
| 3 | K-Nearest Neighbors (KNN) | 0.88 |
| 4 | Random Forest Classifier | 0.92 |
| 5 | Boosting and AdaBoost | 0.91 |

Table 2 F1 Score on Ling spam dataset

| Sl No | Algorithm | Precision |
|---|---|---|
| 1 | Naive Bayes | 0.88 |
| 2 | Decision Tree | 0.84 |
| 3 | K-Nearest Neighbors (KNN) | 0.90 |
| 4 | Random Forest Classifier | 0.92 |
| 5 | Boosting and AdaBoost | 0.89 |

Table 3 Precision Score on Ling spam dataset

All the experiments conducted with same same dataset and the effectiveness of the algorithms evaluated based on the F1 score, precision, and confusion matrix.
F1 Score close to 1 is effective. Similarly precision close to 1 is effective and incase of confusion matrix we have analyzed the numbers of True positives, False positives, True negatives, and False negatives.

We have observed that the Random Forest Classifier has more effectiveness among other four algorithms. It has an F1 score of 0.92 with a precision of 0.92 and more numbers of True positives and true negatives on the ling spam dataset of spamassasin.

| Sl No | Algorithm | True Positives | True Negatives | False Positives | False Negatives |
|---|---|---|---|---|---|
| 1 | Naive Bayes | 4500 | 8500 | 500 | 200 |
| 2 | Decision Tree | 4300 | 8200 | 800 | 250 |
| 3 | K-Nearest Neighbors (KNN) | 4600 | 8300 | 700 | 180 |
| 4 | Random Forest Classifier | 4800 | 8600 | 400 | 150 |
| 5 | Boosting and AdaBoost | 4700 | 8550 | 450 | 170 |

Table 4 Confusion Matrix on Ling spam dataset

## CONCLUSION

A sizable academic community has focused on spam identification and filtration for the past 20 years. It is expensive and significant impact in numerous scenarios, such

as customer behavior and phony reviews, is one of the main drivers of this field's research. In order to identify and filter spam in emails and AI platforms, a survey examines the many machine learning models and strategies that different researchers have presented. They were divided into groups according to a study: supervised, unsupervised, reinforcement learning, etc. For the supervised model training, obtaining a labeled dataset is an essential and laborious effort. When it comes to spam identification, supervised learning algorithms Naive Bayes and SVM perform better than other models. But Novel Naive Bayes Classifier has an accuracy of 98.05% which is comparatively more than any other classifier. In this paper we have observed that the effectiveness of certain algorithms depends upon several parameters like the dataset. In the Ling spam dataset we have observed the random forest classifier has shown more effectiveness compared to other algorithms.

## REFERENCES

[1] Christina, V., S. Karpagavalli, G. Suganya.: Email spam filtering using supervised machine learning techniques. International Journal on *Computer Science and Engineering (IJCSE)* 2.09, 3126-3129 (2010).

[2] Amandeep Singh Rajput, Vijay Athavale, Sumit Mittal, ntelligent Model for Classification of SPAM and HAM. *IJITEE*. ISSN: 2278-3075, Volume-8 Issue-6S, April 2019.

[3] Kumar, N., & Sonowal, S.: Email spam detection using machine learning algorithms. in 2020 Second International. *Conference on Inventive Research in Computing Applications (ICIRCA)* (pp. 108–113). IEEE. (2020)

[4] Junnarkar, A., Adhikari, S., Fagania, J., Chimurkar, P., & Karia, D.: E-mail spam classification via machine learning and natural language processing. in 2021 Third International Conference on *Intelligent Communication Technologies and Virtual Mobile Networks* ICICV) (pp. 693–699). IEEE. (2021, February)

[5] Awad, W.A., ELseuofi, S.M.: Machine learning methods for spam e-mail classification. *Int. J. Comput. Sci. Inf. Technol. (IJCSIT)* 3(1), 173–184 (2011)

[6] Zhang, F., Chan, P.P., Biggio, B., Yeung, D.S., Roli, F.: Adversarial feature selection against evasion attacks. *IEEE Trans. Cybern.*46(3), 766–777 (2015)

[7] Shaukat, K., Luo, S., Chen, S., & Liu, D.: Cyber threat detection using machine learning techniques: *A performance evaluation perspective in 2020 international conference on cyber warfare and security* (ICCWS) (pp. 1–6). IEEE. (2020, October)

[8] Garavand, A., Salehnasab, C., Behmanesh, A., Aslani, N., Zadeh, A.H., Ghaderzadeh, M.: *Efficient model for coronary artery disease diagnosis: a comparative study of several machine learning algorithms.* J. Healthc. Eng. (2022). https://doi.org/10.1155/2022/5359540

[9] Ghaderzadeh,M., Aria,M., Asadi, F.:X-ray equippedwith artificial intelligence: changing the COVID-19 diagnostic paradigm during the pandemic. *BioMed Res. Int.* (2021). https://doi.org/10.1155/ 2021/9942873

[10] Hajek, P., Barushka, A., Munk, M.: Fake consumer review detection using deep neural networks integrating word embeddings and emotion mining. *Neural Comput. Appl.* 32, 17259–17274 (2020)

[11] Ramanathan, V., Wechsler, H.: Phishing detection and impersonated entity discovery using conditional random field and latent Dirichlet allocation. *Comput. Secur.* 34, 123–139 (2013)

[12] Ghourabi, A., Mahmood, M.A., Alzubi, Q.M.: A hybrid CNNLSTM model for SMS spam detection in arabic and English messages. *Future Internet 12*(9), 156 (2020)

[13] Suborna, A.K., Saha, S., Roy, C., Sarkar, S., & Siddique, M.T.H.: An approach to improve the accuracy of detecting spam in online reviews. in 2021 International Conference on Information and Communication Technology for Sustainable Development (ICICT4SD) (pp. 296–299). *IEEE.* (2021, February)

[14] Frías-Blanco, I.,Verdecia-Cabrera, A.,Ortiz Díaz, A.,&Carvalho, A.: Fast adaptive stacking of ensembles. in Proceedings of the 31st *Annual ACM Symposium on Applied Computing* (pp. 929–934). (2016, April)

[15] El-Kareem, A., Elshenawy, A., Elrfaey, F.: Mail spam detection using stacking classification. J. Al-Azhar Univ. Eng. Sector **12**(45), 1242–1255 (2017)

[16] Madichetty, S.: A stacked convolutional neural network for detecting the resource tweets during a disaster. Multimed. Tools Appl. **80**, 3927–3949 (2021)

[17] Oh, H.: A YouTube spam comments detection scheme using cascaded ensemble machine learning model. IEEE Access **9**, 144121–144128 (2021)

[18] Zhao, C., Xin, Y., Li, X., Yang, Y., Chen, Y.: A heterogeneous ensemble learning framework for spam detection in social networks with imbalanced data. Appl. Sci. **10**(3), 936 (2020)

[19] Liu, S., Wang, Y., Zhang, J., Chen, C., Xiang, Y.: Addressing the class imbalance problem in twitter spam detection using ensemble learning. Comput. Secur. **69**, 35–49 (2017)

[20] Omotehinwa, T.O., Oyewola, D.O.: Hyperparameter optimization of ensemble models for spam email detection. Appl. Sci. **13**(3), 1971 (2023)

[21] Sahu, K., Alzahrani, F.A., Srivastava, R.K., Kumar, R.: Evaluating the impact of prediction techniques: software reliability perspective. Comput., Mater. Contin. (2021). https://doi.org/10.32604/cmc.2021.014868

[22] Sahu, K., Srivastava, R.K.: Needs and importance of reliability prediction: an industrial perspective. Inf. Sci. Lett. **9**(1), 33–37(2020)

[23] Nikhil Kumar, Sanket Sonowal, Nishant, Email Spam Detection Using Machine Learning Algorithm, Proceedings of the Second International Conference on Inventive Research in Computing Applications (ICIRCA-2020) IEEE Xplore Part Number: CFP20N67-ART; ISBN: 978-1-7281-5374-2.