

A Spam Transformer Model for SMS Spam Detection

¹ Kanipakam Bhanu Moorthy, ² P.Mohan, ³ K.Vishnu Vardan Varma, ⁴ Pradeep Burri,
CSE Department,

^{1,2,3,4} Assistant Professor, Dhruva Engineering Collage, Hyderabad.
Shree Engineering Collage, Hyderabad.

Abstract: Spam emails, also referred to as non-self, are commercial or harmful unsolicited emails, sent to attack either a particular entity or an organization or a community of individuals. In addition to marketing, these It which contain ties to websites hosting phishing or malware set up to steal sensitive details. In this post, a review on the feasibility of using an anomaly anomaly negative selection algorithm (NSA) It introduces the detector applied to spam filtering. The high efficiency and low false detection of the NSA is Pace. Via three detection stages, the built system intelligently works to eventually decide Legitimacy of an email depending on the information collected in the training process. The unit works by Elimination is analogous to the functionality of T-cells in biological processes by negative selection. It It has been found that efficiency tends to increase with the addition of more datasets, this culminated in a 6% improvement in the identification rate of True Positive and True Negative thus maintaining an actual detection rate. 98.5% spam and ham identification score. The model has been correlated further with related models Studies and the outcome suggest that the proposed method results in an improvement of 2% to 15% in the right system. Spam and ham identification score.

Keywords: spam; ham; phishing; identification of anomalies; Negative range.

1. Introduction

For quite some time now, email has become an increasingly valuable contact tool, allowing virtually immediate exposure to every part of the globe through internet connections. Nearly 5 billion email accounts were actively in use in 2017, as reported by Tschabitscher [1], and this is projected to rise to over 5.5 billion by the end of 2019. The possibility that more than 270 billion emails are exchanged every day is also illustrated by Tschabitscher [1], but roughly 57% of these are only spam emails [1]. There are a range of current methods of machine learning and strategies that strongly mimic the filtering of spam or phishing emails by biological immune systems, but their efficiency has become a major concern. Most of the methods manage to efficiently prevent spam, but trade

They often restrict some of the emails that are not spam, classified as ham. This is a concern, since it may result in the consumer missing valuable details.

1.1. Common Threats

Various forms of email threats, such as email spoofing, phishing, and phishing variations, such as spear phishing, duplicate phishing, whaling, hidden redirect, etc., are routinely bombarded by users worldwide. Email spoofing also includes forging the email header (The From section) such that a real

person appears to have received the post. Email spoofing is a ploy used in spam campaigns and when they feel it is sent from someone they know; people prefer to open an email [2]. Email phishing is a type of spoofing that deceives the recipient with genuine messages [3].

To bypass anti-spam systems, malicious attackers have even tried to conceal the text behind images. It is a form of obfuscation whereby the message text is processed as a JPEG or GIF image and presented in the email. Which avoids the identification and blocking of spam messages by text-based spam filters?

1.2. An Email Architecture

The headers and the body of the email are made of emails. Next, the TCP/IP Header includes the IP address of the source and destination, then the SMTP envelope, containing the email transaction areas, and the email addresses of the source and destination (but this section is not accessible to the email clients); and then the SMTP headers, where the email addresses are stored.

Erent identifiable email sections such as 'topic', 'from' and 'to' fields exist (accessed by the email clients to connect the details to the user). This is the part where the fraudsters tinker with, since the real source and destination emails are stored in the SMTP envelope, which is not directly accessible to the recipient, as described above. Finally, the body of the document is the email message, which optionally includes connections and attachments that may be harmful in nature. Not only is bulk email annoying for regular email consumers, it often generates a big computer protection challenge that causes efficiency losses of billions of dollars [4]. Moreover, it is still the main platform for phishing [5, 6] and distributing harmful malware, such as viruses and worms [4].

1.3. Complexities resulting from Junk Emails

Spam emails may have several problematic effects on people, organizations, and the community in general, as mentioned above. Leung and Liang [7] found that phishing warnings sometimes result in a significant negative stock return. The negative effects on companies whose legal email communications are deemed spam by anti-spam systems [8] have been identified by other researchers. Through downloading malicious attachments, spam emails may inflict serious reputational harm, as well as identity theft of

an individual; stolen information may later be used to blackmail the victims [9]. Botnets [10] can also be circulated via spam emails. High profile incidents of abused firms have been daily occurrences. Due to a phishing email scheme, confidential financial details of workers of a United States Bus Corporation dropped into the possession of scammers in 2018[11]. Then, a recent whaling attack cost over USD 21 million [12] to the French cinema chain 'Path', also in 2018. These are only a few instances of spam emails becoming a common concern.

1.4. Shortcomings of Non-Automated Spam Filtration Methods

There are a range of non-automated spam recognition mechanisms available that do not focus on concepts of machine learning, however in combating contemporary spam assault dynamics and dynamism, these mechanisms face major bottlenecks. We may quickly illustrate some of the flaws in these schemes in this portion. Over the years, the blacklisting of sender addresses has become a common alternative. But this system alone (primary lone protection against spam emails) has proved ineffective; as the spammers are able to modify the sending address and the process of upgrading the database is always slow [13].

Another common method is the heuristic strategy, where a series of rules is added to incoming emails to label them as spam or ham. To build the rule collection, regular expressions are also used. However, if the scammers are able to access the rule set, they will plan their communications quite effectively beforehand to bypass the filtering mechanism. Keyword matching, country-based filtering, and relisting are other established mechanisms, to name a couple, but these are all sure from restrictions that modern spam gangs can easily manipulate. Another main problem common to virtually all of these architectures is that the False Positive rate (FP) often rises significantly with the growth in SPAM detection for these schemes, resulting in low overall results.

2. Related Work

Considerable work has already been done in this area, and new detection methods are constantly suggested due to the significance of the subject. Nosier et al. [17] suggested a strategy focused on characters. This technique uses a classifier for a multi-neural network. On the basis of a normalized weight obtained from the ASCII meaning of the word characters, each neural network is trained. An intruder may, however, camouflage the phrases, e.g., to avoid identification by writing the words with a slightly different spelling

or by using graphics. This results in comparatively low detection rates that are right. A rule-based system was developed by Ask et al. [18], where 23 carefully chosen characteristics were defined from a privately accumulated spam dataset. One ranking was then allocated to each of the parameters. The cumulative score was compared to a threshold value in order to tag the email as spam or ham. Multilayer Perceptron (MLP), Naïve Bayesian Classifier, and C4.5 Decision Tree Classifier were three machine learning concepts; however, the analysis was performed on a small sample with only 750 spam and ham communications. An e for an e

In terms of time and memory footprint, an efficient performance metric has not been identified yet. Another research indicates that at the term level [19], text mining of emails can be performed. The mining method begins by pre-processing the collection of documents and extracting from the documents the related words. Each paper is then described as a series of words characterizing the document and annotations. The number of occurrences of the words is provided by this procedure. One of the disadvantages, though, is that huge messages will not be treated. Work to detect spamming accounts has already been undertaken.

The Eros method (Early Identification of Spamming) utilizes an algorithm specially developed to identify spamming accounts early on. The identification strategy introduced by Eros amalgamates content-related detection with inter-account contact trend-based features [20]. In the future, this study could be generalized so that it can help the real-time signaling of the account of a spammer. A modern hybrid model, incorporating traditional Negative Selection Algorithms, and a new study [21]

Deferential Evolution was proposed. In the random generation process of the NSA, the proposed model has the unusual feature of applying Details 2019, 10, 209 4 of 17Di Eventual Evolution. The model also maximizes the distance of the detector produced while minimizing detector overlap [21]. However, the problems of picture spamming and click jacking are not discussed by this work. "Through using character variations to disguise the word, attackers often try to avoid word-based filtering systems, such as spelling "mortgage" as "M*o*r*t*g*a*g*e. Another common instance is the term ('Viagra', 'Viagra', 'V I am g r a' or 'VI<bra/>agar') 'Viagra'. This procedure restricts the e

The performance in most content-based approaches. However, manually generated regular expressions (regex) may be of great benefit in the detection of messages obfuscated by spammers by different patterns. A standard expression in this sense is a compact way of representing collections of terms or

phrases that fulfill a certain pattern [22]. This can then be combined locally with a filtering scheme centered on content. To filter spam texts, Rexnord's et al. [22] used such a technique. To automatically produce regular expressions for a specified dataset, they created a novel genetic programming algorithm, called Discoveries. One theoretical expansion to this scheme will be to extract regular expressions from the full text of the messages instead of only the 'topic header information' currently implemented. Click jacking, also known as I Frame Overlay or UI redressing, is a form of attack in which malicious scripts or connections that are not ordinarily accessible overlay a field or toggle. Among the hacker group, the methodology has become very popular. This allows people to click on links or buttons that they are unable to see, typically because the color of the icon is the same as the background color of the website.

3. Proposed Methodology

By building a memory of the past actions of spam emails, the model suggested in this study is educated. In subsequent incoming communications, this would not encourage the same form of behavior, since the model has been inoculated against a new behavior by the recipient. Negative Selection is called this method. The architecture of the Negative Selection algorithm was based on the mammalian learned immune system's self-on-self discrimination behavior [33], as seen in Figure 1. Data 2019, 10, x FOR PEER REVIEW 5 of 17 uses word-based similarity for matching through Euclidian distance. Taking into account header details, such as the source IP, may further boost its efficiency.

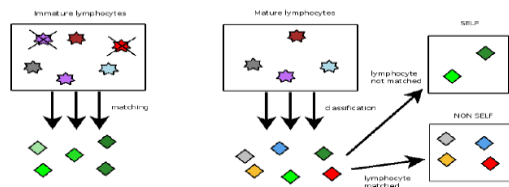


Figure 1. Self and non-self-agents [34], Reproduced with permission from Heba Elshandidy, *Geniuses ONLY! Artificial Immune Systems – PART II*, published by word press, 2011.

The idea behind the Negative Selection method is that unfamiliarity is expected or that it varies from what is common. In anomaly and alteration detection algorithms, this is also a crucial step. The aim is accomplished by creating a model of deviations, alterations, or unknown (non-normal or non-self) data by producing patterns that do not conform to or complement current (self or normal) patterns accessible to an existing entity. By checking for

matches to the non-normal trends, the prepared non-normal model is then used to track current natural data or new data sources. By providing awareness of self and non-self-behavior, Detrimental Selection or Artificial Selection differentiates between normality and exception. In a sequence of learning processes known as teaching, this information may be programmed into a method or created. The teaching phase can be carried out by 'learning' from the contents of die rent databases of self and non-self.

4. Algorithm of Negative Sorting (NSA) Demonstrated

Highly distributed, computationally intelligent approaches or structures focused on evaluation of the actions and association of antigens and antibodies in a biological environment are Artificial Immune Systems (AIS). Negative Selection Algorithms (NSA), a sub-area of AIS, emulates the way dangerous antigens are identified and disposed of by the human body. An antigen may be identified as a material that induces antibodies to be generated against it by the immune system. An antigen may be an environmental material such as pesticides, microbes, viruses (non-self-antigen) or may be created in the body (self-antigen). The immune system does not accept the substance and instead fails to remove the substance [35].

The NSA is produced on the basis of a Thymus system that induces a group of mature T-cells capable of linking or matching to non-self-antigens only. These T-cells are 'trained' to take action if anything out of the ordinary approaches the system or if it senses a change in the anticipated sequence. Antibodies are produced in biological processes to cope with undesirable antigens. The 'binding' results in the non-self-antigen being killed.

The NSA begins by producing a collection of self-strings, S , that describe the system's normal state. The next move is to create a series of detectors, D , which can only recognize or connect with the S complement, that is, S' (non-self). These detectors function in a similar way to mature T-cells that, as soon as a match is detected, are capable of inducing antibodies. Binding happens with the antigen or with spam keywords or blacklisted IPs in this scenario. Inside the detectors, the central logic functions close to the mechanism of biological antibodies. In order to divide them into 'self' or 'non-self', the algorithm can then be extended to new results.

For both spam and ham datasets, the model is trained to construct the knowledge base needed for intelligent activity. When confronted with spam keywords (non-self-antigen) the detectors inside

respond directly rather than with ham keywords (self-antigen). The T-cells in the immune system often go through a maturing phase to learn how to react when presented with self-antigens and non-self-antigens, as described above. In order to overcome undesirable circumstances, the immune system needs more than one structure in operation (e.g., the release of an immature T-Cell into the blood stream). Via the installation of several detectors that serve as a firewall against spam communications, the proposed device also follows such a trend.

5. Design of the Framework

To be trained with self (non-spam or ham) datasets as well as with spam datasets, a model is prepared. These datasets have been investigated historically and are either categorized as spam or non-spam. The well-known Enron email datasets was used for training and constructing spam and non-spam databases. In the data archive of Carnegie Mellon University, USA, you will find the complete raw collection. The six sets used for this analysis were downloaded from the Department of Informatics, University of Economics and Industry of Athens, Greece.

In the following framework, the six datasets downloaded are arranged: And email is transferred into the trainer model where it is analyzed to collect keywords and categorize its contents. As well as the originating IP addresses of these emails, the email addresses are also retrieved. Statistics are provided with a list of terms and their occurrence in the latest collection of emails at the end of each process. This is recorded in different libraries that, as further datasets become accessible in the future, may be modified. It has been reported that more recently revised databases have resulted in better spam identification rates and less false positives. Flagging a real spam message as spam is considered True Positive, whereas False Positive is defined as labeling a valid message as a spam message. The lack of genuine addresses, which is a big problem, results in false positives. The model requires to be trained with modified datasets wherever possible to boost the efficiency of the algorithm and reduce the false positive and false negative levels, so that it is conscious of the behavior of newly arriving risks. In total, 50,409 emails were used for training and research purposes, a sub-set of Enron email corpus. As illustrated in Table 1, 33,792 emails were used for training purposes, or just under 66 percent. The remaining 17,157 emails were used as the test dataset, or slightly over 34 percent.

Table 1. Enron email dataset.

Dataset	Number of Spam	Number of Ham
Enron ₁	1513	3735
Enron ₂	1496	4361
Enron ₃	1500	4012
Enron ₄	4500	1500
Enron ₅	3675	1500
Enron ₆	4500	1500
Total	17,184	16,608

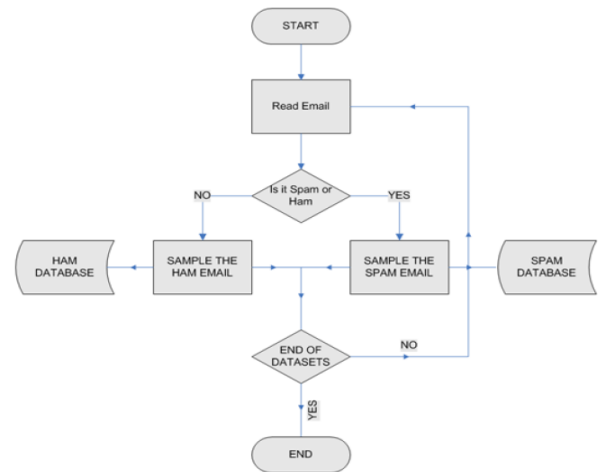


Figure 2. A flowchart of the overall design.

6. Results and Discussion

As can be shown from Figures 3 and 4 as well as Table 3, by minimizing False Positives, FP and False Negatives, FN, training the algorithm with more datasets enhances the efficiency. As seen in Table 3, a total of 17,157 emails were screened from the Enron datasets. With only Enron1, 15,981 emails (True Positive, TP combined with True Negative, TN) were correctly found, whereas 1176 emails were wrongly defined. As can be shown from the same graph, moreover, the resulting inclusion of further datasets raised the proportion of correctly classified emails. After the final dataset was added, Enron6, 16,912 emails were correctly found. True Positive (TP) is the real spam emails in the sense of this report, whereas True Negative (TN) is the actual ham emails.

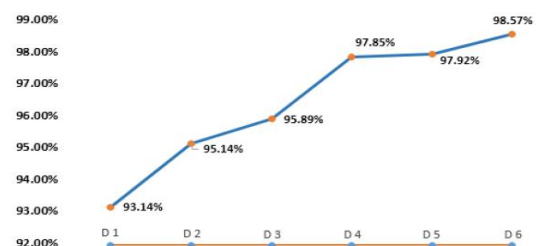


Figure 3. Progressive increase in the correct detection rate.

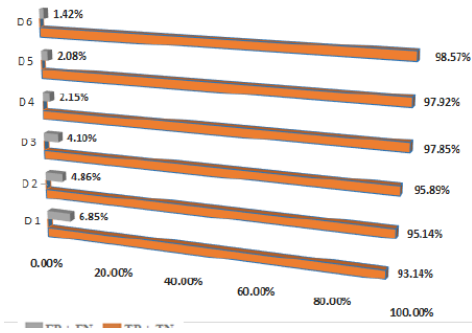


Figure 4. Comparison between (False Positive, FP + False Negative, FN) and (True Positive, TP + True Negative, TN) percentages.

Table 3. Progressing enhancement with the introduction of additional datasets

Datasets	Trained by Number of Datasets	Total Email Scanned: 17,157			
		TP + TN	FP + FN	TP + TN Percentage	FP + FN Percentage
Dataset ₁	Enron ₁	15,981	1176	93.14	6.85
Dataset ₂	Enron ₁ + Enron ₂	16,323	834	95.14	4.86
Dataset ₃	Enron ₁ + Enron ₂ + Enron ₃	16,453	704	95.89	4.10
Dataset ₄	Enron ₁ + Enron ₂ + Enron ₃ + Enron ₄	16,789	368	97.85	2.15
Dataset ₅	Enron ₁ + Enron ₂ + Enron ₃ + Enron ₄ + Enron ₅	16,801	357	97.92	2.08
Dataset ₆	Enron ₁ + Enron ₂ + Enron ₃ + Enron ₄ + Enron ₅ + Enron ₆	16,912	245	98.57	1.42

Both real and false percentages are presented in figure 3. This arrangement results in a slightly lower cumulative incidence of False Positive (FP) and False Negative (FN) than many other proposed structures. The mechanism is often able to capture double or triple word spam sentences, as well as word obfuscation, in addition to being able to identify typical spam terms. It can also be found from Figure 4 that although the percentages of True Positive (TP) and True Negative (TN) combined steadily rose, the percentages of False Negative and False Positive combined decreased. Figure 5 shows graphically the relative contribution of each detector to the accurate classification of an email from Enron1 to Enron6 as spam or ham. The first detector tested the IP address of the source and its existence in the body or header of the email. The second compared the terms of the email body to those of a predefined spam token dataset, while the final detector determined if there were some spam words in at least 30% of the email body's row. As can be shown, the relative contribution of detectors 1 and 2 improved with the inclusion of more datasets, while the reverse was true for detector 3. This means that several of the spam emails were caught during the first two stages of

identification with the improved testing of the algorithm.

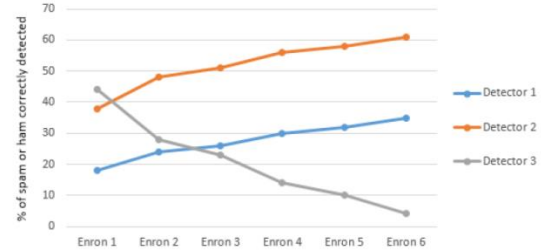


Figure 5. The relative contribution of each of the detectors in identifying spam or ham (True Positive and Negative) upon addition of more datasets.

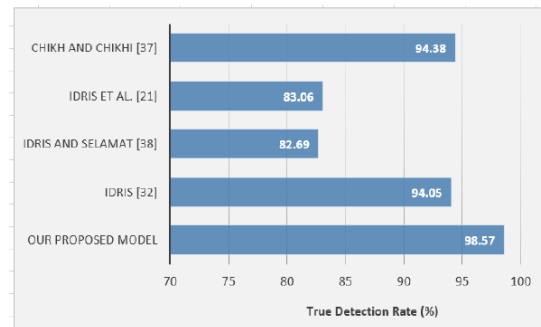


Figure 6. A benchmark comparison with similar studies.

7. Conclusions

Spam is a major concern that is not only irritating to end-users, but also financially detrimental and a danger to protection. Algorithms and processes for machine learning have been shown to be very effective. Based on anomaly detection systems and machine learning concepts, the work described in this paper indicates that the addition of further datasets substantially raises the right detection rate, from 93.14 percent based on one dataset to nearly 98.57 percent after adding the last dataset (Enron6). A new method is the amalgamation of IP-based filtering with a Negative Sorting Algorithm in a controlled approach. We have contrasted our proposed system to other NSA implementations and noticed that in terms of real spam and ham identification, our proposed system outperforms other NSA implementations.

References

1. Tschabitscher, H. *How Many Emails Are Sent Every Day*. 2015. Available online: <https://www.lifewire.com> (accessed on 11 June 2019).
2. Gupta, S.; Pilli, E.S.; Mishra, P.; Pundir, S.; Joshi, R.C. *Forensic Analysis of Email Address Spoofing*. In *Proceedings of the 5th International Conference on Confluence 2014: NGIT Summit, Noida, India, 25–26 September 2014*; pp. 898–904.
3. Smadi, S.; Aslam, N.; Zhang, L. *Detection of Phishing Emails Using Data Mining Algorithms*. In *Proceedings of the 9th International Conference on Software, Knowledge, Informacion*

- Management and Applications, Kathmandu, Nepal, 15–17 December 2015; p. 4.*
4. Bratko, A.; Filipic, B.; Cormack, G.; Lynam, T.; Zupan, B. Spam filtering using statistical data compression models. *J. Mach. Learn. Res.* 2006, 7, 2673–2698.
5. Jagatic, T.; Johnson, N.; Jakobsson, M.; Menczer, F. Social Phishing. *Commun. ACM* 2007, 50, 94–99. [CrossRef]
6. Shan, T.L.; Narayana, G.; Shanmugam, B.; Azam, S.; Yeo, K.C.; Kannoorpatti, K. Heuristic Systematic Model Based Guidelines for Phishing Victims. In *Proceedings of the IEEE Annual India Conference, Bangalore, India, 16–18 December 2016*; pp. 1–6.
7. Leung, C.; Liang, Z. An Analysis of the Impact of Phishing and Anti-Phishing Related Announcements on Market Value of Global Firms. Master's Thesis, HKU, Pok Fu Lam, Hong Kong, 2009.
8. Raad, N.; Alam, G.; Zaidan, B.; Zaidan, A. Impact of spam advertisement through e-mail: A study to assess the influence of the anti-spam on the email marketing. *Afr. J. Bus. Manag.* 2010, 4, 2362–2367.
9. Al-Sharif, S.; Iqbal, F.; Baker, T.; Khattack, A. White-Hat Hacking Framework for Promoting Security Awareness. In *Proceedings of the 8th IFIP International Conference on New Technologies, Mobility and Security (NTMS), Larnaca, Cyprus, 21–23 November 2016*.
10. Ghafir, I.; Prenosil, V.; Hammoudeh, M.; Baker, T.; Jabbar, S.; Khalid, S.; Jaf, S. BotDet: A System for Real Time Botnet Command and Control Trac Detection. *IEEE Access* 2018, 6, 38947–38958. [CrossRef]
11. Foley, C. ABC Bus Companies, Inc.—Cyber Incident Notification. 2018. Available online: <https://www.doj.nh.gov/consumer/security-breaches/documents/abc-bus-20180302.pdf> (accessed on 24 May 2019).
12. French Cinema Chain Fires Dutch Executives Over CEO Fraud. Available online: <https://www.bankinfosecurity.com/blogs/french-cinema-chain-fires-dutch-executives-over-ceo-fraud-p-2681> (accessed on 25 May 2019).