

Identification of Dominant Spam Email Features to Improve Detection Accuracy of Machine Learning Algorithms

MSc Research Project
Cybersecurity

Sujay Hegde
Student ID: 20174217

School of Computing
National College of Ireland

Supervisor: Liam McCabe

National College of Ireland
MSc Project Submission Sheet
School of Computing



Student Name: Sujay Hegde

Student ID: 20174217

Programme: MSc Cybersecurity **Year:** 2021

Module: Research Project

Supervisor: Liam McCabe

Submission Due Date: 16/12/2021

Project Title: Identification of Dominant Spam Email Features to Improve Detection Accuracy of Machine Learning Algorithms

Word Count: 7421 **Page Count:** 29

I hereby certify that the information contained in this (my submission) is information pertaining to the research I conducted for this project. All information other than my contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

ALL internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other authors' written or electronic work is illegal (plagiarism) and may result in disciplinary action.

Signature:

Date: 16/12/2021

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST

Attach a completed copy of this sheet to each project (including multiple copies)	<input type="checkbox"/>
Attach a Moodle submission receipt of the online project submission, to each project (including multiple copies).	<input type="checkbox"/>
You must ensure that you retain a HARD COPY of the project, both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on the computer.	<input type="checkbox"/>

Identification of Dominant Spam
Email Features to Improve Detection
Accuracy of Machine Learning
Algorithms

Abstract

The current business environment and personal use of internet shoots larger numbers of individuals and companies, falling prey to phishing attacks and spam emails. The growing dependence on internet allows the cybercriminals to hatch nasty plans against internet users, by releasing spam mails, with attractive contents and convincing them to fall victims to such incidences. In simpler terms, spam emails are unsolicited commercial/bulk e-mails, which have become a big cause of concern for the users of emails and for that matter, the internet.

The application of spam feature detection will be able to influence the nature of prediction models that can be utilized further for the detection of other effective datasets apart from spam. The identification and recognition of the most dominant feature of spam will be efficient in ensuring that the spam developing system is being disclosed which will further contribute to the research process of the spam detection process.

Contents

	I
1. Introduction	3
1.1. Background	3
1.2. Rationale	4
1.3. Research Questions	5
1.4. Research Objective	5
1.5. Contribution to scientific literature	5
1.6. Structure of report	6
2. Related work	7
2.1. Use of Machine learning algorithms for determining spam email features and classifications. 7	
2.2. Spam email detection and review for improving the detection accuracy	10
2.3. Summary	11
3. Research Methodology	12
3.1. Data Summary	12
3.2. Data Visualization	16
3.3. Data Pre-processing	17
4. Design Specification	18
4.1. Data Transformation	18
5. Implementation	19
5.1. Model training and testing	19
5.1.1. Part 1	20
5.1.2. Part 2	21
6. Evaluation and Discussion.....	24
7. Conclusion and Future work	26
7.1. Conclusive Analysis and Future work	26
8. Acknowledgement	26
References	27

Table of Figures

Figure 1.....	12
Figure 2.....	13
Figure 3.....	13
Figure 4.....	14
Figure 5.....	14
Figure 6.....	15
Figure 7.....	15
Figure 8.....	15
Figure 9.....	16
Figure 10.....	16
Figure 11.....	16
Figure 12.....	16
Figure 13.....	17
Figure 14.....	18
Figure 15.....	18
Figure 16.....	19
Figure 17.....	20
Figure 18.....	21
Figure 19.....	21
Figure 20.....	22
Figure 21.....	23
Figure 22.....	24
Figure 23.....	25
Figure 24.....	25

I. Introduction

Spam detection has become a necessity as it is a filter-based software that is being used to ensure there is effective identification of the unsolicited along with the unwanted email that is usually obtained from unwanted email threats. The presence of spam detection etc. helps in ensuring the prevention of incoming messages into the user's inbox, which might be clicked or opened by the users unintentionally. In machine learning, the spam detection processes are a necessity to hold supervised machine learning problems that will be effective in providing the machine-learning model with multiple examples of spam sets and harm messages, which are to be recognised in the future and kept in the memory of the machine-learning model. As the application of machine learning has been largely common in multiple business applications along with technological grounds, it is important to ensure that the features of spam mails are identified thoroughly to improve the overall accuracy rate of the machine learning process to detect spam as needed. The current research study is developed to ensure a clear understanding of the same by analysing multiple resources.

I.1. Background

Spamming is the process that allows the sending of multiple messages, which are unsolicited and are usually sent in large numbers to multiple numbers of recipients entitled for the commercial advertising process and in many cases for non-commercial advertising feedback. The spam detection technique is mostly applied in the email spamming process, even though it can be applied in multiple contexts considering the wide range of messaging systems that are currently being developed and used through different types of platforms (Rahman *et al.* 2020). The application for spam detection is a necessity in business organisations and for personal usage as well. This is because it helps in filtering a large amount of data and information that are not essential and prohibits business organisations to obtain messages that might contain viruses or any other malicious content, which can lead to corruption in the files that are present in the system being used. The quality of business email is more protected by the use of spam filters or spam detection mechanisms. The advancement in the overall application of technology into different filtering processes has led to the inclusion of machine learning to be one of the most profoundly used mechanisms in ensuring the prospects of spam filtering (Faris *et al.* 2019). It is important to note that the process of machine learning has been highly effective in using spam detection techniques.

Automatic email filtering has been considered highly effective in ensuring effective detection of spam but in the current times, the spammers can bypass the filters easily leading to the requirement of addressing more comprehensive and complex spam detection techniques, which led to the application of the machine learning process. The major approaches that have been applied so far are based on junk mail filtering, text analysis, white and blacklist domain recognition, and the community primary techniques (Jawale *et al.* 2018). Apart from email, SMS spam is also a major consideration that is essential to be considered to ensure that there is productive management of communication systems without the interference of messages that could pose harm to the personal use of emails or business purpose applications. This further led to the application of machine learning applications once again to make the overall prospects of spam detection methods to be fruitful and smooth. The application of machine learning into the spam detection process has been majorly based on its two broad categories, which are supervised learning and the unsupervised learning process. The application of this categorisation of the machine learning approach is further based on the TF-IDF vectorization ethanoic (Vinitha and Renuka, 2019). This technique is used for the generation of word clouds that further stands for Turn Frequency Inverse Document Frequency, which is used for text mining

processes and identifying the world's features. These features are collectively used to ensure that the spam messages are detected thoroughly.

Multiple techniques are present that help in using machine learning to proceed with the prospects of spam detection processes. In terms of the filtering process, the machine learning techniques are based on content-based filtering techniques, case base spam filtering methods, the previous line-based spam filtering techniques, and heuristics or rule-based spam filtering processes. It can be noted that the machine learning approach in detecting spam is based on statistical models, which help in classifying the data more fluently (Dada *et al.* 2019). It is important to note that the spam detection process involves the application of a trained machine-learning model that is focused on generally identifying the sequence of words that are available in the email and are closer to the emails, which are found in the form of spam. Despite the difference in the machine learning approach that has the capability of detecting spam, it can be noted that the Naive Bayes algorithm has been determined to be a more prominent one. The data that is applied and used to provide the machine-learning algorithm with a set of examples of the spam ensure that the future spam is recognized in a similar pattern (Awad and ELseofi, 2011). Hence, the application of machine learning and increasing the accuracy rate of the machine learning process involves the usage of spam features that will be collectively unimportant in holding accountability to more accuracy in the overall process of detection.

1.2. Rationale

The given research paper is significant as spamming has become a major cause of concern for business entities as they often contain messages that hold the presence of information that could be harmful to the company. It is important to note that maintaining accuracy in the spam detection process is a necessity as the advancement in technology and programming languages are improving the overall ability of spam emails to be filtered out from the traditional machine learning approach used. Hence, there is an increasing requirement of addressing the accuracy rate in the spam detection process (Kumar and Sonowal, 2020). Intelligent systems for spam detection and identification of the relevant features are important as it helps in ensuring that the massive data that flows that contains hundreds of individual and large numbers of attributes increases the overall problem of detecting spam and makes the entire process to be complex. There have already been multiple identifications of dominant methods that help in collecting data mining processes, which is one of the most important prospects of prediction methods (Ameen and Kaya, 2018). However, spam contains features that are similar to the system details and data that are most regularly used by individuals. This further reduces the overall predictability nature of the machine learning process. Thus, the given research paper is a necessity as it will collectively identify new algorithms that can be used in the machine learning process so that there is a more efficient format of spam detection taken into account. It will allow the identification of the most advanced and dominant features of spam that is essential to be present in it and can be used in the future to recognize them using machine learning algorithms and techniques (Faris *et al.* 2019). The research paper will be effective in determining the possible solutions that might be utilised with the potentiality of occurrence of any error through the application of the new features and approaches.

1.3. Research Questions

The research questions developed are based on providing a clear direction on the different outcomes that must be addressed in the given paper to ensure the effective completion of the research study.

- 1) What are the common features present in spam email contents that allow collective management identification of the spam?
- 2) What are spam features that must be utilised by machine learning to ensure product development of the techniques helping in the detection of spam with efficiency?
- 3) What are possible challenges that are being faced in identifying the most prominent features that will be effective in detecting spam?
- 4) Why there is a requirement of addressing spam detection features to increase the accuracy of current machine learning approaches?

1.4. Research Objective

Some of the objectives that are to be met through the given research paper are based on the pathway and the information flow that has been maintained in drawing the final analysis of the research paper. The objective of the research paper is stated as follows:

- ❖ To identify the different types of dominant spam email features present in multiple email formats
- ❖ To identify and discuss the prospects of machine learning algorithms that will be productive in identifying the dominant features as needed
- ❖ To identify and discuss the improvement techniques that must be considered for the machine learning technologies in respect to the new dominant features identified in the spam email
- ❖ To discuss the possible challenges that might be evolving in utilizing the newly identified features in detection spam by machine learning technique with more accuracy.

1.5. Contribution to the scientific literature

The contribution that the current research paper will be making is focused on bringing better insights on new technologies and features that can be contributed to the scientific domain of machine learning algorithms for spam or prediction methods.

The use of the machine learning algorithm is developed and implemented through the recognition of a specific mechanism and algorithm, which enables the users of the emails and internet to classify and segregate the emails, differentiating between spam and non-spam emails (Basavaraju and Prabhakar, 2010). Although the machine learning approach is considered to be one the best ways to classify spam emails and identify the dominant spam email features, there are confusions on the selection of the most appropriate approach. Multiple opinions are developed, wherein some critics support the Naïve Bayes approach, while some support the support vector machines and Neural Networks, and the rest opts for K-nearest neighbor, Rough sets and the artificial immune system approaches.

1.6. Structure of the report

The research paper is based on an experiment-focused research study, which addresses that the overall outline of the research study will be involving sections that will identify proper segments of machine learning. Apart from the introduction, the entire research paper has been grouped into six segments. Thus, the entire research paper contains seven segments.

- ❖ The introduction is the first segment, which includes content based on a clear introduction on the research topic, and the overall significance of the research topic has been addressed in the section. It also involves research questions and objectives along with an analysis of the research contribution to scientific literature as well.
- ❖ The literature review segment is the second section of the research paper effective in providing a review of different kinds of literature present on the topic identified.
- ❖ The third section is the research methodology, where the research processes and methods have been mentioned with clear identification of the different methods that have been taken into account to collect the specific data to complete the entire research.
- ❖ The fourth section is the design specification section where the research techniques have been defined based on the overall machine learning designing or framework that will be used to bring accuracy to the spam detection process. The design specification is widely based on the description of any new algorithm that has been introduced in the given paper to make the machine learning approach for spam detection to be more effective.
- ❖ The fifth section of the given paper is the implementation section that will be discussing the implementation of the proposed solution to the new machine learning approach to be taken into account.
- ❖ The sixth section of the given paper is the evaluation section that is effective in generating a comprehensive analysis of the result and the overall findings of the study have been developed and presented in the previous section along with a discussion in it.
- ❖ The seventh segment in the given appear is the conclusion and the future work which will be focused on concluding the entire research paper and further representing strategies that will be effective in maintaining future work of the entire study as needed.

2. Related work

This part of the research work conducts the literature review on the subject, by identifying the leading variables in the study. Based on an array of sources of secondary existing information, the literature review embeds a critical analysis of the use of machine learning algorithms in detecting the spam email features and classifications to avoid malicious invasion of online data. The sources of literature have been meticulously detected to ensure that the information provided is authentic and significant. The literature review part comprises two subsections, which focus on the two major points in the study, related to the use of machine learning algorithms in spam email detection and using the same for improving the detection accuracy.

2.1. Use of Machine learning algorithms for determining spam email features and classifications.

The critics state that creating spam emails is not only used as cunning weapons by the cybercriminals but also are a wastage of time, space, communication bandwidth, and storage in a system (Awad and ELseuofi, 2011). The current global statistics show that the problem of spam emails have been growing, lately, specifically, with the outbreak of the COVID 19 pandemic and the majority of the people, around the world, becoming internet-dependent and opting for work from home, where most of the work is done through the sending of emails (Batra et al.2021). The figures show that the current number of emails, sent across every day through the internet, 40% of them are spam, which accounts for almost 15.4 billion per day, costing the internet and email users a financial loss of \$355 million per year (Batra et al.2021).

Although, the critics state that the automatic email filtering process can prove to be effective in this regard, however, the risk remains, as there is right competition between the spammers and spam-filtering models (Fariset al.2019). Despite the internet and email users have used the ways to block the influx of spam emails, the spammers turned out to be smarter than the users, as they used several tricky methods to overcome the filtering methods like using random sender addresses, to ensure that the spam emails, hit the internet and the email ID of the victim (Uddinet al.2019).

Thus, in the current stages of advancements and developments in the interphase of internet use and email access, critics consider the use of knowledge engineering and machine learning algorithms to be effective ways in filtering the spam emails and classifying the same, to enable the email and internet users, overcome the crisis of spam emails. However, comparing the two methods, the critics have found that the use of machine learning algorithms seems to be more useful and efficient and does not require any specific regulations. It is capable of identifying the dominant spam email features and classifications to ensure effective filtering is done.

The critical aspects of determining the most appropriate method of detecting the features of spam email features, along with classifying the spam emails, based on the implications and the outcome, it has been stated by Castillo et al.(2020), that the use of the naïve Bayes e-mail content classification approach is mostly used for the detection of the layer-3 processing, which is needed for reassembling of the features of the spam emails. Further, this approach helps in recognizing the hardware architecture of naïve Bayes inference engine, for conducting an effective spam email controlling, using the two-way classification process (Dada et al.2019).

On the contrary, Jawalet et al.(2018) states that the use of the support vector machines (SVM) approach is helpful, in detecting the features of the dominant spam emails, allowing a proper classification of the

same. The use of this approach helps in extracting the email sender behavior data based on global sending distribution and enables the technical expert to assign a specific value to each of the IP addresses, from where the spam emails are being sent and the systems, which receive such emails (Bhuiyanet *al.*2018).

The studies on the implementation of the support vector machines (SVM) approach are useful, as it effectively and accurately helps in classifying the emails and recognizing the dominant email features, which seems to be faster than the Random Forests (RF) Classifier process. Further, the application of the K-nearest neighbour and identification of the rough sets, allows the experts to conduct the personalized email prioritization (PEP) process, which focuses on the recognition and analysis of the personal social networks, to identify and capture the user groups and acquire the rich features that allow the email users to identify the specific elements of a spam email (Hussainet *al.*2019).

This specific model encapsulates the formulation of a particular machine-learning algorithm to help in the recognition of the spam email features and classifications, through creating a special viewpoint from the users, followed by a formation of a supervised classification framework, for setting the required priorities and measuring the importance of the emails and the significance of the spam email content (Gibson *et al.*2020). Some other proponents in this field adhere to the adoption of the immune-inspired model. This immune-inspired model refers to the identification of a framework, which is quite innate and adaptive to the changing business environment and transactions.

This framework facilitates the identification of the artificial immune system (IA-AIS) and helps in recognizing the problem of identification of unsolicited bulk e-mail messages (SPAM), among all the other emails, sent and received in between two systems. This artificial immune system (IA-AIS), model helps in delivering the most relevant information, about the aspects of integrating the analogous to macrophages, helping in the formation of a comprehensive framework of identifying the features and classifications of spam emails (Gangavarapuet *al.*2020).

It has been stated by Jáñez-Martinoet *al.*(2021), that the implementation of the artificial immune system (IA-AIS), has helped in identifying 99% of the SPAM emails and facilitating the specific parameter configurations. A comparative assessment between the artificial immune system (IA-AIS), and the naive Bayes approach to identify the machine learning algorithm shows that a lot of debate hovers over the capabilities of both the frameworks, and each is said to compete with the other to identify the features of the SPAM emails and classifying the same (Wang *et al.*2021). Although, the proponents of the artificial immune system (IA-AIS), state that this machine learning algorithm has a greater ability, compared to the naive Bayes approach, to identify the features of SPAM emails along with the classification of the emails, to differentiate between the spam and non-spam emails.

Developing a comparative analysis between the artificial immune system (IA-AIS), and the naive Bayes approach, it needs to be stated that the naive Bayes approach is quite old and had been proposed in the year, 1998 (Ablel-Rheemet *al.*2020). This approach is used effectively to classify spam emails, by checking out the features of the same, wherein the concept of probabilities plays takes the lead role in identifying the features and accordingly classifying between the spam and non-spam emails.

Considering the popularity of the naive Bayes approach, there are layers of probability, in terms of checking across the contents of the emails and recognizing the database. The below formula shall explain the effective use of the naive Bayes approach in the filtering of Spam emails, with non-spam emails.

$$S [T] = \frac{C_{Spam}(T)}{C_{Spam}(T) + C_{Ham}(T)}$$

The formula for naive Bayes approach to conduct spam email filtering

(Source: Ablel-Rheemet *al.*2020)

On the contrary, the use of the artificial immune system (IA-AIS), has been considered to be another capable method, wherein the critics find the naive Bayes approach to be limited. The artificial immune system (IA-AIS), also times called the Neural Network is considered to be a computational method, framed upon the biological neural networks. This machine-learning algorithm accounts for the identification of the interconnection between the artificial neurons within the artificial immune system approach.

This algorithm is highly adaptive and can modify its structures based on the changes in the processes and transmission process of the emails. However, it has been found that this machine learning algorithm engages in refraining from operating, wherein the decision function is found accurately, and classifies all the proper training samples.

$$W_{n+1} = W_n + cX \qquad b_{n+1} = b_n + c$$

Artificial immune system (IA-AIS) approach in machine learning algorithm for filtering of Spam emails

(Source: Sattu, 2020)

```
For (each term t in the message) do {
  If (there exists a detector p, based on base
      String r, matches with t) then {
    If (m is spam) then {
      Increase r's spam score by s-rate;
    } else {
      Increase r's ham score by ns-rate;
    }
  } else{
    if (m is spam) then {
      if(detector p recognizes t and edmf(p,t) > threshold) then {
        The differing characters are added to its corresponding entry in the library of character
        generalization rules;
      } else{
        A new base string t is added into the library of base strings;
      }
    }
  }
  Decrease the age of every base string by a-rate;
}
```

Sample of a spam email, checked by the use of artificial immune system (IA-AIS) approach in filtering the same as Spam or non-spam emails

(Source: Sattu, 2020)

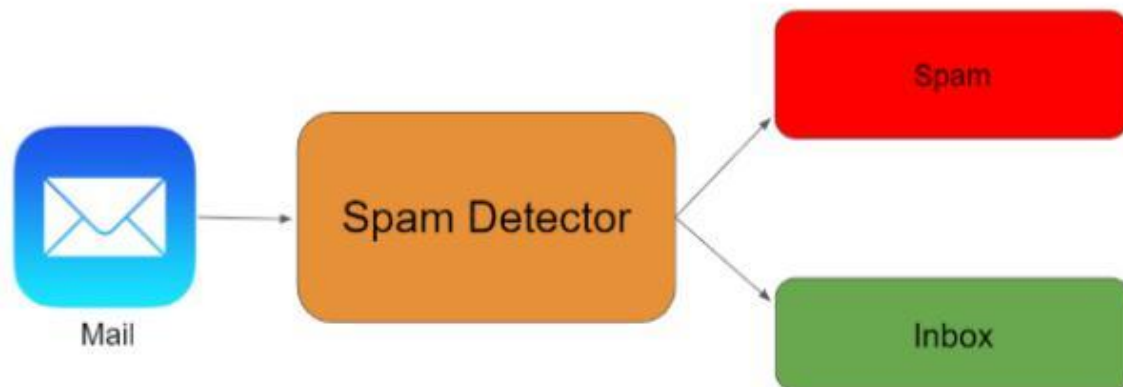
Hence, from the above discussion, it needs to be stated that the selection of the right machine learning algorithm for filtering spam and non-spam emails is a complicated task and needs to be done meticulously. The comparison between the artificial immune system (IA-AIS) approach and the naive Bayes approach is also explained to derive significant conclusions.

2.2. Spam email detection and review for improving the detection accuracy

In congruence to the above discussion, it needs to be stated that the detection of spam emails and phishing attacks, the need for improving the detection accuracy is significant. The use of Natural Language Processing (NLP), helps in processing and converting the texts into useful insights that have become widely accepted and popular (Sheikhiet *al.*2020). Considering the field of AI (artificial intelligence), it has been observed that Natural Language Processing (NLP) is a robust process, delivering complex areas in research, leading to the formation of data, which is contextual.

It needs a proper identification and modification of the contexts of machine-interpretable designs and requires a comparative assessment, facilitating the process of understanding the feature extraction (Mohammed *et al.*2021). To conduct an assessment and identification of spam email features along with the determination of the machine learning algorithms, the process of classification includes the binary and multi-class classification approaches.

Defining the classifications in understanding the language in the systems, it needs to be stated that binary classification refers to the two possible levels of classification, wherein the multi-class classification refers to the process of identifying the cases, wherein more than two labels of assessment are made (Abdullahiet *al.*2021). While detecting the spam emails, within the system, it needs to be stated that the prevention of spam emails and images, from entering into the mail inbox helps in improving the user experience.



Spam detection is the process of identifying the spam and non-spam emails in the mailbox (Source: Alauthman, 2020)

The use of the machine learning algorithm helps in defining the intricacies of using an open-source for detecting spam emails and identifying the target variable. The target variable in the process of detecting the spam and non-spam emails, for the dataset, is considered to be the way of predicting the spam email features (Mustapha *et al.*2020).

	text	spam
0	naturally irresistible your corporate identit...	1
1	the stock trading gunslinger fanny is merril...	1
2	unbelievable new homes made easy im wanting ...	1
3	4 color printing special request additional ...	1
4	do not have money , get software cds from her...	1

The text column that includes the email, spam column, and the target variable to recognize between the spam and non-spam emails

(Source: Yaseen, 2021)

Hence, it is evident that the upsurge in the recognition of measuring the volume of unwanted emails and spam emails that have been creating an intense need for the development of a robust anti-spam framework. These machine learning methods and frameworks, contribute towards the successful detection of spam and non-spam emails, leading to better management of emails (Gangavarapuet *al.*2020). According to the reports from the research by Kaspersky lab, it needs to be stated that there has been an increase in the number of emails, due to the impact of the lockdown and COVID 19 impact.

Further, it needs to be stated that the dependence of the offline methods of operation to the online platform has led to the identification of conducting online reviews, wherein the purchase of products or services has become the primary source of the views and opinions of the users. To conduct an assessment of the spam reviews, it needs to be stated that the email spam reviews are shared on the sites, to ensure proper recognition and assessment of promoting and demoting the products and services (Krithiga and Ilavarasan, 2020). In consideration of the above discussion, It needs to be stated that the approaches of machine learning algorithms are considered to be the most suitable form of reading the features of spam emails and classifying between the spam and non-spam emails, to ensure that the spam and phishing email attacks can be restricted.

The use of the rule-based classifier refers to the process of developing a framework of rules that help in classifying the spam and non-spam emails, based on the features of these emails (Jánez-Martino *et al.*2020). This method enables the technical experts to create a rule, which might be written or non-written and needs to be followed, considering the need for developing a proper classification (Madhavanet *al.*2021). However, the biggest controversy lies in the selection of the most appropriate approach to machine learning and an effective assessment and classification between spam and non-spam emails.

2.3. Summary

Summarizing the points in the literature review, it needs to be stated that the use of machine learning algorithms does play a crucial role in classifying the spam email features. It does help in improving the process of detection accuracy and ensures that the risks of spam emails need to be addressed and mitigated, duly, with the implementation of machine learning algorithms.

3. Research Methodology

This chapter discusses the usage of two datasets for spam domain identification and classification. Two approaches are used namely supervised machine learning algorithms like a decision tree and Naïve Bayes along with the deep learning model of bidirectional gated recurrent unit (GRU).

3.1. Data Summary

The machine learning algorithm is characterized by the elements of identifying the diverse machine learning methods, wherein the classification of the emails, in between the spam and non-spam emails are considered. Some critics state that the use of the supervised learning approach enables accurate and powerful assessments of the datasets, followed by conducting a spam review, based on the features and classification groups of the emails. This process requires two datasets, which include the training data and the testing data. Both the data are used to train the classifier and evaluate the overall performance of the classifier, respectively.

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 517401 entries, 0 to 517400
Data columns (total 2 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   file        517401 non-null object
1   message     517401 non-null object
dtypes: object(2)
memory usage: 7.9+ MB
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5572 entries, 0 to 5571
Data columns (total 4 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   Category    5572 non-null   object
1   Message     5572 non-null   object
2   text        5572 non-null   object
3   spam        5572 non-null   object
dtypes: object(4)
memory usage: 174.2+ KB
```

Figure 1. Data Frame Summary -1

The figures above depict the data frame summary for the two-dataset used for the email spam domain identification and classification. The first dataset consists of the email list along with the email domain and email aspects whereas the second dataset consists of the category, message, and text in the emails for the classification.


```
allen-p/_sent_mail/1.
Message-ID: <18782981.1075855378110.JavaMail.evans@thyme>
Date: Mon, 14 May 2001 16:39:00 -0700 (PDT)
From: phillip.allen@enron.com
To: tim.belden@enron.com
Subject:
Mime-Version: 1.0
Content-Type: text/plain; charset=us-ascii
Content-Transfer-Encoding: 7bit
X-From: Phillip K Allen
X-To: Tim Belden <Tim Belden/Enron@EnronXGate>
X-cc:
X-bcc:
X-Folder: \Phillip_Allen_Jan2002_1\Allen, Phillip K.\'Sent Mail
X-Origin: Allen-P
X-FileName: pallen (Non-Privileged).pst
```

Here is our forecast

Figure 2. Data Frame Summary -2

The pattern shown below searches for substrings called metadata, which has the form as:

Begin with a beginning of a line - signified by caret ^ (also, for this we need to specify re.M flag)

Contain any number (but at least one) of any signs after this beginning of a line - signified by a dot (any sign) and a plus (at least one, but no upper limit)

Finish with a colon, which we do not take as a part of these substrings (so-called look-ahead) - this is the meaning of '(?:)' part

```
metadata_pattern = re.compile(r'^.+(?::)', re.M) # this way we will get greedy behavior
metadata_names = metadata_pattern.findall(data.message[0])
metadata_names

['Message-ID',
 'Date: Mon, 14 May 2001 16:39',
 'From',
 'To',
 'Subject',
 'Mime-Version',
 'Content-Type',
 'Content-Transfer-Encoding',
 'X-From',
 'X-To',
 'X-cc',
 'X-bcc',
 'X-Folder',
 'X-Origin',
 'X-FileName']
```

Figure 3. Code Snippet-1

```

# we place ? after the quantifier (+)
metadata_pattern = re.compile(r'^.+?(?=:)', re.M)
metadata_names = metadata_pattern.findall(data.message[0])
metadata_names

['Message-ID',
 'Date',
 'From',
 'To',
 'Subject',
 'Mime-Version',
 'Content-Type',
 'Content-Transfer-Encoding',
 'X-From',
 'X-To',
 'X-cc',
 'X-bcc',
 'X-Folder',
 'X-Origin',
 'X-FileName']

```

Figure 4. Code Snippet-2

The following loop extracts strings that begin right after each of the metadata categories we extracted (stored in the `metadata_names` list) plus a colon and whitespace, and continue until the end of the line (marked by a dollar sign).

Again, we need to pass it the `re.M` flag for the multiline special character (caret previously, dollar sign now) to work as we intend it to.

```

for metadata_name in metadata_names:
    pattern = re.compile(r'(?<=%s:\s).+$'%metadata_name, re.M)
    print(metadata_name, pattern.search(data.message[0]))

Message-ID <re.Match object; span=(12, 57), match='<18782981.1075855378110.JavaMail.evans@thyme>'>
Date <re.Match object; span=(64, 101), match='Mon, 14 May 2001 16:39:00 -0700 (PDT)'>
From <re.Match object; span=(108, 131), match='phillip.allen@enron.com'>
To <re.Match object; span=(136, 156), match='tim.belden@enron.com'>
Subject None
Mime-Version <re.Match object; span=(181, 184), match='1.0'>
Content-Type <re.Match object; span=(199, 227), match='text/plain; charset=us-ascii'>
Content-Transfer-Encoding <re.Match object; span=(255, 259), match='7bit'>
X-From <re.Match object; span=(268, 283), match='Phillip K Allen'>
X-To <re.Match object; span=(290, 330), match='Tim Belden <Tim Belden/Enron@EnronXGate>'>
X-cc None
X-bcc None
X-Folder <re.Match object; span=(356, 409), match='\\Phillip_Allen_Jan2002_1\\Allen, Phillip K.\\'Se'>
X-Origin <re.Match object; span=(420, 427), match='Allen-P'>
X-FileName <re.Match object; span=(440, 467), match='pallen (Non-Privileged).pst'>

```

Figure 5. Code Snippet-3

```

weekday_pattern = re.compile(r'\A\w+(?=(,|?))', re.M) # from the beginning
monthday_pattern = re.compile(r'(?<=,|s)\d+\b', re.M) # from the first
month_pattern = re.compile(r'(?<=\d\s)\w+(?=\s\d)', re.M)
year_pattern = re.compile(r'(?<=[A-Z][a-z]{2})\s\d+(?=\s\d\d:)', re.M)
hour_pattern = re.compile(r'(?<=\d{4})\s\d\d(?=\s\d\d)', re.M)
minute_pattern = re.compile(r'(?<=\d\d)\s\d\d(?=\s\d\d)', re.M)
second_pattern = re.compile(r'(?<=\d\d:\d\d)\s\d\d(?=\s)', re.M)

```

Figure 6. Code Snippet-4

The figure depicts the data transformation and data cleaning from the input text of the email messages and the alphanumeric character and the first comma, first comma, whitespace, composed of one or more digits to a word boundary.

```

time_multi_pattern = re.compile(
    pattern=r'(?P<Weekday>\A[A-Z][a-z]{2})\W+?'
    flags=re.M
)

time_multi_pattern.groupindex

mappingproxy({'Hour': 5,
             'Minute': 6,
             'Month': 3,
             'Monthday': 2,
             'Second': 7,
             'Weekday': 1,

```

Figure 7. Code Snippet-5

The figure above depicts the time stamping of the emails messages based on the hour, month, month day, weekday, and year.

	Category	Message
0	ham	Go until jurong point, crazy.. Available only ...
1	ham	OK lar... Joking wif u oni...
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...
3	ham	U dun say so early hor... U c already then say...
4	ham	Nah I don't think he goes to usf, he lives aro...

Figure 8. Spam and Ham Classification

The figure depicts the email messages along with their category based on the ham and spam classification. The categorization is based on the sender and receiver email domain. The data set is pre-processed based on the text message contained in the emails using the NLP approaches and machine learning and deep learning models.

3.2. Data Visualization

Messages not longer than 100000 characters:

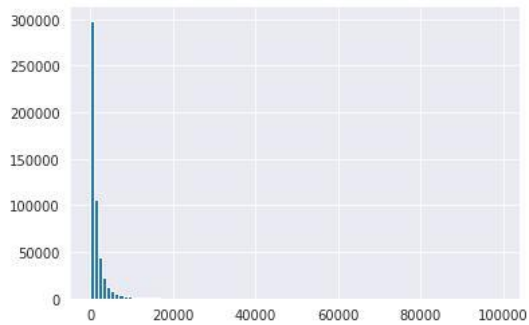


Figure 9. Data Visualization Graph -1

Messages not longer than 10000 characters:

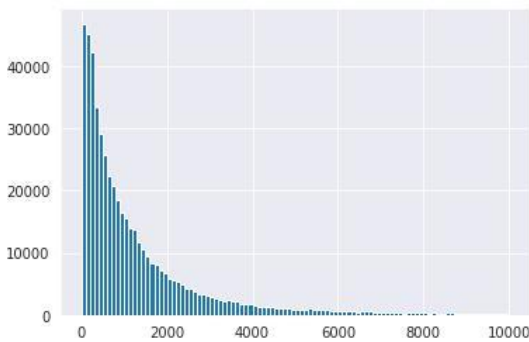


Figure 10. Data Visualization Graph -2

Messages not longer than 1000 characters:

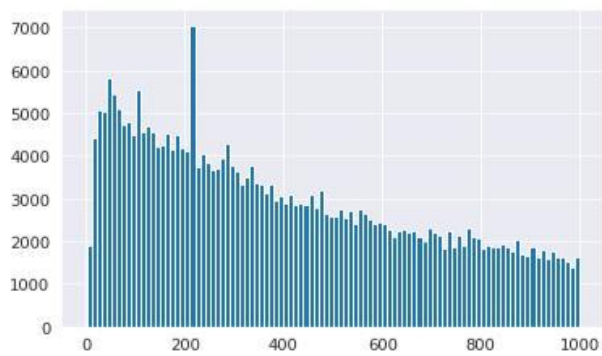


Figure 11. Data Visualization Graph -3

Messages not longer than 100 characters:

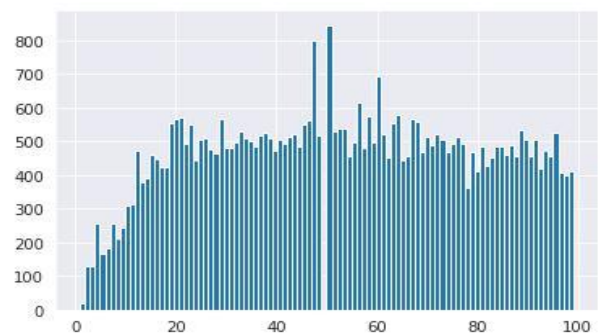


Figure 12. Data Visualization Graph -4

3.3. Data Pre-processing

Data processing is an important factor in many planning and recovery issues. Some data have a similar effect, some have a misleading effect and some do not affect identification or reduction and help in selecting the right size and small features. The scheduling or reversing problem involves more time and less performance when a large number of features are used, but less time for food and smaller size and more performance with better features. The non-deterministic polynomial (NP) problem is the selection of complete materials that can achieve the maximum performance of the partition or model.

Subsequently, the data was used for classification and feature selection in this data mining process. Impact factors contributing to the volatility of the email domain identification were also estimated using values included in the range of maximum impact on data-mine conditions. Another method used to obtain information is to extract the text using NLP (Natural Language Processing) approach.

Sender Classifier is analyzed using the email domain, it is done by taking the two most common e-mail senders in the dataset and building a model, which, based on the message content, tries to predict the person who sent it.

```
# Extract the labels of the two most productive e-mail senders

label_0, label_1 = list(data['From'].value_counts().index)[:2]
label_0, label_1

('kay.mann@enron.com', 'vince.kaminski@enron.com')

# A function to remove stopwords
def remove_stopwords(txt):
    for stopword in stopwords_eng:
        while stopword in txt:
            txt = txt.replace(stopword, '')
    return txt

# Number of messages per one class - We chose 1280, because after splitting it 8:1:1 into train:val:test
n_per_class = 1280

# Take that many messages sent by this person after shuffling them (.sample() method)
messages_0 = data.query('From==@label_0')['Message Content'].sample(frac=1)[:n_per_class].values
messages_1 = data.query('From==@label_1')['Message Content'].sample(frac=1)[:n_per_class].values

# Remove stopwords by applying the function defined above
messages_0 = [remove_stopwords(s) for s in messages_0]
messages_1 = [remove_stopwords(s) for s in messages_1]
```

Figure 13. Code Snippet -6

The above figure consists of functions for the removal of the stop words which are considered to be insignificant based on their value. The messages contained in the emails are segregated based on the training, testing, and evaluating data subsets.

4. Design Specification



Figure 14. Design Specification Diagram

The figure above represents the data analysis along with data pre-processing and exploratory data analysis. The figure describes the steps to process the data before training, testing, and evaluating the datasets for the spam email domain identification and classification.

4.1. Data Transformation

```
# Analysing the shape of the classes
for X, y in [train_X, train_y], [val_X, val_y], [test_X, test_y]:
    print(X.shape, y.shape)
    print(np.bincount(y.astype(np.int32)))
```

```
(2048, 256) (2048,)
[1024 1024]
(256, 256) (256,)
[128 128]
(256, 256) (256,)
[128 128]
```

```
## some config values
embed_size = 100 # how big is each word vector
max_feature = 50000 # how many unique words to use
max_len = 2000 # max number of words in a question to use
```

```
from keras.preprocessing.text import Tokenizer
tokenizer = Tokenizer(num_words=max_feature)

tokenizer.fit_on_texts(x_train)

x_train_features = np.array(tokenizer.texts_to_sequences(x_train))
x_test_features = np.array(tokenizer.texts_to_sequences(x_test))

x_train_features[0]
```

Figure 15. Code Snippet-7

The data transformation includes the analysis of the data shape based on the training, testing, and evaluating data subsets. The text in the emails is pre-processed using the tokenizing technique for the data features analysis.

5. Implementation

Future spam domain classification often appears in the collection of historical information and heuristic tests in the data sets. Positive results from the various machine learning models in Chapter 3 and Chapter 4 of this thesis are taken from the systematic approach provided in Section 5.1. of the machine learning and deep learning models. An extension of this quantitative assessment model in collaboration with linguistic features is presented in this chapter. This chapter covers feature extraction and feature selection for spam domain identification, deep learning modeling development using reference data, feature integration and prediction formatting, and, finally, the artistic results obtained in the exploitation of these features.

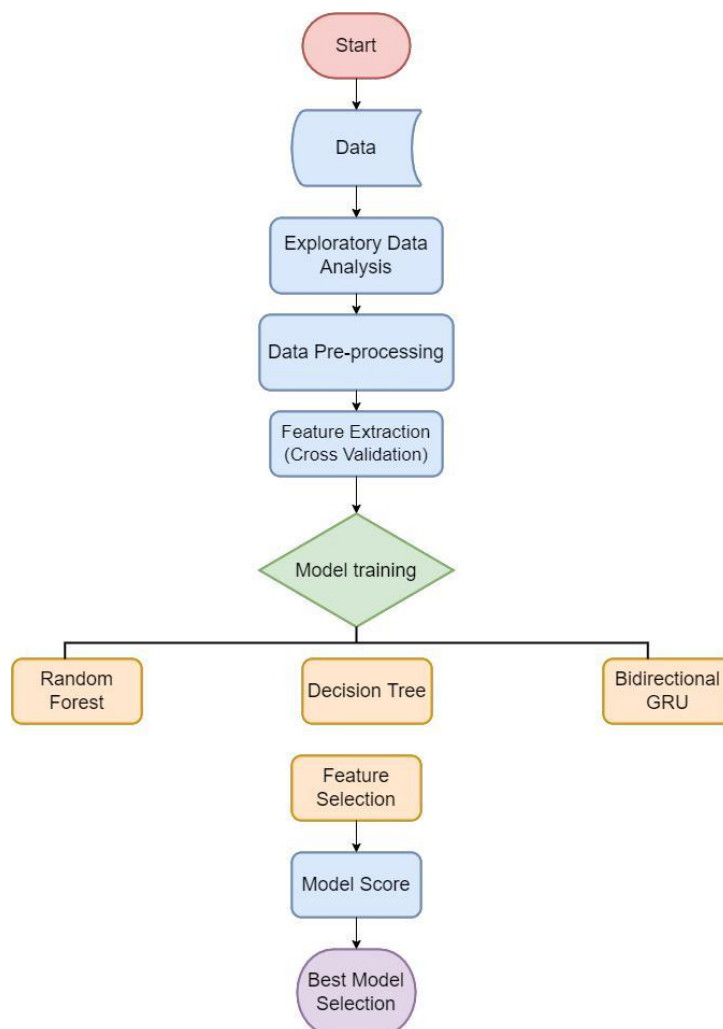


Figure 16. Flow Diagram

The machine learning and deep learning algorithm utilised for the spam domain specification with the selected two datasets are mentioned below.

1. Decision Tree model
2. Naïve Bayes model
3. Bidirectional Gated Recurrent Unit (GRU) model

5.1. Model training and testing

5.1.1. Part I

Additionally, it has been found that the use of the Neural networks within the machine learning algorithms, the use of machine learning framework would help in identifying spam and non-spam emails.

```
# Maximum number of words, we are going to embed
max_features = 2048

# Number of embedding dimensions in the word embedding space constructed by the embedding layer
embed_dim = 256

# The length of the sequence (message) - The longer messages will be cropped to 256, while the shorter ones will
maxlen = 256

# Take messages from both classes and shuffle them before feeding to tokenizer
messages_all = messages_0+messages_1
np.random.shuffle(messages_all)

# The tokenizer ascribes a number to each token (word) in the sequence
tokenizer = Tokenizer(num_words=max_features)
tokenizer.fit_on_texts(messages_all)
word_index = tokenizer.word_index # This dictionary translates each word to its index (corresponding number)

# Transform messages into sequences of numbers corresponding to its particular words
seqs_0 = tokenizer.texts_to_sequences(messages_0)
seqs_1 = tokenizer.texts_to_sequences(messages_1)

# Pad sequences, i.e. make them exactly 256 tokens long (as described above)
seqs_0 = pad_sequences(seqs_0, maxlen=maxlen)
seqs_1 = pad_sequences(seqs_1, maxlen=maxlen)
```

Figure 17. Code Snippet-8

2-layer stacked model with bidirectional GRU as the base model followed by two Dense layers, regularized with batch normalization, L2 dropout


```
Model: "sequential"
```

Layer (type)	Output Shape	Param #
embedding (Embedding)	(None, 256, 256)	524288
bidirectional (Bidirectional)	(None, 256, 64)	55680
bidirectional_1 (Bidirectional)	(None, 64)	18816
dense (Dense)	(None, 64)	4160
batch_normalization (Batch Normalization)	(None, 64)	256
dropout (Dropout)	(None, 64)	0
dense_1 (Dense)	(None, 32)	2080
batch_normalization_1 (Batch Normalization)	(None, 32)	128
dropout_1 (Dropout)	(None, 32)	0
dense_2 (Dense)	(None, 1)	33

```

Total params: 605,441
Trainable params: 605,249
Non-trainable params: 192

```

Figure 18. Model Information

The use of these machine learning techniques helps in learning and identifying spam mails and phishing messages by analysing loads of such messages, throughout a vast collection of computers (Ora, 2020). The use of the machine learning technique leads to the identification of the frameworks, which are robust and more organized in form. This process is a significant contribution towards the development of a compact framework, it is useful to iterate the processes of identifying and classifying between spam and non-spam emails.

5.1.2. Part 2

```

# Splitting the data into training and testing data subset, we segregated 80% training data and 20% as testing data
from sklearn.model_selection import train_test_split
emails_train, emails_test, target_train, target_test = train_test_split(data.text, data.spam, test_size = 0.2)

# Checking the shape of training data subset
emails_train.shape

(4457,)

# Checking the shape of testing data subset
emails_test.shape

(1115,)

```

Figure 19. Code Snippet-9

The uses of the Naïve Bayes (NB) approach have been highly successful in this process.

The selection of features and identification of the spam emails are conducted, through the stages of gathering and crawling the dataset.

Once, the gathering and the crawling are done, the features of the spam emails are extracted from the dataset, by applying the engineering approach, within the system.

```

# Fitting Naive Bayes to the Training set
from sklearn.naive_bayes import GaussianNB
classifier = GaussianNB()
classifier.fit(x_train_features, train_y)

# Predicting the Test set results
y_pred = classifier.predict(x_test_features)

# Making the Confusion Matrix
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(test_y, y_pred)
print(cm)

[[ 31 952]
 [  0 132]]

from sklearn.metrics import accuracy_score
nb_acc = accuracy_score(test_y, y_pred)
nb_acc*100

14.618834080717487

```

Figure 20. Naive Bayes Confusion Matrix

The figure above depicts the confusion matrix for the Naïve Bayes model for the classification of the spam email domain along with the Naïve Bayes model accuracy of 15.61% which is very poor.

This entire process leads to measuring the performance of the classifier, which helps in recognizing the features of the spam emails and show a different supervisory learning technique within the entirety of the process.

Another form of using the machine learning algorithm is the use of the decision tree classifier framework along with the rule-based classifier approach. Both the approaches are useful in the detection of spam and non-spam emails.

```
# Fitting Decision tree to the Training set
from sklearn.tree import DecisionTreeClassifier
classifier = DecisionTreeClassifier(criterion = 'entropy', random_state = 0)
classifier.fit(x_train_features, train_y)

# Predicting the Test set results
y_pred = classifier.predict(x_test_features)
```

```
# Making the Confusion Matrix
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(test_y, y_pred)
print(cm)
```

```
[[937  46]
 [ 47  85]]
```

```
from sklearn.metrics import accuracy_score
dt_acc = accuracy_score(test_y, y_pred)
dt_acc*100
```

```
91.65919282511211
```

Figure 21. Decision Tree Confusion Matrix and Accuracy

The figure above depicts the confusion matrix for the Decision tree model for the classification of the spam email domain along with the Decision tree model accuracy of 91.65% which is very poor.

The use of the decision tree classifier framework helps the researcher to develop a hierarchical decomposition of the training data space and is essentially used to recognize and acknowledge the authenticity of the review.

The process is integrally located and connected with the unique features of the system, allowing the tester to review the dataset that is present and assess the impact of the inverse document frequency, to facilitate the recognition of the spam and non-spam emails.

6. Evaluation and Discussion

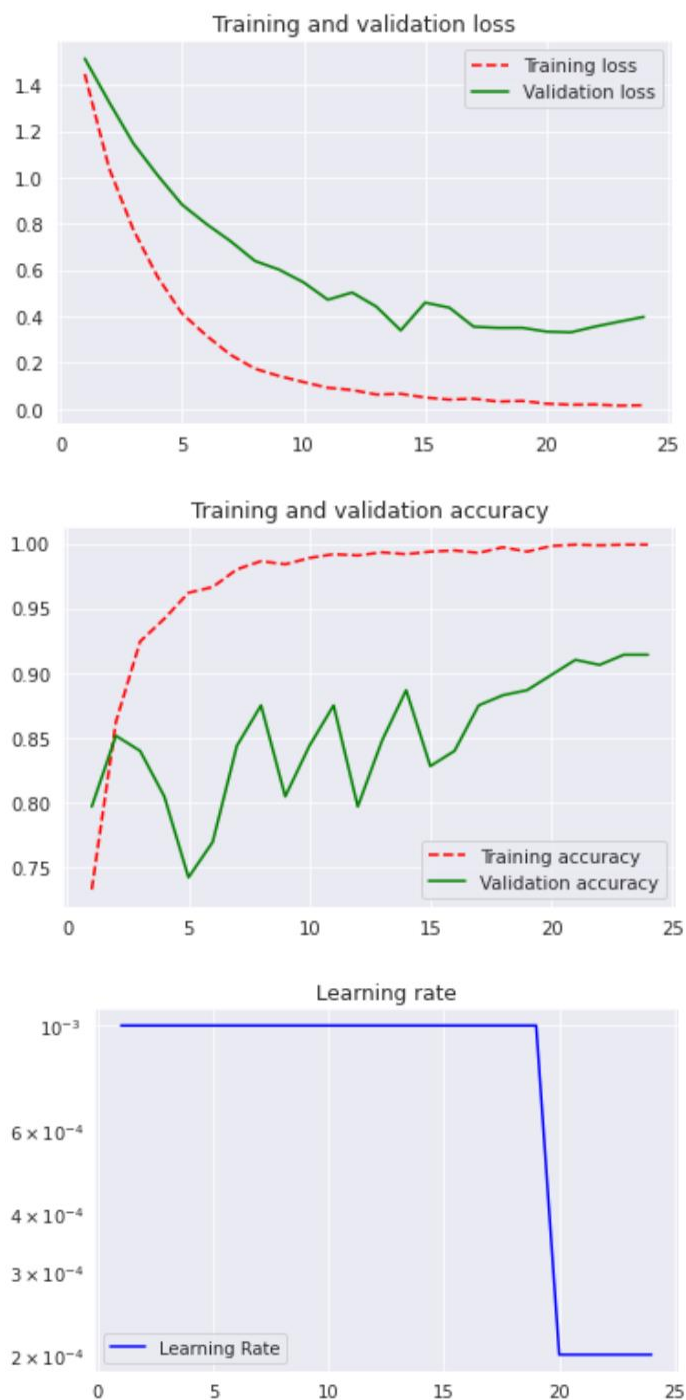


Figure 22. Training Accuracy and Loss Plots

The above three figures depict the variation of the training and testing loss plot and accuracy plot along the learning rate curve. Through the plots, we can observe that validation loss is higher as compared to the training dataset. Similarly, the training accuracy score is higher as compared to the evaluating accuracy score.

The learning rate curve describes the variation of the learning rate for the number of epochs for the gated recurrent unit (GRU) network.

Classification Report				
	precision	recall	f1-score	support
0	1.00	0.03	0.06	983
1	0.12	1.00	0.22	132
accuracy			0.15	1115
macro avg	0.56	0.52	0.14	1115
weighted avg	0.90	0.15	0.08	1115

Figure 23. Naive Bayes Classification Report

The figure above depicts the classification report for the Naïve Bayes model with a 15% accuracy score for the training dataset, a precision score of 1, a recall value of 1, and the F1-Score of 0.06. The model performance of the Naïve Bayes model is poor as compared to the Decision tree model.

Classification Report				
	precision	recall	f1-score	support
0	0.95	0.95	0.95	983
1	0.65	0.64	0.65	132
accuracy			0.92	1115
macro avg	0.80	0.80	0.80	1115
weighted avg	0.92	0.92	0.92	1115

Figure 24. Decision Tree Classification Report

The figure above depicts the classification report for the Decision tree model with a 92% accuracy score for the training dataset, a precision score of 0.95, a recall value of 0.95, and the F1-Score of 0.64. The model performance of the Decision tree model is better as compared to the Naïve Bayes model.

7. Conclusion and Future work

This chapter summarizes the major advances of this study. It also recognizes the major contribution of mathematical models developed with this concept. In conclusion, it suggests further improvements to existing models as well as potential future research topics that may be of interest for further research in the future. The utilization of two data represents the domain identification of the spam emails along with the spam and ham email classification.

The GRU model is used for the email domain identification along with the NLP (Natural Language Processing) techniques. The supervised machine learning models, namely the Decision tree model and the Naïve Bayes model are used for the email classification with 91% and 16% accuracy respectively.

7.1. Conclusive Analysis and Future work

Additionally, this process of using the machine learning algorithm is also accompanied by intricate complications, followed by challenges and issues in future work projects. Some of the leading issues and challenges concerning the identification of spam and non-spam emails, the unavailability of datasets, along the limited data attributes are some of the recognized issues. The lack of adequate datasets is one of the leading challenges in reviewing the features of spam and non-spam emails along with classifying the respective emails. Further, critics have found that the limited data attributes also account for the identification of the challenges in the process of identifying emails, as either spam or non-spam emails.

8. Acknowledgement

I would like to thank my Supervisor for his help and advice with this thesis completion. I would also like to thank my colleague; it would not have been possible without them. I also appreciate all the support I have received from the rest of my family. In conclusion, I would like to thank the college for the scholarship that allowed me to conduct this thesis.

References

- Abdullahi, M., Mohammed, A.D., Bashir, S.A. and Abisoye, O.O., 2021, February. A Review on Machine Learning Techniques for Image-Based Spam Emails Detection. In *2020 IEEE 2nd International Conference on Cyberspace (CYBER NIGERIA)* (pp. 59-65). IEEE.
- Ablel-Rheem, D.M., Ibrahim, A.O., Kasim, S., Almazroi, A.A. and Ismail, M.A., 2020. Hybrid Feature Selection and Ensemble Learning Method for Spam Email Classification. *International Journal*, 9(1.4).
- Alauthman, M.O.H.A.M.M.A.D., 2020. Botnet spam E-mail detection using deep recurrent neural network. *Int. J*, 8(5), pp.1979-1986.
- Ameen, A.K. and Kaya, B., 2018, September. Spam detection in online social networks by deep learning. In *2018 International Conference on Artificial Intelligence and Data Processing (IDAP)* (pp. 1-4). IEEE.
- Awad, W.A. and ELseuofi, S.M., 2011. Machine learning methods for spam e-mail classification. *International Journal of Computer Science & Information Technology (IJCSIT)*, 3(1), pp.173-184.
- Basavaraju, M. and Prabhakar, D.R., 2010. A novel method of spam mail detection using text based clustering approach. *International Journal of Computer Applications*, 5(4), pp.15-25.
- Batra, J., Jain, R., Tikkiwal, V.A. and Chakraborty, A., 2021. A comprehensive study of spam detection in e-mails using bio-inspired optimization techniques. *International Journal of Information Management Data Insights*, 1(1), p.100006.
- Bhuiyan, H., Ashiquzzaman, A., Juthi, T.I., Biswas, S. and Ara, J., 2018. A survey of existing e-mail spam filtering methods considering machine learning techniques. *Global Journal of Computer Science and Technology*.
- Castillo, E., Dhaduvai, S., Liu, P., Thakur, K.S., Dalton, A. and Strzalkowski, T., 2020, May. Email Threat Detection Using Distinct Neural Network Approaches. In *Proceedings for the First International Workshop on Social Threats in Online Conversations: Understanding and Management* (pp. 48-55).
- Dada, E.G., Bassi, J.S., Chiroma, H., Adetunmbi, A.O. and Ajibuwa, O.E., 2019. Machine learning for email spam filtering: review, approaches and open research problems. *Heliyon*, 5(6), p.e01802.
- Faris, H., Ala'M, A.Z., Heidari, A.A., Aljarah, I., Mafarja, M., Hassonah, M.A. and Fujita, H., 2019. An intelligent system for spam detection and identification of the most relevant features based on evolutionary random weight networks. *Information Fusion*, 48, pp.67-83.
- Gangavarapu, T., Jaidhar, C.D. and Chanduka, B., 2020. Applicability of machine learning in spam and phishing email filtering: review and approaches. *Artificial Intelligence Review*, 53(7).
- Gibson, S., Issac, B., Zhang, L. and Jacob, S.M., 2020. Detecting Spam Email with Machine Learning Optimized with Bio-Inspired Meta-Heuristic Algorithms. *IEEE Access*.
- Halim, Z., Waqar, M. and Tahir, M., 2020. A machine learning-based investigation utilizing the in-text features for the identification of dominant emotion in an email. *Knowledge-Based Systems*, 208, p.106443.
- Hussain, N., Mirza, H.T., Hussain, I., Iqbal, F. and Memon, I., 2020. Spam review detection using the linguistic and spammer behavioral methods. *IEEE Access*, 8, pp.53801-53816.
- Hussain, N., TurabMirza, H., Rasool, G., Hussain, I. and Kaleem, M., 2019. Spam review detection techniques: A systematic literature review. *Applied Sciences*, 9(5), p.987.

- Jáñez-Martino, F., Alaiz-Rodríguez, R., González-Castro, V. and Fidalgo, E., 2021, August. Trustworthiness of spam email addresses using machine learning. In *Proceedings of the 21st ACM Symposium on Document Engineering* (pp. 1-4).
- Jáñez-Martino, F., Fidalgo, E., González-Martínez, S. and Velasco-Mata, J., 2020. Classification of spam emails through hierarchical clustering and supervised learning. *arXiv preprint arXiv:2005.08773*.
- Jawale, D.S., Mahajan, A.G., Shinkar, K.R. and Katdare, V.V., 2018. Hybrid spam detection using machine learning. *International Journal of Advance Research, Ideas and Innovations in Technology*, 4(2), pp.2828-2832.
- Kumar, N. and Sonowal, S., 2020, July. Email Spam Detection Using Machine Learning Algorithms. In *2020 Second International Conference on Inventive Research in Computing Applications (ICIRCA)* (pp. 108-113). IEEE.
- Madhavan, M.V., Pande, S., Umekar, P., Mahore, T. and Kalyankar, D., 2021. Comparative Analysis of Detection of Email Spam With the Aid of Machine Learning Approaches. In *IOP Conference Series: Materials Science and Engineering* (Vol. 1022, No. 1, p. 012113). IOP Publishing.
- Mohammed, M.A., Ibrahim, D.A. and Salman, A.O., 2021. Adaptive intelligent learning approach based on visual anti-spam email model for multi-natural language. *Journal of Intelligent Systems*, 30(1), pp.774-792.
- Mustapha, I.B., Hasan, S., Olatunji, S.O., Shamsuddin, S.M. and Kazeem, A., 2020. Effective Email Spam Detection System using Extreme Gradient Boosting. *arXiv preprint arXiv:2012.14430*.
- Noekhah, S., bintiSalim, N. and Zakaria, N.H., 2020. Opinion spam detection: Using multi-iterative graph-based model. *Information Processing & Management*, 57(1), p.102140.
- Ora, A., 2020. *Spam detection in short message service using natural language processing and machine learning techniques* (Doctoral dissertation, Dublin, National College of Ireland).
- Rahman, R.U., Verma, R., Bansal, H. and Tomar, D.S., 2020. Classification of spamming attacks to blogging websites and their security techniques. In *Encyclopedia of Criminal Activities and the Deep Web* (pp. 864-880). IGI Global.
- Sattu, N., 2020. *A Study of Machine Learning Algorithms on Email Spam Classification* (Doctoral dissertation, Southeast Missouri State University).
- Shehnepoor, S., Salehi, M., Farahbakhsh, R. and Crespi, N., 2017. NetSpam: A network-based spam detection framework for reviews in online social media. *IEEE Transactions on Information Forensics and Security*, 12(7), pp.1585-1595.
- Sheikhi, S., Kheirabadi, M.T. and Bazzazi, A., 2020. An effective model for SMS spam detection using content-based features and averaged neural network. *International Journal of Engineering*, 33(2), pp.221-228.
- Uddin, S., Khan, A., Hossain, M.E. and Moni, M.A., 2019. Comparing different supervised machine learning algorithms for disease prediction. *BMC medical informatics and decision making*, 19(1), pp.1-16.
- Vinitha, V.S. and Renuka, D.K., 2019, April. Performance Analysis of E-Mail Spam Classification using different Machine Learning Techniques. In *2019 International Conference on Advances in Computing and Communication Engineering (ICACCE)* (pp. 1-5). IEEE.
- Wang, C., Zhang, D., Huang, S., Li, X. and Ding, L., 2021, May. Crafting Adversarial Email Content against Machine Learning Based Spam Email Detection. In *Proceedings of the 2021 International Symposium on Advanced Security on Software and Systems* (pp. 23-28).

Yaseen, Q., 2021. Spam Email Detection Using Deep Learning Techniques. *Procedia Computer Science*, 184, pp.853-858.