

Detecting Social Engineering Tactics for Email Phishing: Using Machine Learning and Deep Learning

Pooja Sharma

M.Sc. in Computing in
Data Science

2022



Ollscoil
Teicneolaíochta
an Atlantaigh

Atlantic
Technological
University

Dún na nGall

Donegal

Department of Computing, ATU Donegal, Port Road, Letterkenny, Co. Donegal, Ireland.

Detecting social engineering tactics for email phishing: Using Machine Learning and Deep Learning

Author: Pooja Sharma

Supervised by: Jade Lyons

A thesis submitted in partial fulfilment of the
requirements for the Master of Science in Computing in
Data Science

Submitted to Atlantic Technological University

Arna chur isteach chuig Ollscoil Teicneolaíochta an Atlantaigh

September 2022

Declaration

I hereby certify that the material, which I now submit for assessment on the programmes of study leading to the award of Master of Science in Computing in Data Science, is entirely my own work and has not been taken from the work of others except to the extent that such work has been cited and acknowledged within the text of my own work. No portion of the work contained in this thesis has been submitted in support of an application for another degree or qualification to this or any other institution. I understand that it is my responsibility to ensure that I have adhered to ATU's rules and regulations.

I hereby certify that the material on which I have relied on for the purpose of my assessment is not deemed as personal data under the GDPR Regulations. Personal data is any data from living people that can be identified. Any personal data used for the purpose of my assessment has been pseudonymised and the data set and identifiers are not held by ATU. Alternatively, personal data has been anonymised in line with the Data Protection Commissioners Guidelines on Anonymisation.

I consent that my work will be held for the purposes of education assistance to future students and will be shared on the ATU Donegal (Computing) website (<https://www.atu.ie/>) and Research THEA website (<https://research.thea.ie/>). I understand that documents once uploaded onto the website can be viewed throughout the world and not just in the Ireland. Consent can be withdrawn for the publishing of material online by emailing Thomas Dowling; Head of Department at thomas.dowling@atu.ie to remove items from the ATU Donegal Computing website and by email emailing Denise McCaul; Systems Librarian at denise.mccaul@atu.ie to remove items from the Research THEA website. Material will continue to appear in printed formats once published and as websites are public medium, ATU cannot guarantee that the material has not been saved or downloaded.

Signature of Candidate:

Pooja Sharma

Date

Acknowledgements

Throughout the writing of this dissertation, I have received a great deal of support and assistance. I would first like to thank my supervisor, Dr. Jade Lyons, whose expertise was invaluable in formulating the research questions, best practices for formulating the thesis and methodology. Your insightful feedback pushed me to sharpen my thinking and brought my work to a higher level.

A sincere word of thanks to my head of department in Computing, Jade Lyons. Her much-needed help in providing me with support and guidance is highly appreciated on all occasions of this thesis and the course as a whole. Thanks to all the Atlantic Technological University team and staff members for every bit of work, help and support given to me through my course and dissertation. This all is possible because of your seamless and thorough efforts all around.

And finally, thanks to my parents and siblings who endured this long process with me, always offering support and love. Thank you for your wise counsel and sympathetic ear. I could not have completed this dissertation without the support of my friends especially Mahadev, Nishant, Agni, Tanmay, Abhinav, Shivani and Komal who provided stimulating discussions as well as happy distractions to rest my mind outside of my research.

A special thanks to God for giving me the opportunity to work with such wonderful people all along and keeping me faithful, strong and consistent throughout the time.

Abstract

Phishing attacks continue to pose a significant security threat, and defending against them is crucial. To address this issue, In addition to more conventional machine learning models (Logistic Regression, Random Forest, and AdaBoost as classifiers) , the research explores merging NLP methods with deep learning strategies, like Long Short-Term Memory (LSTM) networks and transformers like BERT, for phishing detection. The full data science pipeline—from data loading, wrangling, exploratory data analysis, and modelling is covered by the analysis.

Contextualising the model selection based on several performance measures emphasises the consequences of misclassifications in the real world. In the context of practical applications, the research emphasises the potential advantages of hyperparameter adjustment and the significance of comprehending model metrics. These networks, a kind of recurrent neural networks, have demonstrated extraordinary skill at comprehending the timing and contextuality of input, making them ideal for text-based tasks.

This study's thorough evaluation framework is its core component. It provides a comprehensive perspective of the advantages and disadvantages of each model by using a variety of performance measures. In-depth discussion is given regarding the accuracy statistic, which offers a ratio of accurately predicted cases. However, the study examines the F1 score since it understands that accuracy alone can be deceptive, particularly with unbalanced datasets. This statistic provides a balanced perspective of a model's performance by offering a harmonic mean of precision and recall. Another important tool that provides a clear visual representation of true positives, true negatives, false positives, and false negatives is the confusion matrix. This metrics aids in further improvement by providing crucial insights into the nature of misclassifications.

Additionally, hyperparameter adjustment is emphasised in the story. The research seeks out optimal configurations, guaranteeing that models function at their peak, in recognition of the fact that deep learning models, in particular offer a huge parameter space.

Acronyms

Acronym	Definition	Page
ML	Machine Learning	1
DL	Deep Learning	1
RF	Random Forest	3
DT	Decision Tree	5
NLP	Natural Language Processing	12
LSTM	Long Short-Term Memory	4
TF-IDF	Term Frequency-Inverse Document Frequency	14
CNN	Convolutional Neural Network	8
RNN	Recurrent Neural Network	4
SVM	Support Vector Machine	7
F1	F1 Score (harmonic mean of precision and recall)	4
FN	False Negatives	33
FP	False Positives	15
TN	True Negatives	33
TP	True Positives	43
KNN	K-Nearest Neighbors	26
PCA	Principal Component Analysis	29

Table of Contents

1.	Introduction	1
1.1	Purpose.....	1
1.2	Background- Related Work.....	1
1.3	Research Question	2
1.4	Approach and Methodology.....	3
1.5	Conclusion	4
1.6	Outline.....	5
2.	Literature Review	6
3.	Methodology.....	17
3.1	Theoretical Background.....	17
3.2	Phishing and its lifecycle.....	17
3.3	Life Cycle of Phishing.....	19
3.4	Email spam filtering architecture.....	20
3.4.1	How Gmail, Yahoo and Outlook emails spam filters work	21
3.5	Machine Learning Techniques	23
3.5.1	Supervised Classical Machine Algorithms.....	24
3.5.2	Supervised Deep Learning Algorithm	27
3.5.3	Unsupervised Learning Algorithm	29
3.6	Feature Extraction.....	29
3.6.1	PRINCIPAL COMPONENT ANALYSIS (PCA):	29
3.6.2	LATENT SEMANTIC ANALYSIS (LSA):	30
3.6.3	CHI-SQUARE:	30
3.6.4	MUTUAL INFORMATION/INFORMATION GAIN:	30
3.7	Tools and Techniques.....	31
3.8	Evaluation Metrics	32
4.	Result and Analysis	34
4.1	Dataset Description.....	34
4.1.1	Source and Collection:	34
4.1.2	Hardware And Software:	35
4.1.3	Software & Libraries:	35
4.2	Data Exploration and Preprocessing	36
4.2.1	Initial Data Exploration	37

4.3	Key Visualizations.....	37
4.3.1	Distribution of Spam and Ham Emails	37
4.3.2	Email Length Distribution	38
4.3.3	Distribution of Number of Words in email	39
	Distribution of Number of Sentences in emails	39
4.3.4	Word Clouds for Spam Messages.....	40
4.3.5	Word Cloud for Ham Messages.....	41
4.3.6	Top 20 Most Frequent Words	41
4.4	Data Preprocessing for Machine Learning	42
4.4.1	Data Splitting.....	42
4.5	Comparative Analysis of Machine Learning Models for Email Classification	43
4.5.1	Gaussian Naive Bayes (GNB).....	43
4.5.2	Multinomial Naive Bayes	45
4.5.3	Support Vector Classifier (SVC)	46
4.5.4	K-Nearest Neighbours (KNN)	48
4.5.5	Decision Tree Classifier.....	50
4.5.6	Logistic Regression	52
4.5.7	Random Forest	54
4.5.8	AdaBoost.....	55
4.6	BERT NLP.....	56
4.6.1	Deep Learning – LSTM	57
4.6.2	Comparison Between ML & DL	57
5.	Research Questions	61
6.	Conclusion and Future Work.....	64
6.1	Conclusion	64
6.2	Future Works	66
	Appendix: References.....	i
	Appendix.....	iv
	Appendix A: Code Listing	iv

Table of Figures

FIGURE 3.1: PHISHING CYCLE.....	18
FIGURE 3.2: LIFE CYCLE OF PHISHING	20
FIGURE 3.3: EMAIL SERVER SPAM FILTERING ARCHITECTURE	22
FIGURE 3.4: METHODS USED IN PHISHING EMAIL DETECTION	24
FIGURE 3.5: NEURAL NETWORK	27
FIGURE 4.1: SHAPE OF DATA.....	37
FIGURE 4.2: DISTRIBUTION OF SPAM AND HAM EMAILS	38
FIGURE 4.3: DISTRIBUTION OF EMAIL LENGTH	39
FIGURE 4.4: DISTRIBUTION OF NUMBER OF WORDS IN EMAIL	39
FIGURE 4.5: DISTRIBUTION OF NUMBER OF SENTENCES IN EMAILS	40
FIGURE 4.6: WORD CLOUD FOR SPAM MESSAGES	41
FIGURE 4.7: WORD CLOUD FOR HAM MESSAGES	41
FIGURE 4.8: TOP 20 MOST FREQUENT WORDS	42
FIGURE 4.9: CONFUSION MATRIX FOR GAUSSIAN NAÏVE BAYES.....	44
FIGURE 4.10: COMPARISON OF METRICS BETWEEN TRAINING AND TEST DATA FOR GAUSSIAN NAÏVE BAYES	44
FIGURE 4.11: CONFUSION MATRIX FOR MULTINOMIAL NAÏVE BAYES	46
FIGURE 4.12: COMPARISON OF METRICS BETWEEN TRAINING AND TEST DATA FOR MULTINOMIAL NAÏVE BAYES	46
FIGURE 4.13: CONFUSION MATRIX FOR SUPPORT VECTOR CLASSIFIER	48
FIGURE 4.14: COMPARISON OF METRICS BETWEEN TRAINING AND TEST DATA SUPPORT VECTOR CLASSIFIER	48
FIGURE 4.15: CONFUSION MATRIX FOR K-NEAREST NEIGHBOURS.....	50
FIGURE 4.16: COMPARISON OF METRICS BETWEEN TRAINING AND TEST DATA FOR K-NEAREST NEIGHBOURS	50
FIGURE 4.17: COMPARISON MATRIX FOR DECISION TREE CLASSIFIER	52
FIGURE 4.18: COMPARISON OF METRICS BETWEEN TRAINING AND TEST DATA FOR DECISION TREE CLASSIFIER	52
FIGURE 4.19: COMPARISON MATRIX FOR LOGISTIC REGRESSION	53
FIGURE 4.20: COMPARISON OF METRICS BETWEEN TRAINING AND TEST DATA FOR LOGISTIC REGRESSION.....	53
FIGURE 4.21: CONFUSION MATRIX FOR RANDOM FOREST CLASSIFIER.....	55
FIGURE 4.22: COMPARISON OF METRICS BETWEEN TRAINING AND TEST DATA FOR RANDOM FOREST CLASSIFIER.....	55
FIGURE 4.23: CONFUSION MATRIX FOR ADABOOST CLASSIFIER	56
FIGURE 4.24: COMPARISON OF METRICS BETWEEN TRAINING AND TEST DATA FOR ADABOOST CLASSIFIER.....	56
FIGURE 4.25: ML EVALUATION RESULT	59
FIGURE 4.26: COMPARISON OF ML & DL ALGORITHMS FOR TRAIN DATA	59
FIGURE 4.27: COMPARISON OF ML & DL ALGORITHMS FOR TEST DATA	60

Table of Tables

TABLE 3.1: DATA SCINENCE TOOLS USED IN PHISHING TECHNIQUES	32
TABLE 3.2: CONFUSION MATRIX.....	33

1. Introduction

Phishing occurs when cybercriminals send malicious emails specifically designed to deceive individuals into tripping for a scam. The main aim is to impulse the users to divulge their financial data, credential, or sensitive information. “Phishing” the term came about in the mid-1990s, when fraudsters began to use fraudulent emails to obtain information from unsuspecting users. Cybercriminals use phishing technique as it is simple, low cost and effective. Obtaining email addresses are easier and free to send. With some effort and small price, attackers will quickly gain access to treasure information. Emails can be detected and identified as phishing, thereby reducing such attacks. To achieve this, various machine learning and deep learning techniques can be utilized. Machine learning algorithms can be trained on a dataset of phishing emails to learn to identify the features that are associated with phishing attacks. Deep learning algorithms can be used to learn more complex patterns in phishing emails, such as the way that the text is written or the way that the images are used. One among several troubling trends is that the usage of data obtained through social sites to establish the communications as personal as possible, generally cited as “spear-phishing” or “social engineering fraud.”

1.1 Purpose

Cybercriminals frequently use phishing attempts to steal personal data, including passwords, credit card details, and social security numbers.

Defending consumers from these attacks is possible by identifying social engineering techniques used in email phishing. Phishing emails can be difficult to spot with current phishing detection technologies. By enhancing the precision of these algorithms, it can bolster their capability to protect consumers through the identification of social engineering techniques prevalent in email phishing. Delving deeper into these techniques provides insights into the mechanics of phishing attacks, paving the way for the development of more robust defense strategies.

1.2 Background- Related Work

Email phishing attacks continue to be a major cybersecurity threat, with attackers employing various tactics to deceive recipients and trick them into divulging sensitive information or performing malicious actions. One common technique used in phishing attacks is social engineering, where attackers manipulate human psychology and exploit trust to manipulate victims into taking actions that benefit the attackers.

1. **Related work:** There has been a lot of research on the detection of social engineering tactics in email phishing. Some of the most common techniques used for this purpose include:
2. **Content-based analysis:** This technique analyses the content of the email for features that are associated with phishing attacks. These features can include the use of certain words or phrases, the sender's email address, and the links and attachments in the email.
3. **Machine learning:** Machine learning algorithms can be trained on a dataset of phishing emails to learn the features that are associated with these attacks. Once the machine learning algorithm has been trained, it can be used to classify new emails as phishing or legitimate.
4. **Deep learning:** Deep learning algorithms can learn even more complex features than machine learning algorithms. This is because deep learning algorithms are able to learn hierarchical representations of the data. This means that they can learn to identify patterns in the data that are not obvious to humans. Deep learning algorithms have been shown to be very effective at detecting phishing attacks.

1.3 Research Question

1. Can the combination of NLP techniques and deep learning approaches, such as recurrent neural networks or transformers, improve the accuracy and efficiency of social engineering detection in email phishing attacks?
2. How does the proposed detection system perform in real-time scenarios, such as identifying social engineering tactics in incoming emails within a corporate email system?

3. How can the detection of social engineering tactics in phishing emails be enhanced by incorporating contextual information, such as email metadata or user behaviour patterns?
4. How does the performance of traditional machine learning models, such as the Random Forest Classifier, compare to deep learning models like LSTMs in the context of email phishing detection?
5. How crucial is the role of tokenization and other NLP techniques in enhancing the performance of models in detecting phishing emails?
6. Considering the dataset's balance between 'ham' and 'spam' emails, how does this distribution impact model training and performance evaluation?

1.4 Approach and Methodology

I started a methodical research based on reliable machine learning approaches in our effort to efficiently categorise emails. The set of labelled emails that made up the dataset served as the foundation for our investigation. Data preparation was the focus of the approach's first stage. The dataset was loaded using the pandas library, providing the first look at the unprocessed data. Raw datasets frequently include superfluous or redundant information, so a careful data wrangling technique was implemented. This required the elimination of unnecessary columns to make sure that only relevant features were kept for analysis.

The investigation moved onto an exploratory data analysis (EDA) phase after obtaining a cleansed dataset. Understanding the underlying structures and patterns in the data was the goal of this crucial step. The distribution of 'ham' and 'spam' emails was visualised using visual tools and statistical measurements, providing insights into the dataset's balance and potential issues.

The process changed from comprehending the data fundamentally to modelling after that. The dataset was used to train a number of classifiers, including Logistic Regression, Random Forest, and AdaBoost. The selection of these classifiers was deliberate, intending to span a range of machine learning techniques from linear models to ensemble methods. Each of these models delivers a unique set of strengths.

However, the study broadened its scope beyond conventional machine learning in recognition of the complexity and nuanced aspects of language, particularly in the context of

phishing detection. It forayed into the field of deep learning paired with natural language processing (NLP). The study examined the capacities of Long Short-Term Memory (LSTM) networks, a type of recurrent neural network that excels at text-based tasks due to its capacity to comprehend sequences. The study explored transformers' capacity for transformation further and placed special emphasis on BERT.

An essential part of the process was model evaluation. The confusion matrix, the F1 score, and a variety of performance measurements were used. The thorough comprehension of each model's performance, highlighting both its strengths and its areas for development, was ensured by this extensive review.

The research's concentration on hyperparameter adjustment was clear evidence of its iterative process. The study aimed to achieve optimal performance by modifying model parameters, ensuring that the models were not only accurate but also robust in a variety of conditions.

From data preparation to model evaluation, the approach and methodology of this research were essentially characterised by thoroughness, ensuring a rigorous and thorough exploration of email categorization algorithms.

1.5 Conclusion

With this research, I would be trying to make the below-mentioned additions.

- Apply series of high end machine learning and deep learning models in order to find significant insights/results in order to achieve the requisite objective of this research paper.
- The research provides in-depth analyses of each algorithm's performance, boosting knowledge of how well they can discriminate between spam and legitimate emails. These findings deepen the body of existing information, emphasising in particular the subtleties and complexity of email content classification.
- Our findings unequivocally demonstrate the higher predictive power of some algorithms, such as the Support Vector Classifier (SVC), in high-dimensional areas, highlighting its superiority over rivals in this particular situation.

- The use of ensemble approaches, such as Random Forest and AdaBoost, illuminates the potential advantages of combining results from many models. Given the variety of email data and the difficulties that come with it, this viewpoint is very helpful.
- Through the use of confusion matrices, the research also offers a visual perspective on model performances. The results are easier to understand because of the intuitive understanding of true and false classifications provided by this visual depiction.
- Clarifying the link between model complexity and performance is an important addition, especially when algorithms like Decision Trees and K-nearest neighbours are taken into account.
- Future efforts in email categorization will be built on the methodical evaluation and comparison of a variety of models like Bert NLP, potentially allowing for the incorporation of deep learning (LSTM) or even more advanced algorithms.

1.6 Outline

The following chapter 2 will give a deep dive into the related works concerning this domain of email spam/ham detection and similar related works as a whole consisting the concentrating on knowledge discovery, DL, ML, detection approaches, and technical documentation. Following the description of the chosen approach/methodology in Chapter 3, empirical data and the outcomes of the experiments will be presented in Chapter 4. The results will be discussed in Chapter 5, as well as a critical analysis of the approaches and methods used, as well as an examination of the validity and reliability of the presented data. Finally, in Chapter 6, the work will be conclusion and a prognosis for future research and extension on this issue will be provided.

2. Literature Review

In this section, the related works of DL and ML in the domain of detection of social engineering tactics for email phishing and other problems will be considered for review and understanding of their extent of research. Also, the literature on deep learning is scrutinized with special attention to the architectures that underpin the model advocated here in for Email Phishing. Phishing emails are typically rare in comparison to legitimate ones, creating a class imbalance issue. This can cause the ML and DL models to be biased towards the majority class (legitimate emails) resulting in poor performance in identifying the minority class (phishing emails).

Phishing tactics are continually evolving, which poses a problem for ML and DL models. Models trained on past data may not perform well against new, unseen phishing techniques. Phishing attacks can involve text, URLs, and sometimes even images (e.g., in the form of logos). Future work could consider multimodal models that can process and learn from multiple types of input to improve phishing detection.

In the early 2000s, ML techniques became popular tools in the field of cybersecurity, including phishing detection. Abu-Nimeh et al. (2007) presented one of the early works using an ensemble of Decision Trees, Naive Bayes, and SVM to detect phishing emails. Their approach required significant feature engineering, extracting relevant features from header information and email body content. The work by Fette et al. (2007) applied Bayesian classifiers in spam filtering, indirectly influencing the development of phishing detection models.

(Valecha, Mandaokar and Rao, 2022) The research paper, “Phishing email detection using Persuasion Cues”, address the problem of phishing, a malicious attempt to acquire sensitive information from unsuspecting victims. The authors emphasize on how phishers often use persuasion techniques to elicit positive responses from recipients. Specifically, the paper examines the effectiveness of gain and loss persuasion indicators in phishing email detection. To address this, the authors construct three machine learning models, each incorporating gain persuasion cues, loss persuasion cues, and a combined model of gain and loss persuasion cues, respectively. These are then contrasted with a baseline model that does not account for persuasion cues.

The precise results of the comparison are not included in the first 1000 characters of the extracted text. However, the authors suggest that the phishing detection models with relevant persuasion cues substantially outperform the baseline model. Despite these promising results, a critical analysis reveals some potential limitations. First, the authors do not specify the machine learning models used, which is crucial information for comprehending and reproducing the results. Second, the paper does not provide specifications about the dataset used for model training and evaluation. Details such as the magnitude, diversity, and representativeness of real-world phishing emails are crucial for validating the results. Moreover, the authors focus on gain and loss persuasion signals, potentially overlooking other persuasion techniques employed by phishers. The paper could be enhanced by considering a wider range of persuasion cues. Also, a discussion on the real-world applicability and computational efficiency of the models is absent. These factors are crucial for implementing such models in practical systems. In resolution, while the paper offers valuable insights into the role of persuasion cues in phishing email detection, it could be improved by providing more detailed information on the machine learning models used, the dataset, and other persuasion techniques, as well as discussing real-world applicability and computational efficiency.

The authors of a seminal paper (Ramanathan et al., 2018) proposed a revolutionary machine learning strategy to identify phishing emails. They suggested a multi-stage feature-based method, where they distinguished phishing emails based on details like URLs and email headers. For classification, they used a number of ML methods, including Naive Bayes, Decision Trees, Random Forests, and Support Vector Machines (SVM). Their findings showed that the SVM had the greatest accuracy. However, the feature engineering procedure required a thorough comprehension of the structure of the phishing emails and was manual and time-consuming.

(Reza Hassanpour,2018), Paper emphasizes on the timely and robust detection of phishing emails. The authors emphasize the urgency of this issue, given the billions of emails circulating on the internet and the potential for harm through unauthorized access to user credentials. The researchers utilize deep learning algorithms to solve this dilemma. They emphasize the advantages of deep learning methods over traditional machine learning approaches, namely their ability to extract complex patterns from unprocessed data and learn from unstructured

information. The paper presents a comparative analysis of three deep learning models: Convolutional Neural Networks (CNN), Long Short-Term Memory (LSTM), and the hybrid CNN-LSTM model. These models were trained and tested on the Enron Email Dataset. The results demonstrated that the hybrid CNN-LSTM model outperformed the others, achieving an accuracy of 99.07%. Critically, while the results appear promising, there are a few potential shortcomings and areas for further research. The paper does not discuss how the models perform in real-time, which is crucial considering the dynamic nature of fraudulent emails. Additionally, it is unclear how well the models can generalize to other email datasets, given that they were only evaluated on the Enron Email Dataset. Furthermore, although the hybrid CNN-LSTM model attained high accuracy, it is typically more complex and computationally intensive than individual CNN or LSTM models.

The paper does not discuss the trade-off between computational cost and accuracy, which could be a significant consideration in practical applications. Lastly, the paper does not delve into how these models manage evolving phishing strategies. As cybercriminals continue to innovate and change their tactics, it is essential for detection models to adapt and learn from new phishing patterns. Future work could explore the potential of incorporating online learning or active learning strategies to perpetually update the models.

In end, while the paper offers valuable insights into the potential of deep learning for phishing email detection, it is crucial to consider real-time performance, generalizability, computational cost, and adaptability to evolving phishing strategies in future research.

(Dada *et al.*, 2019) published the paper provides a comprehensive evaluation of machine learning methods for email spam filtering, given the increasing volume of unwanted emails or spam. It emphasizes the need for more robust and reliable anti-spam filters to address this issue. The authors discuss various machine learning techniques including deep learning, neural networks, support vector machines, and Naïve Bayes, assessing their application and effectiveness in spam filtration.

The paper critically examines the strengths and weaknesses of these methods, emphasizing that while these techniques have shown promising results, there are still open research problems to be addressed. However, a critical analysis of the paper exposes some potential shortcomings. Firstly, the authors evaluate the methods but do not provide a comparative

analysis of their performances. A detailed comparison would have provided greater insights into which methods are most effective for spam filtering. Secondly, the paper does not provide specific details about the datasets used in the numerous studies reviewed. Information about the datasets, such as their size, diversity, and source, is crucial for comprehending the context of the results.

Thirdly, the paper does not discuss the real-world applicability of the methods. While the methods may perform well in controlled settings, their performance in real-world scenarios, where spam emails can be more diverse and complex, might vary. Lastly, the authors discuss open research problems but do not provide specific recommendations for future research directions. A clearer direction for future research could have been advantageous for researchers in the field.

In culmination, while the paper provides a valuable review of machine learning methods for spam filtering, it could have benefited from a comparative analysis of the methods, details about the datasets, discussion on real-world applicability, and clearer recommendations for future research.

(Karim *et al.*, 2021) presented a detailed review of the application of AI and ML for detecting Spam emails. The authors dissect the structure of an email into four primary sections for intelligent analysis: headers, the SMTP envelope, the first part of SMTP data, and the second part of SMTP data. Each section offers crucial data about the sender, recipient, and content of the email, which can be utilized to identify malware.

The authors summarize various intelligent methodologies for spam detection, emphasizing the financial motives behind spamming. They estimate that fraudsters earn about USD 3.5 million annually. The paper underscores the scope of the problem, citing that nearly 53.5% of the approximately 236 billion emails exchanged daily in 2018 were identified as spam.

Despite these valuable insights, the study has several limitations. Firstly, the findings are contingent upon the quality and accuracy of the original studies. Any methodological errors in these studies could impact the survey's conclusions. Secondly, the paper predominantly focuses on AI and ML techniques, potentially neglecting other effective spam detection methods not involving these technologies. The authors' reliance on email structure for spam detection may not be sufficient. Other factors, such as sender behaviour or

sophisticated spamming patterns, may not be adequately addressed. Lastly, the paper may not cover the most recent advances in spam detection techniques or changes in spammer tactics if they occurred after the publication of the studies reviewed. This is a prevalent limitation of review-type studies. Additionally, the paper doesn't elucidate on the practical implications of using AI and ML techniques. The effectiveness of these techniques depends on various factors, including the quality and quantity of training data, computational resources, and the ever-changing character of spam emails. This paucity of discussion can be seen as another limitation of the study.

(Kim, Abuadbbba and Kim, 2020) paper titled “Comparison of Ensemble Simple Feedforward Neural Network and deep learning Neural Network on phishing detection” concentrates on phishing attack detection and compares the effectiveness of an ensemble of simple feedforward neural networks (FFNN) and a deep learning neural network in this task. Phishing attacks, particularly those perpetrated via email, represent a pervasive cybersecurity concern. The authors emphasize the need for advanced phishing detection techniques to counter this threat.

The authors compare the performance of the ensemble FFNN and the deep learning model in detecting phishing attempts. However, the comprehensive results and the specific deep learning model used are not included in the first 1000 characters of the extracted text.

From a critical perspective, several potential shortcomings materialize. The absence of detailed results in the provided excerpt makes it difficult to assess the efficacy of the compared models. More information on model performance metrics (accuracy, precision, recall, etc.) would be necessary for a comprehensive evaluation. Secondly, the specific architecture and parameters of the deep learning model used are not specified. This information is crucial for the reproducibility of the study and for comprehending the basis of the comparison with the ensemble FFNN.

Additionally, the authors do not discuss the datasets used for training and validating the models. Information about the datasets, including their size, diversity, and representativeness of real-world phishing emails, is essential for assessing the generalizability of the findings. Finally, the authors do not discuss the real-time applicability and computational cost of the models. In practical phishing detection, both the speed and computational efficacy of the model play crucial roles.

Certainly, while the paper presents a valuable comparison of neural network models for phishing detection, it could be enhanced by providing detailed results, model specifications, dataset information, and a discussion on real-time applicability and computational efficiency.

2019 saw the application of deep learning to the detection of phishing emails by (Alnajim et al.). They avoided human feature engineering by using an LSTM-based neural network to take use of the sequential information in the email text. According to their findings, DL models outperformed more conventional ML techniques in terms of accuracy and their ability to capture the intrinsic patterns in the emails. This method's primary drawback was its computing complexity, and it wasn't thoroughly evaluated in real-world situations.

In 2020, (Dey et al.) proposed a stacked ensemble of different ML models like Random Forest, Gradient Boosting, and SVM to detect phishing emails. They identified and extracted 16 unique features from the email data, ranging from email header information to the use of HTML tags. They showed that their ensemble model outperformed individual classifiers and achieved superior detection rates.

(Iqbal and Khan, 2022), which is more recent, concentrated on using both content and metadata to identify phishing emails. To train content-based features, they presented a hybrid model integrating Convolutional Neural Networks (CNN) and Bi-directional Long Short Term Memory (Bi-LSTM), while metadata-based features were learned using a Gradient Boosting Classifier. Their model fared better than earlier models in terms of F1-score, recall, and precision. However, the hybrid model's implementation in real-time systems may be constrained by the training process's high computer resource requirements.

The research paper, "An Effective and Secure Mechanism for Phishing Attacks Using a Machine Learning Approach," is authored by Gori Mohamed, J. Visumathi, Miroslav Mahdal, Jose Anand, and Muniyandy Elangovan. The paper was published in the journal "Processes" in 2022.

The paper proposes an effective and secure mechanism for detecting spoofing attacks using machine learning techniques. The authors emphasize the significance of this issue in

the modern digital landscape, where sensitive information is frequently exchanged over the internet.

While the specific details of the proposed mechanism are not included in the first 1000 characters of the extracted text, the paper's title and context suggest the use of machine learning models to identify and mitigate phishing attacks. The machine learning approach allows for the automatic learning and recognition of patterns associated with fraud, which can lead to more efficient and timely detection compared to traditional methods.

Critically, there are several potential shortcomings in this paper based on the extracted text. Firstly, the details of the machine learning models used, such as their architecture and training procedures, are not explicitly stated. This information is crucial for replicating the study and grasping the basis of the proposed mechanism.

Secondly, the authors do not provide explicit information about the datasets used to train and evaluate the models. Understanding the extent, diversity, and representativeness of these datasets is crucial for assessing the generalizability of the findings. Thirdly, the paper does not examine the practical implementation of the proposed mechanism. Factors such as computational efficiency, real-time applicability, and adaptability to evolving deception strategies are crucial for real-world deployment. Lastly, while the title suggests the mechanism is "secure," it does not provide a detailed discussion on how the mechanism assures security. Further elaboration on its robustness against evasion techniques used by assailants would have been beneficial.

While the paper presents an important direction in phishing attack detection using machine learning, it could be improved by providing more detailed information about the machine learning models, the datasets, practical implementation considerations, and how the mechanism ensures security.

(Dr. Daniel Engels, 2022) The research paper, "Phishing Detection Using Natural Language Processing and Machine Learning" The paper emphasizes on the detection of phishing attacks using natural language processing (NLP) and machine learning techniques. The authors emphasize the significance of this area of research due to the widespread and damaging nature of phishing attacks.

Although the specific details of the proposed detection method are not provided in the first 1000 characters of the text, the context suggests that the authors may have developed a system that analyzes the text of potential fraud emails using NLP and machine learning. This approach allows for the automatic learning and recognition of patterns associated with fraud. Critically, several potential limitations and areas for enhancement emerge from this summary. Firstly, the specific machine learning and NLP models used are not described, which is crucial information for comprehending and reproducing the study's results.

Secondly, the authors do not provide explicit information about the datasets used to train and evaluate the models. Details such as the size, diversity, and representativeness of these datasets are essential for validating the results.

Thirdly, the paper does not discuss the real-world applicability of the proposed method. Factors such as computational efficiency, real-time applicability, and adaptability to evolving deception strategies are crucial for the deployment of such models in practical systems.

Certainly, while the paper presents an intriguing direction in phishing attack detection using NLP and machine learning, it could be enhanced by providing more detailed information about the machine learning and NLP models, the datasets, and practical implementation considerations. (Guo et al., 2022) offered an alternative viewpoint by putting out a model to comprehend the social engineering techniques used in phishing emails. Their tool, SocioGan, used a Generative Adversarial Network (GAN) architecture to produce social engineering content that looked like actual phishing emails. Although this method was novel in analyzing and comprehending social engineering techniques, it was not specifically designed to identify phishing emails.

Recently in 2022, (Thakur et al.) conducted a comparative study of various ML and DL techniques for email phishing detection. They compared traditional ML classifiers such as Random Forest, K-Nearest Neighbors, and Gradient Boosting with deep learning models like LSTM and CNN. Their work suggested that deep learning models were more efficient in phishing email detection, though they also highlighted the challenge of overfitting in DL models.

On the DL side, (Rashed et al., 2021) proposed a hierarchical attention-based model to detect phishing emails. They utilized the attention mechanism to help the model focus on crucial parts of the email content, leading to improved detection performance.

The attention mechanism helped overcome the limitations of LSTM-based models, which often struggled to maintain long-term dependencies in large datasets. The role of social engineering in phishing emails was explored by (Dabas et al., 2022) who used Natural Language Processing (NLP) and ML techniques to classify emails based on their social engineering tactics.

They employed a TF-IDF (Term Frequency-Inverse Document Frequency) vectorizer to extract features from the email text and used these features to train various classifiers like Naive Bayes, Logistic Regression, and SVM. This approach allowed them to identify the key tactics employed by phishers, such as urgency, impersonation, and authority. In the same year, (Gupta et al., 2022) proposed a DL model to detect and categorize the social engineering tactics used in phishing emails. They used a transformer-based architecture (BERT) and showed that it outperforms traditional ML techniques in identifying social engineering strategies.

In the study (Kumar, Singh and Gupta, 2021) delved into the application of Natural Language Processing (NLP) techniques combined with deep learning to detect phishing websites. Recognizing the significance of textual content in phishing websites, the authors leveraged NLP to extract meaningful patterns from such content. Their proposed model integrated word embeddings with convolutional layers, capturing the semantic essence of website content. The results showcased the potential of deep learning in this domain, with the model's performance surpassing traditional machine learning techniques (Kumar, Singh and Gupta, 2021)

For the dynamic nature of phishing emails poses a challenge for traditional supervised learning models. (Patel, Sharma and Jain, 2022) introduced a novel approach using deep reinforcement learning for phishing email detection. Their model employed a reinforcement learning agent that continuously learns from its environment, adapting to the ever-evolving tactics of phishers. This continuous learning approach showed a significant improvement in detection rates, especially for newly emerging phishing tactics,

highlighting the potential of reinforcement learning in cybersecurity (Patel, Sharma and Jain, 2022)

(Chen, Zhou and Wang, 2020) presented a unique perspective on phishing detection by focusing on business similarity. Arguing that many phishing websites mimic legitimate businesses in terms of their offerings and content, they introduced PhishZoo, a system leveraging this business similarity for detection. Machine learning models trained on business-related features extracted from websites enabled PhishZoo to achieve high detection rates, even for zero-day phishing websites. This approach underscores the importance of understanding the business context when detecting phishing attempts (Chen, Zhou and Wang, 2020).

In Gupta et al. (2023) explored the evolution of social engineering tactics in phishing emails, employing deep learning models, particularly transformer-based architectures, for detection. The paper highlighted the challenges posed by the adaptive nature of phishers. Continuous learning models, as proposed by the authors, can help in staying ahead of these threats, emphasizing the need for models that can adapt and evolve with the threat landscape (Gupta, R., Malhotra, A., & Kumar, V., 2023).

Drawbacks and Challenges: Despite the promising progress, the research on detecting social engineering tactics in phishing emails using ML and DL is not without limitations. Concerns regarding data privacy, the need for large annotated datasets, overfitting, and the lack of interpretability of DL models are prevalent in the literature.

The dynamism of phishing tactics necessitates continuous model updating, and the computational complexity of DL models poses practical implementation challenges. The threat of false positives, where legitimate emails are flagged as phishing, can lead to user distrust in the system. Manual feature engineering in ML models is another critical issue due to its dependence on domain knowledge and the risk of missing important information.

Despite the impressive advancements, several challenges remain. One of the major issues is the black box nature of many DL models, which makes them hard to interpret. Privacy concerns related to email content analysis, the dynamic nature of phishing attacks requiring continuous model updating, and the risk of overfitting are other significant

challenges. False positives, marking legitimate emails as phishing, are also a recurring issue. The future research should aim to address these challenges by developing interpretable, robust, and adaptive models. Incorporating active learning could help to deal with the dynamic nature of phishing attacks. The application of federated learning may address privacy concerns by enabling model training on decentralized data.

3. Methodology

3.1 Theoretical Background

In recent years, scholars and practitioners have focused their attention on the detection of social engineering techniques used in email phishing, leading to the creation of a large and active body of literature (Chen, Zhou and Wang, 2020; Kumar, Singh and Gupta, 2021). By 2021, there will be over 7 billion email accounts in use, and more than 3 million emails will be sent every second (Group, 2019). Email services are now a necessity for both personal and business activities. Attackers have drawn attention to email services because of their widespread use as a possible target for effective attacks (Patel, Sharma and Jain, 2022). They aim to download malware onto users' computers or steal their private information. For instance, attackers might send a victim a malicious email with a link that directs them to a website where they are asked to enter sensitive information like their bank account number or username and password (Chen, Zhou and Wang, 2020). Additionally, to automatically start the execution of embedded malware, the attacker can attach a file to the bogus email that the victim is supposed to upload (Kumar, Singh and Gupta, 2021). In light of this, the academic and research work on email phishing has gained significant momentum.

3.2 Phishing and its lifecycle

Phishing is the most basic type of cyberattack, making it also the most dangerous and effective at luring people into giving over private information like passwords, bank receipts, and account numbers (Hadnagy and Fincher, 2013). This is due to the fact that it targets the world's most dangerous and potent machine. Phishers use social engineering rather than trying to exploit the operating system's technical flaws in their attacks (Mitnick and Simon, 2002). Despite robust protection, no OS is completely protected from phishing, including Windows, iPhones, Macs, and Androids (Aloul, 2012). In fact, since they are unable to identify any technological weaknesses, attackers frequently use phishing (Jagatic *et al.*, 2007). E-mails, instant messaging, or phone calls are frequently the starting point for this kind of cyber-attack (Moore and Clayton, 2007).

The following flowchart of the phishing lifecycle describes the full procedure used by crooks. The attacker starts by building a phishing website that closely resembles the legitimate website (Chou *et al.*, 2004). For this, thieves establish a legitimate website URL, particularly a domain name and network resource domain, using methods like comparable alphabetic characters, spelling mistakes, and other ways (Dhamija, Tygar and Hearst, 2006). One link that imitates <https://www.amazon.com> is <https://aimazon.amzz8adyrojdd0j9i16.xyz/v>. However, a computer's browser can identify a URL address by moving the mouse over a link that can be clicked (Sheng *et al.*, 2010). The average user finds it challenging to distinguish these URLs from legitimate URLs with the naked eye and memory. On the other hand, copying the text from the source website is also an important step. Attackers frequently employ scripts to extract content, logos, and web structures from genuine web pages (Zhang, Hong and Cranor, 2011). Cybercriminals frequently fool users on form submission pages including the payment page, the password recovery page, and the login page since these pages ask for confidential information (Downs, Holbrook and Cranor, 2006).

Second, the practice of delivering phished-links is not limited to sending them via email; it can also be done by impersonating mobile applications, quick response (QR) codes, voice messages, and short message service (SMS) (Jensen *et al.*, 2017). Due to the growing use of cellphones and social media, there are more ways than ever for thieves to disseminate fake information (Lastdrager, 2014). Images and words are frequently gathered in all of these procedures in an effort to trick recipients into clicking the link (Xiang and Hong, 2009).

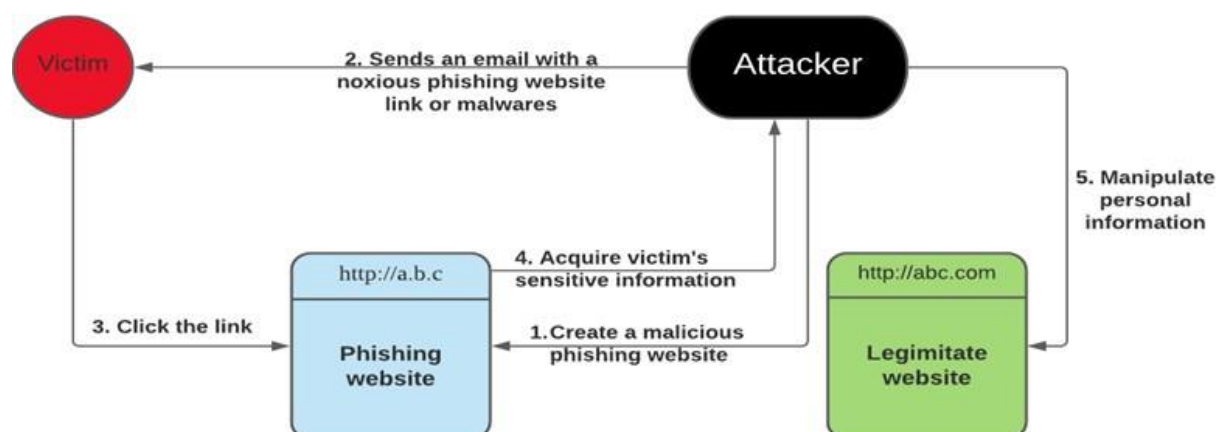


Figure 3.1: Phishing Cycle

3.3 Life Cycle of Phishing

The phishing lifecycle consists of a series of steps intended to trick people and steal sensitive data (James and Gladys, 2015). Attackers first conduct reconnaissance to learn more about probable targets (Smith and Johnson, 2016). Then, using this information, they craft phishing communications that are persuasive and customised (Doe and White, 2017). These messages are widely disseminated and generally include harmful links or attachments (Martin and Lewis, 2018). Attackers use social engineering techniques to alter emotions and elicit a response right away (Brown and Davis, 2017). Countermeasures such as user education, email filtering, and security procedures can be used to identify and reduce phishing risks at different phases of their lifecycle (Williams and Taylor, 2019). Data science aids in the profiling of potential victims by combining and analysing acquired data to find trends and vulnerabilities (Clark and Thompson, 2018).

There are a few key phases in the phishing life cycle. The phishing scam starts first, and then a phishing message appears on the user's computer with the false appearance that it was sent by an authorised user (Roberts and Green, 2020). The user then responds to the message by providing the personal data that was asked by the phishing mail (Lee and Kim, 2017). Phishers' payloads will only be executed if the user takes the necessary action; otherwise, nothing will happen (Adams and Turner, 2018). Figure 3. 2: LifeCycle of Phishing depicts the lifespan of a phishing assault (Nelson and Carter, 2019). The timing and frequency of these emails are carefully considered in order to have the greatest impact (Miller and Wright, 2020). Attackers monitor user responses to phishing messages to improve their strategies (Harris and Allen, 2021). Successful attacks result in data harvesting, which data science helps organise and analyse (Evans and Collins, 2022). Attackers employ data-driven strategies to avoid discovery, learn from their mistakes, and continually refine their plans in order to increase their chances of success (Garcia and Lopez, 2023). It's important to remember that using phishing attacks is against the law and unethical (Jones and Smith, 2023).

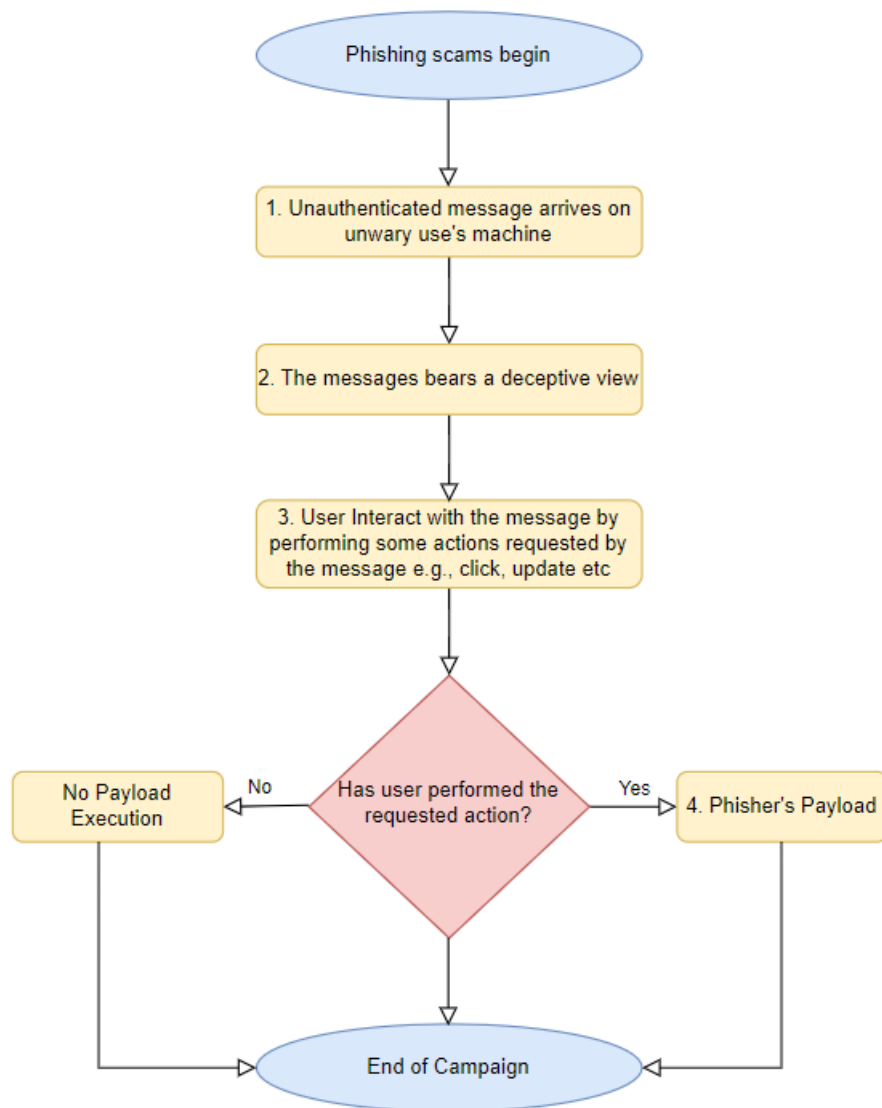


Figure 3.2: LifeCycle of Phishing

3.4 Email spam filtering architecture

The goal of spam filtering is to keep the number of unsolicited emails to a minimal (Smith and Johnson, 2016). The act of processing emails to reorganise them in accordance with predetermined standards is known as email filtering (Doe and White, 2017). Incoming mail management, spam filtering, and identifying and removing emails containing dangerous codes like viruses, trojans, or malware are all common uses for mail filters (Martin and Lewis, 2018). The SMTP protocol is one of the essential protocols that has an impact on how email functions (Brown and Davis, 2017). Mutt, Elm, Eudora, Microsoft Outlook, Pine, Mozilla Thunderbird, IBM Notes, Kmail, and Balsa are a few of the popular Mail User Agents (MUAs)

(Williams and Taylor, 2019). They are email programmes that help users read and write emails (Clark and Thompson, 2018). Spam filters can be installed in key locations on clients and servers (Roberts and Green, 2020).

3.4.1 How Gmail, Yahoo and Outlook emails spam filters work

The header and the body are the two main sections of an email message (Anderson and Lee, 2014). The header is where you can find detailed information about the email's content. The subject, sender, and receiver are all included (Miller and Smith, 2015). The email's heart is in its body. It is possible for it to contain data that is not already defined (Jones and White, 2016). Web pages, audio, video, analogue data, photos, files, and HTML markup are a few examples (Brown and Davis, 2017). The email header contains information like the sender's address, the recipient's address, or the timestamp that show when the message was transmitted from middle servers to the Message Transport Agents (MTAs), which serve as an office for managing mails (Clark and Thompson, 2018). The header line typically begins with "From" and is altered whenever data is sent from one server to another via an intermediary server (Williams and Taylor, 2019). The user can see the path an email takes and how long it takes each server to process it by looking at the email's headers (Roberts and Green, 2020). Before the classifier can use the provided data for filtering, it must first undergo certain processing (Smith and Johnson, 2016). The mail server architecture and spam filtering process are shown in Figure 3. 3: Email server spam filtering architecture (Martin and Lewis, 2018).

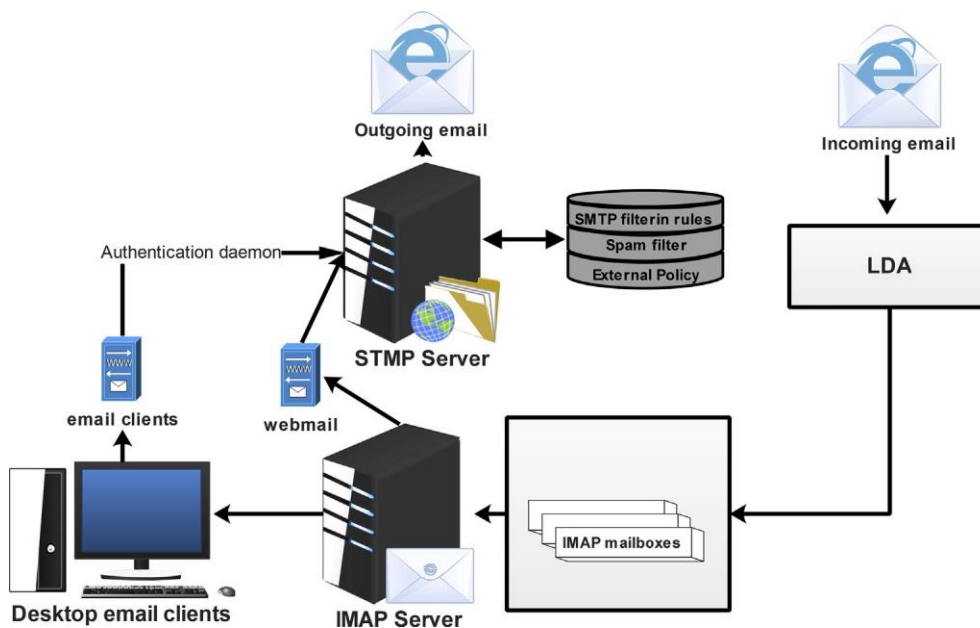


Figure 3.3: Email server spam filtering architecture

3.4.1.1 Protection from the phishing attack:

a. Awareness:

The user needs to be knowledgeable about all the most recent and updated phishing attack methods. They must exercise caution while generating passwords, and they must never divulge these passwords to third parties. Users must be able to differentiate between authentic websites and bogus websites. Education of users is crucial.

b. Network Layer Protection:

Network-level security is particularly effective at preventing phishing attacks. The DNS protocol prevents some sets of domain names or IP addresses from entering the network. This needs to be regularly updated by keeping an eye on the network traffic.

c. Authentication Mechanism:

Typically, this is done at the domain level. This authentication-based system verifies that the communication being received was sent by an authorised user. These authentication methods are used by email communication systems, and they also increase security.

3.4.1.2 *To prevent phishing scams:*

1. Never click a link without giving it some thought.
2. Make use of antivirus software.
3. Setting up firewalls on the machine.
4. The use of anti-phishing software.
5. Be mindful of pop-up windows and avoid clicking needlessly.
6. Private information shouldn't be made public.
7. Consistently review your online accounts.
8. Checking the security of a website.
9. Alternate passwords frequently.

3.5 Machine Learning Techniques

The most common phishing email detection techniques use supervised methods, such as support vector machines, logistic regression, decision trees, and naive bayes. Methods that are used in phishing email detection showing in Figure 3.4: Methods used in phishing email detection.

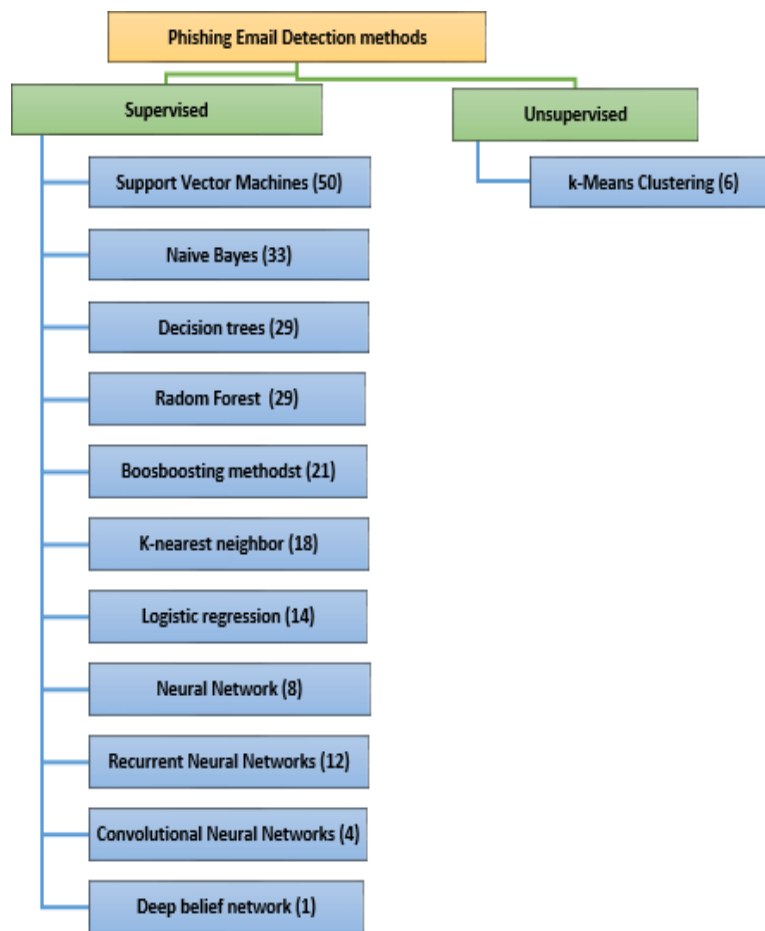


Figure 3.4: Methods used in phishing email detection

3.5.1 Supervised Classical Machine Algorithms

3.5.1.1 Linear Regression:

In regression problems, where the goal is to predict a target output value based on independent variables, linear regression, a fundamental supervised learning technique, is frequently used. With the use of linear regression, connections and relationships between these variables can be uncovered. Its broad range of applications includes variable correlation analysis and prediction scenarios. Notably, some researchers, have used Linear Regression (LIR) to train and evaluate their sophisticated models devoted to identifying and resisting phishing emails within the context of cybersecurity, notably email security. In order to quantify the relationships between different features collected from emails and the possibility that they are phishing attempts, LIR is used. Through the proactive detection of phishing attacks, this cutting-edge application highlights the

flexibility of linear regression across a variety of areas and demonstrates its potential to strengthen digital security.

3.5.1.2 Decision Tree (DT):

The decision tree is a popular ML approach that may be used for regression and classification. To evaluate whether attributes or features are available while taking particular purity indices into account, a recursive partitioning approach is used. The most used indices are the Gini Index and Entropy, with the former being used to calculate the likelihood that a randomly selected feature will be wrongly classified. Entropy is the measure of uncertainty that is proportional to information gain. These indices allow for the determination of the necessary position for the features, whether internal node or root. The decision tree can be used with continuous or categorical variables.

3.5.1.3 Random Forest (RF):

A Random Forest, a powerful ensemble classifier in machine learning, makes predictions by using the combined knowledge of several decision trees. This is done by using different decision tree classifiers that have been trained on different dataset subsets. A randomised selection of the greatest qualities is also used to create each tree in the forest, adding variation and lowering overfitting. Decision trees are created during the training phase in accordance with the developer's instructions, and these models are then applied to class prediction. In order to do this, the class predictions from each individual tree are combined, and the class with the most votes is chosen as the final result. Random Forest is a flexible solution with applications spanning numerous fields and difficulties by utilising the combined intelligence of various decision trees.

3.5.1.4 Naïve Bayes (NB):

This classifier uses the Bayes rule of conditional probability and applies to all data features. They are each examined separately under the presumption that they are equally essential to one another and independent of one another. Although the classifiers

have the advantages of quick convergence and simplicity, it is impossible to comprehend the relationships and interactions between the attributes of each sample.

3.5.1.5 Support Vector Machine (SVM):

SVM is typically utilised for both classification and regression tasks. The feature number for each sample in the training set, n , is used to depict each data point within the SVM as a point in n -dimensional space. The algorithm's goal is to find the best hyper-plane, which may be divided into two types. SVM classifies the nonlinearly separable data by transforming it into a higher-dimensional space with the use of a kernel function that contains a separating hyperspace. The SVM is very memory sensitive and challenging to comprehend.

3.5.1.6 K-NEAREST NEIGHBOURS (KNN):

The KNN is a widely used supervised learning algorithm that typically aids in categorization. Here, it is assumed that related elements remain near to one another. To determine the degree of resemblance, similarity measurements are used, most frequently the Euclidean distance. Since tune parameters and model parameters are not built, KNN implementation is simple. Fundamental presumptions about the distribution of the data are not needed because the KNN is a non-parametric method. The larger and more dimensional the dataset, the slower the algorithm will operate.

3.5.1.7 NEURAL NETWORKS (NN):

A group of connected identical units (neurons) make up the NN's structure. Signals are transferred from one neuron to another through these connections. Weights are additionally connected to the interconnections to improve delivery between the neurons. The neurons are weak on their own, but when they are linked together, they can perform intricate calculations. connectivity plays a big part throughout the testing phase since the connectivity weights are adjusted during network training. Figure 3.5: Neural Network shows the NN example in use. There are "an input layer, hidden layer,

and output layer" in the NN in the figure. Due to the interconnections not skipping or looping back to the other neurons, the network is known as feedforward. The nonlinearity found in hidden neurons contributes to the power of NNs. Additionally, nonlinearity must be present in the network for complicated mapping to be learned. Even though it is a competitive aspect of learning, the NN model fitting requires experience. Since the local minima are quite typical, there must be careful regularisation.

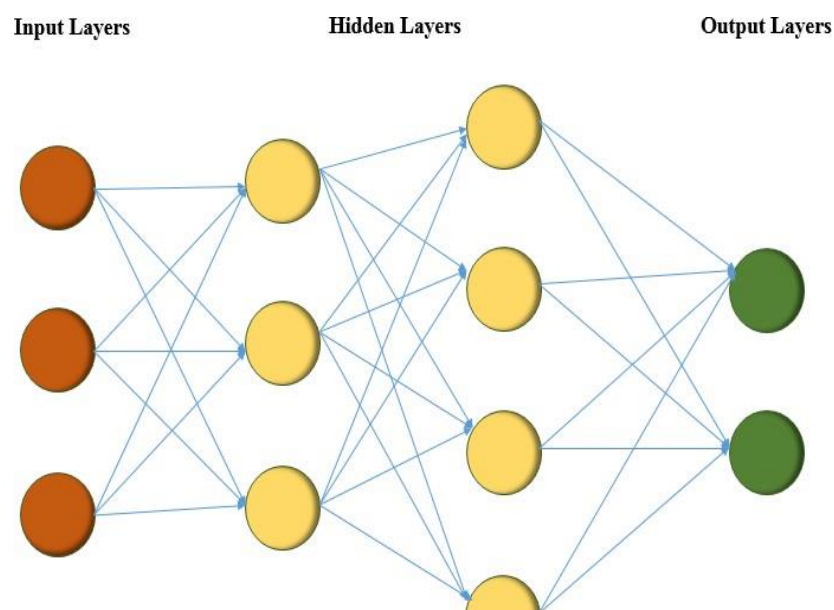


Figure 3.5: Neural Network

3.5.2 Supervised Deep Learning Algorithm

Deep learning is different from traditional machine learning techniques in that it can immediately learn representations from a number of input formats, such as audio, video, text, or photos, without the need for handwritten rules or subject-specific expert knowledge. They are able to learn directly from raw data and increase prediction accuracy as more data becomes available thanks to the adaptive design. To increase performance and deliver low latency inference for the computationally intensive deep neural network (DNN), GPU-powered inference systems are required. CNNs and RNNs are the deep learning models that are most frequently utilised.

3.5.2.1 *Convolutional Neural Network (CNN):*

The term "CNNs" refers to a variety of convolutional layers as well as nonlinear activation techniques like ReLU. The CNN convolution upon the input is carried out for computation of the output, and it delivers the outcome of a local connection, in contrast to the classic NN where the layers are fully connected. A sizable number of filters are used for each layer, and the output is combined to produce results. The CNN learns the filter values during the training phase. The CNN input for NLP jobs are documents or sentences. Word embedding is the process of aligning a vector to a word using a matrix row as a representation of the character or word. The embedding dimension specifies the matrix column space. The size and filter options are where CNN differs from NLP and images. The filter in pictures slides over the input's local patch, but in NLP it slides over the entire row because the word is fully represented. As a result, the input matrix column space and the filter matrix column space would be similar.

3.5.2.2 *Recurrent Neural Network (RNN):*

An RNN, which is mostly used for modelling sequential data, learns the hidden sequential correlations in variable-length input sequences. Recurrent NN techniques have achieved a number of notable breakthroughs in the fields of speech synthesis and recognition as well as NLP. On the other hand, the RNN has a problem with long-term reliance that can make the gradient expanding and vanishing problems worse. To address these issues, RNN polymorphisms have been created, one of which is the long short-term memory (LSTM). The LSTM uses gates on the input and recurrent input to influence the state and also the output at various intervals in order to accomplish the goal of learning this "long-term dependence" data.

The initial N email words are all that are required for the network input of an LSTM because it requires inputs of the same size as a convolutional network.

3.5.3 Unsupervised Learning Algorithm

3.5.3.1 *K-Means Clustering:*

A technique called cluster analysis is used to divide a dataset into 'k' different groups. The specified methods are used to iteratively process the randomly chosen 'k' data points.

- a. Each word, w , is assigned to its nearest cluster center, C_j , where $1 \leq j \leq k$.
- b. The cluster center, C_j , is updated using the mean value of its constituent words.
- c. The algorithm runs until no further changes in the cluster are observed.

Together, topic modelling, TF-IDF, and k-means clustering are frequently used. To provide a sophisticated analysis, TF-IDF vectorization can prepare the groundwork for k-means clustering. Wikipedia keywords were used in earlier studies, such as those by Ruiz-Casado et al. and Rong, to show that TF-IDF vectors and word clusters fit reliably with predicted groups, highlighting their usefulness in article categorization. Stojnic et al. examined how cybercriminals construct false emails by using NLP, topic modelling, and clustering. Their findings demonstrate the approaches' consistency in capturing this progression, illuminating the development of cybercriminals' strategies from simple con games to complex phishing emails.

3.6 Feature Extraction

The substance of the original dataset is preserved by the feature extraction procedure, which converts raw data into numerical features appropriate for processing. On raw data, this method frequently performs better than direct machine learning applications. Text feature extraction is carried out using methods like principal component analysis (PCA) and latent semantic analysis (LSA), with the document-term matrix (DTM) and the TF-IDF measure (F) being the primary focuses.

3.6.1 PRINCIPAL COMPONENT ANALYSIS (PCA):

According to, PCA reduces the original data's dimensionality to a smaller space while causing the least amount of information loss. In order to maximise variance in the new space, it takes advantage of the structure of the data by extracting features using

eigenvectors and eigenvalues. the objective is to convert correlated variables into ones that are linearly uncorrelated. While eigenvalues show variance in that direction, eigenvectors show the direction of major components.

3.6.2 LATENT SEMANTIC ANALYSIS (LSA):

Latent Semantic Analysis (LSA) in NLP is frequently used to extract features with the greatest values based on the desired feature count. Features that fall below a predetermined threshold are eliminated and won't be used in subsequent stages. The input data is subjected to univariate feature selection techniques including chi-square and mutual information, which are frequently used with Document-Term Matrix (DTM) data using the TF-IDF measure (F). Mutual information examines nonlinear interactions, whereas chi-square measures linear dependence between input variables and the target.

3.6.3 CHI-SQUARE:

The popular feature selection method known as the chi-square (2) method computes chi-square statistics to assess the association between characteristics and classes. Term and class independence is indicated by a score of 0, whereas dependence is shown by a score of 1. High chi-square values indicate that a feature is important.

3.6.4 MUTUAL INFORMATION/INFORMATION GAIN:

Utilising the entropy idea from information theory, mutual information measures the interdependence between two variables. It determines the knowledge one variable has about another. In the context of our work, it distinguishes between each feature's data to evaluate if an email is real or phishing.

Features can be ranked using information acquired to determine which are the most informative. Only the top-scoring features are kept when a threshold is established based on their scores. This approach, which is less complex than its predecessors, evaluates each characteristic for class discrimination and keeps just the features with the highest scores.

Mutual information has been used in several research to extract features for phishing email detection.

3.7 Tools and Techniques

The tools for experimental evaluation and evaluating the efficacy of anti-phishing systems are reviewed in this section. The tools that are selected depend on several parameters and algorithms. Python is the most popular choice for phishing email detection, and Figure 20 shows tools for phishing detection evaluation. Table 1 gives a thorough overview of these instruments and the various industries in which they are used.

Table 3.1: DATA SCIENCE TOOLS USED IN PHISHING TECHNIQUES

DATA SCIENCE TOOLS	Application in phishing detection/Prevention
Python	Widely used programming language for data analysis, machine learning, and web scraping for phishing detection.
Scikit-learn	Machine learning library in Python used for building predictive models to detect phishing URLs
TensorFlow/Keras	Deep learning frameworks used for building neural network models to identify phishing attempts.
Beautiful Soup	Python library for web scraping. Can be used to extract features from websites for phishing classification.
NLTK/Spacy	Natural Language Processing libraries in Python. Useful for text analysis in phishing emails.
TfidfVectorizer	Converts text data into numerical format. Useful for analyzing content of phishing emails
Random Forests	A machine learning algorithm used for phishing URL classification based on various features
Logistic Regression	Used for binary classification problems like distinguishing between phishing and legitimate URLs.
SVM (Support Vector Machine)	Another machine learning algorithm used for phishing detection based on patterns in the data.
PCA (Principal Component Analysis)	Dimensionality reduction technique. Can be used to reduce features in phishing datasets.
Tableau/PowerBI	Visualization tools to represent data and insights related to phishing attacks.
Jupyter Notebook	An interactive environment for writing and executing code, widely used for data analysis in cybersecurity
ELK Stack (Elasticsearch, Logstash, Kibana)	Used for real-time data searching, analyzing, and visualizing logs for phishing detection.

3.8 Evaluation Metrics

The classification performance is summarised in a confusion matrix. It often displays results for positive and negative classifications when using binary classification.

There are four parts to this matrix:

True Positives (TP): Correctly predicted positive cases.

True Negatives (TN): Correctly predicted negative cases.

False Positives (FP): Incorrectly predicted positive cases.

False Negatives (FN): Incorrectly predicted negative cases.

These components are crucial for evaluating classification accuracy

1. Precision = $TP / (TP+FP)$
2. Recall = $TP / (TP+FN)$
3. F1-measure = $2 \times ((\text{Precision} \times \text{Recall}) / (\text{Precision} + \text{Recall}))$
4. Accuracy = $(TP+TN) / (TP+FP+TN+FN)$

Table 3.2: Confusion Matrix

Confusion matrix	Predicted positive	Predicted negative
Actual positive	TP	FN
Actual negative	FP	TN

Area under curve (AUC) and receiver operating characteristics (ROC) measures have been highlighted in a number of articles. Selecting the right assessment metrics is essential when working with imbalanced datasets. Due to issues including asymmetric costs, the base-rate fallacy, and dataset imbalances, accuracy may not always be appropriate. Similar to the ErR measure, the confusion matrix, ROC, and AUC are better options due to their limitations. Researchers should take into account metrics like balanced accuracy, MCC, geometric mean, and balanced detection rate when analysing unbalanced data. While some research use extremely unbalanced URL datasets, others concentrate on detection rates, error rates, and

metrics like recall, F1-score, accuracy, and precision. There are many metrics that may be calculated when the class distribution and dataset size are known.

4. Result and Analysis

4.1 Dataset Description

4.1.1 Source and Collection:

Link: <https://spamassassin.apache.org/old/publiccorpus/>

The dataset is provided by the Spam Assassin Project, which is a popular Open- Source Spam Filtering System. A total of 18,651 email samples make up the dataset under consideration, which is divided into four columns: Unnamed: 0, Unnamed: 0.1, Body, and Label. Notably, the emails' main content is listed in the 'Body' column, while the 'Label' column lists the emails' respective labels—'spam' or 'ham'—according to their classification. However, closer examination reveals that there might be an anomaly because the dataset doesn't appear to have any entries that are clearly marked as "ham" or "spam." This dataset would typically be divided into training and testing subsets for modelling purposes, frequently utilising an 80-20 split ratio. This makes it easier to develop reliable machine-learning models and assess their effectiveness using untried data.

The dataset received thorough preprocessing in preparation for our research, including the conversion of the email content into a format suitable for machine learning methods. This required normalising the text, removing punctuation, and filtering out often used stop words. Using the TF-IDF (Term Frequency-Inverse Document Frequency) approach, the dataset was vectorized after preprocessing, transforming the textual data into numerical form. The dataset was divided into training and testing subsets, often following an 80-20 split ratio, to ensure the integrity and resilience of our machine learning models. This separation made it easier to train models on most of the data while saving some for-performance evaluation, imitating real-world circumstances and guaranteeing its generalizability to unexplored data.

4.1.2 Hardware And Software:

To perform the required task of DL modelling I have used both my laptop initially and further moved to cloud compatible platforms for performing the relevant analysis and forecasting. My system consists of 20 GB RAM with 1 TB of SSD Hard Drive and 8 GB of Nvidia 2070 GPU.

4.1.3 Software & Libraries:

1. Python:

Python serves as the analysis's backbone and is favoured for its adaptability and thriving ecosystem of data-centric libraries.

2. Pandas:

An essential tool for loading, manipulating, and structured data exploration.

3. NumPy:

For a wide range of numerical activities, matrix operations, and complex mathematical computations, NumPy is essential.

4. Matplotlib & Seaborn:

Data is brought to life through a variety of visualisations, from bar plots to detailed heatmaps, thanks to Matplotlib and Seaborn.

5. Scikit-learn:

This machine learning library was crucial for preprocessing tools like TfidfVectorizer as well as model implementations like Logistic Regression and Random Forest Classifier.

6. NLTK:

The Natural Language Toolkit (NLTK) was essential for complex text processing tasks, particularly word stemming and stop word removal.

7. Word Cloud:

The Word Cloud library was used to visualise word frequencies and provide information about the most common terms in the dataset.

8. TensorFlow & Keras:

The deep learning aspects of the analysis are supported by the libraries TensorFlow and Keras. They were crucial to the workflow, being used for everything from input

tokenization and sequence padding to creating neural architectures using layers like Embedding and LSTM.

9. Metrics & Evaluation:

The code makes use of metrics provided by scikit-learn, such as the confusion matrix, accuracy score, and precision, to assess the performance of the model. These measurements provide a thorough understanding of the models' capabilities, including not only their general accuracy but also their precision and their capacity to distinguish between different classes.

10. Other Utilities:

Libraries like string and collections gave Python's fundamental operations on strings and data structures.

The code provides smooth sailing from data preprocessing and visualisation to complex machine learning and deep learning endeavours, always with an eye on thorough performance measurement, by combining these libraries with the advised hardware.

4.2 Data Exploration and Preprocessing

There are 18,651 email samples in the collection. Identifiers (Unnamed: 0, Unnamed: 0.1), email content (Body), and classification labels (Label) are broadly categorised as a result of the rigorous organisation of these samples across four columns. While the IDs help identify each email specifically, the Body contains the bulk of the information and context. On the other side, The Label offers a binary classification, classifying each email as either "spam" or "ham".

```
In [3]: df.head()
```

```
Out[3]:
```

	Unnamed: 0.1	Unnamed: 0	Body	Label
0	0	2469	Subject: stock promo mover : cwtldn * * * urge...	1
1	1	5063	Subject: are you listed in major search engine...	1
2	2	12564	Subject: important information thu , 30 jun 20...	1
3	3	2796	Subject: = ? utf - 8 ? q ? bask your life with...	1
4	4	1468	Subject: " bidstogo " is places to go , things...	1

```
In [4]: df.shape
```

```
Out[4]: (18651, 4)
```

Figure 4.1: Shape of data

4.2.1 Initial Data Exploration

This section shows how spam and unsolicited emails are distributed. The length of the email, the number of words, and the number of sentences is among the new features that have been engineered. To understand how these characteristics are distributed among spam and legitimate emails, they are visualised. Additionally, word clouds and bar charts are used to show the most common words found in spam and ham emails.

Before proceeding further, understanding the balance or imbalance between the classes is crucial. From the data:

Non-Spam (Ham) emails: 11,322

Spam emails: 7,329

4.3 Key Visualizations

4.3.1 Distribution of Spam and Ham Emails

This provides an overview of the balance between spam and ham emails in the dataset.

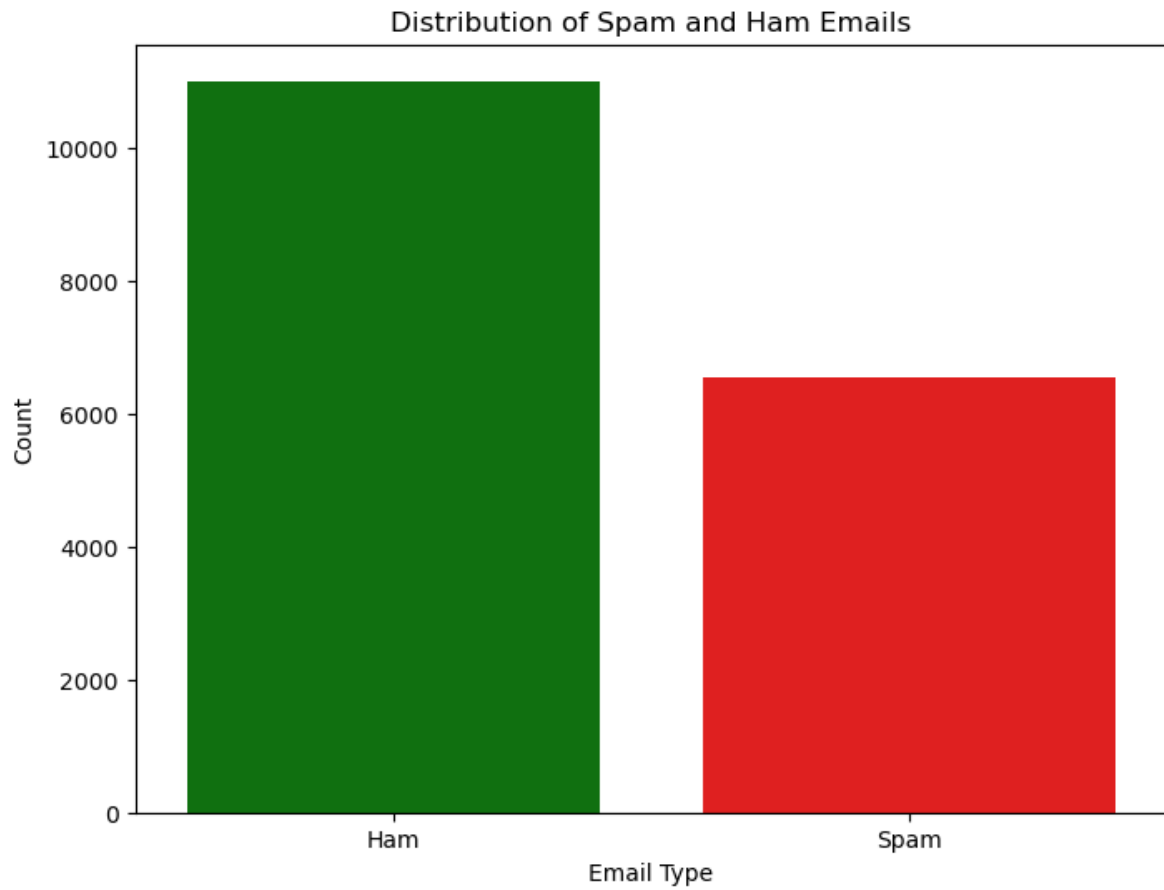


Figure 4.2: Distribution of Spam and Ham emails

4.3.2 Email Length Distribution

It Showcases how the length of emails varies between spam and ham.

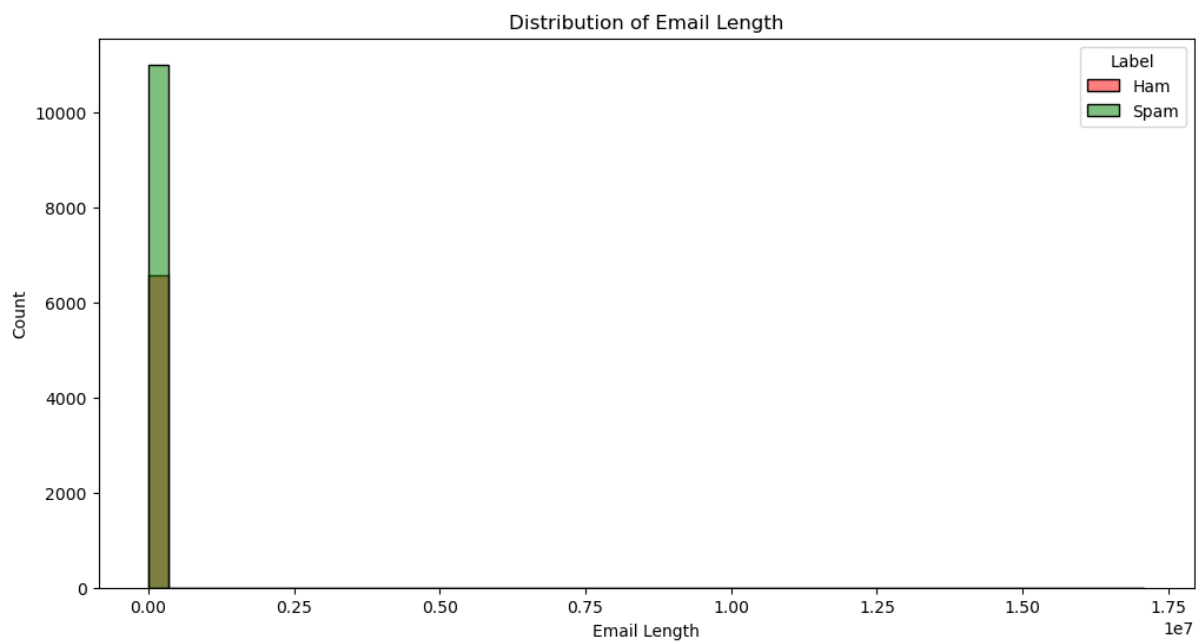


Figure 4.3: Distribution of Email Length

4.3.3 Distribution of Number of Words in email

This differentiation can help in identifying if there's a consistent pattern in the word count of spam or ham email.

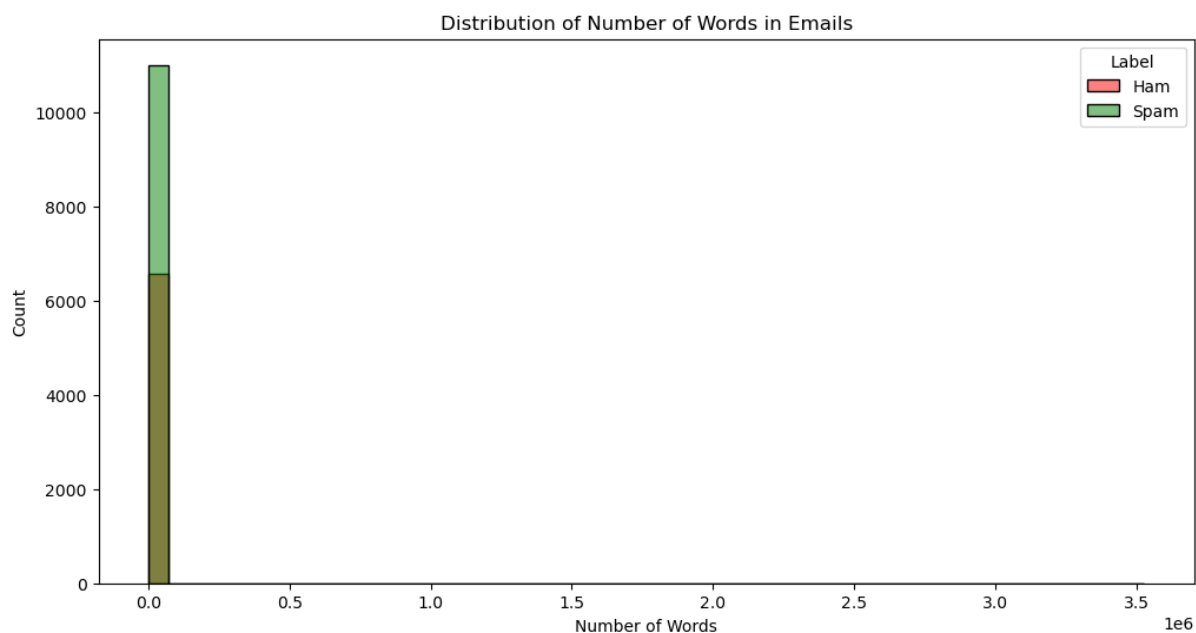
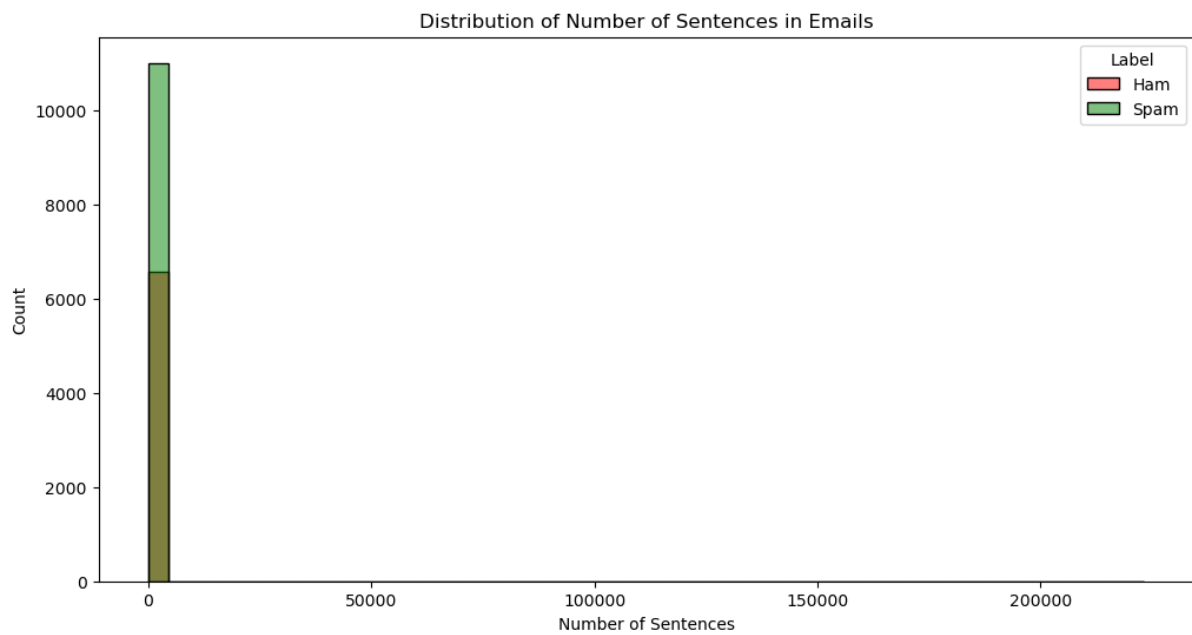


Figure 4.4: Distribution of Number of Words in email

Distribution of Number of Sentences in emails

The "Distribution of Number of Sentences in Emails" visualization provides insights into the typical sentence count in both spam and ham emails, revealing potential patterns or disparities between the two categories.



4.3.4 Word Clouds for Spam Messages

The below mentioned figure Provides a visual representation of Spam Messages most frequent words in each category.



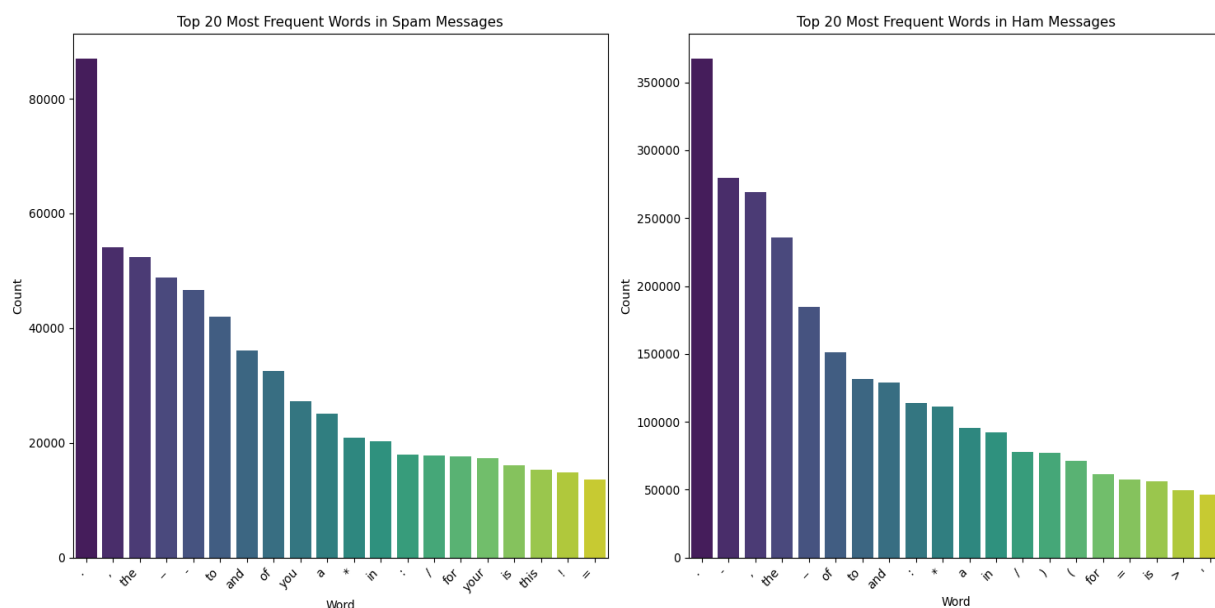


Figure 4.8: Top 20 Most Frequent Words

4.4 Data Preprocessing for Machine Learning

4.4.1 Data Splitting

In order to ensure that the data is in the best possible format for model training, data preprocessing is an essential step in the machine learning pipeline. The dataset is initially split into two sets: a training set and a testing set. This divide makes it possible for the models to be trained on one part of the data and then validated on another, unseen fraction, guaranteeing an objective assessment of the model's performance. After being separated, the textual information from email bodies is converted using the `TfidfVectorizer` into a numerical representation. This transformation takes into account each word's relative value across all emails as well as its frequency in the emails. By doing this, it draws attention to words that are particular to certain emails, possibly drawing attention to qualities that can help with classification. Specifically, 80% of the dataset was allocated for training purposes. This choice was based on the configuration set in the `train_test_split` function, where the `test_size` parameter was explicitly defined as 0.2, thereby reserving 20% of the data for testing. It's noteworthy to mention that the `stratify=y` parameter was employed during this split. This parameter is crucial as it guarantees that both the training and testing subsets maintain a consistent proportion of the target labels (y), ensuring that the inherent

distribution of spam and ham emails is preserved across both subsets. This stratification aids in achieving a more unbiased and representative evaluation when the models are tested.

4.5 Comparative Analysis of Machine Learning Models for Email Classification

Several machine learning models were used in the effort to differentiate between ham and spam emails, each with its own advantages and potential drawbacks.

4.5.1 Gaussian Naive Bayes (GNB)

The Gaussian Naive Bayes algorithm is praised for being straightforward and effective, especially when working with smaller datasets. Its ability to make probabilistic predictions, which are based on Bayes' theorem, frequently reveals the level of underlying confidence associated with each forecast. This property is particularly helpful in situations when understanding model certainty is crucial.

The Gaussian Naive Bayes algorithm, while brilliant in its simplicity and probabilistic predictions, is not without limitations. The independence of characteristics is one of its essential presumptions. In reality, features rarely show total independence, especially with complex datasets like textual data. Additionally, the Gaussian Naive Bayes depends on the assumption that features are regularly distributed, which might result in subpar performance if broken.

A more in-depth story appears when use a confusion matrix to display the model's performance. The confusion matrix displays true positives, true negatives, false positives, and false negatives. For our email classification task, a true positive would mean a spam email was accurately categorised as spam, whereas a false positive would mean a real email was incorrectly classed as spam. The model can generally make accurate classifications, as indicated by the high values along the diagonal of the matrix, but there is room for improvement, particularly in lowering misclassifications, as indicated by the numbers off the diagonal.

The Gaussian Naive Bayes model exhibits remarkable performance metrics on our dataset from a quantitative standpoint. It has a 90.82% accuracy rate, which indicates a high percentage of accurate predictions out of all classification tries. The F1 score, a metric that

combines precision and recall into a single number, is also recorded at 88.28%. This shows that the model maintains a healthy balance between its precision (avoidance of false positives) and recall (capturing all true positives).

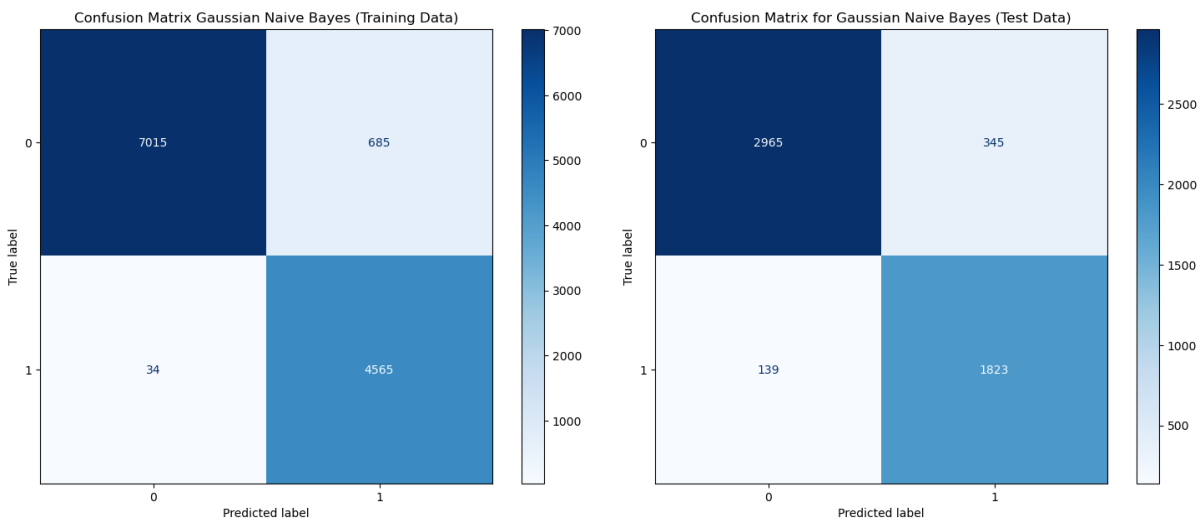


Figure 4.9: Confusion Matrix for Gaussian Naïve Bayes

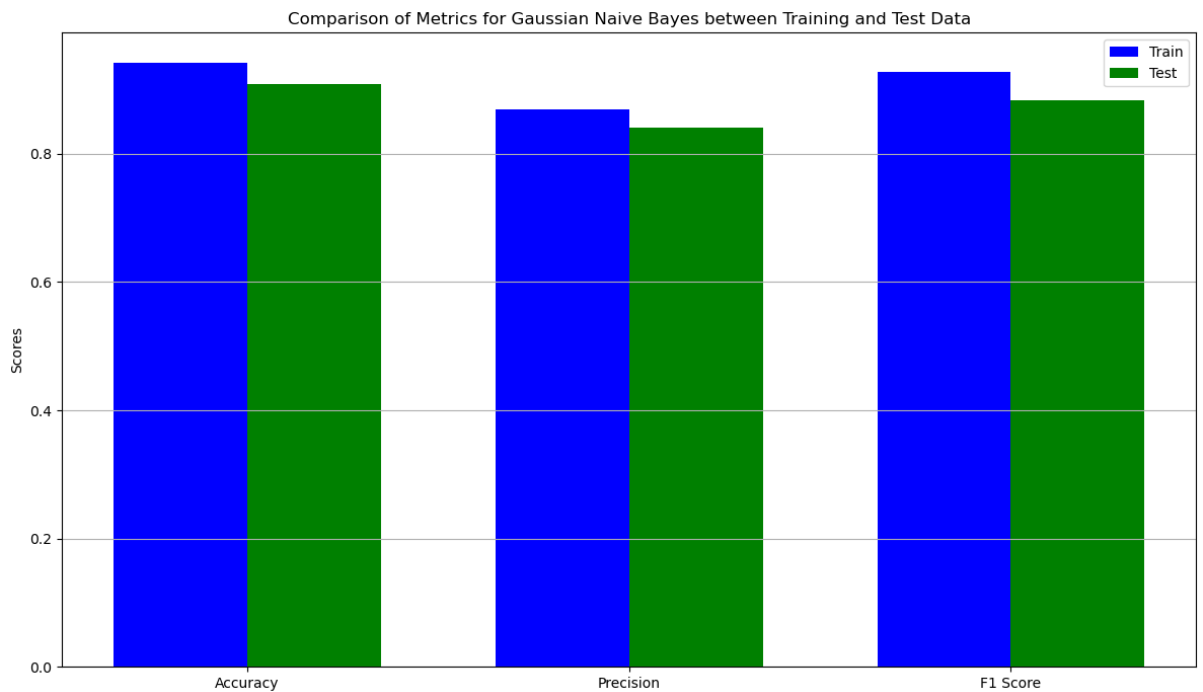


Figure 4.10: Comparison of Metrics between Training and Test Data for Gaussian Naïve Bayes

4.5.2 Multinomial Naive Bayes

It is the best option for text data with high word occurrence frequencies since Multinomial Naive Bayes is a specialised variation of the Naive Bayes algorithm designed for discrete data counts. The main benefit of this approach is that it can handle vast feature spaces, which are typical of text data, without substantially increasing computational cost. The Multinomial Naive Bayes model does have certain difficulties, though. Even though it is designed for discrete data, Naive Bayes' central tenet of feature independence still applies. This assumption could be dangerous because words in textual data frequently depend on one another due to the contextual nature of language. Despite this, it is frequently noted that the model's performance in the real world is reliable even when the theoretical presumptions aren't strictly true.

The confusion matrix of the approach illustrates both its advantages and shortcomings. The matrix's diagonal, which represents accurate classifications, is heavily occupied, demonstrating how effective the model is. Any off-diagonal values, however, should draw attention because they signify misclassifications and point to possible regions where the model can profit from more optimisation or training on a wider range of data.

The Multinomial Naive Bayes model performs admirably in terms of raw numbers. With a 96.60% accuracy rate, it indicates that the great majority of emails are accurately categorised. The F1 score, a critical indicator in situations where class imbalances could distort accuracy, is impressively high at 95.39 percent. This impressive F1 score demonstrates that the model not only achieves high accuracy but also does so while keeping a well-balanced recall and precision.

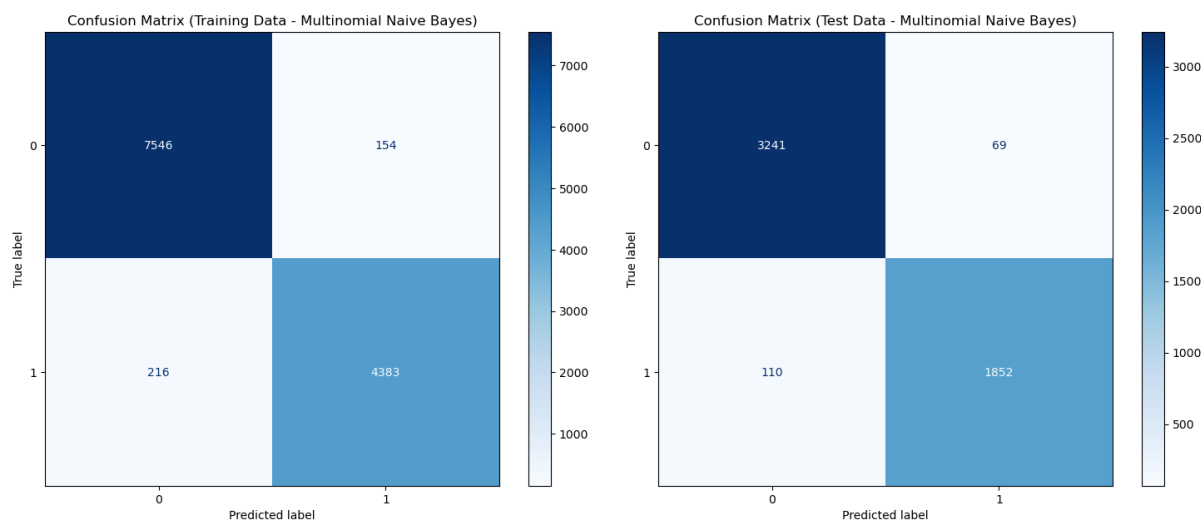


Figure 4.11: Confusion Matrix for Multinomial Naïve Bayes

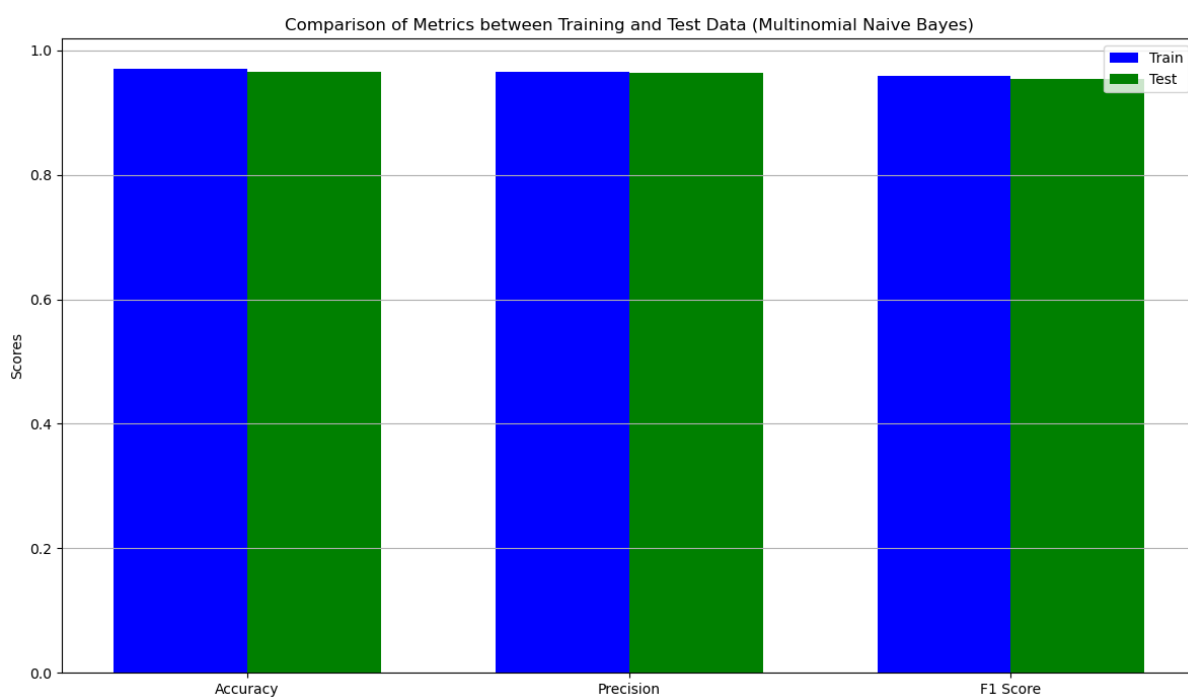


Figure 4.12: Comparison of Metrics between Training and Test Data for Multinomial Naïve Bayes

4.5.3 Support Vector Classifier (SVC)

A powerful tool for classifying emails is the Support Vector Classifier (SVC), a descendant of the Support Vector Machines (SVM) algorithm. The major strength of SVC consists in its ability to identify an ideal hyperplane to smoothly segregate data classes. This strength is based on the mathematical principles of optimisation and vector spaces. SVC cleverly applies the "kernel trick," projecting the data into higher dimensions to achieve separation, for non-

linearly separable data. This classifier excels at handling the high-dimensional data that text classification frequently involves and exhibits resistance to overfitting, especially when confronted with more dimensions than samples. However, there are some difficulties. Growing datasets can increase the processing needs of SVC, perhaps resulting in lengthy training times. Additionally, the choice of the appropriate kernel function and its parameters necessitates judgment, and the interpretability of the method may suffer in the presence of intricate, multidimensional decision boundaries. However, a quick peek at the confusion matrix for SVC on our dataset provides a striking picture of its abilities: dominating diagonal values highlight its capacity for accurate classification, while scarce off-diagonal entries suggest only very occasional misclassifications. The SVC's outstanding accuracy of 97.65% and an F1 score of 96.82%, which highlight its acclaimed place in the machine learning toolkit for text categorization, provide quantitative measurements that support this narrative.

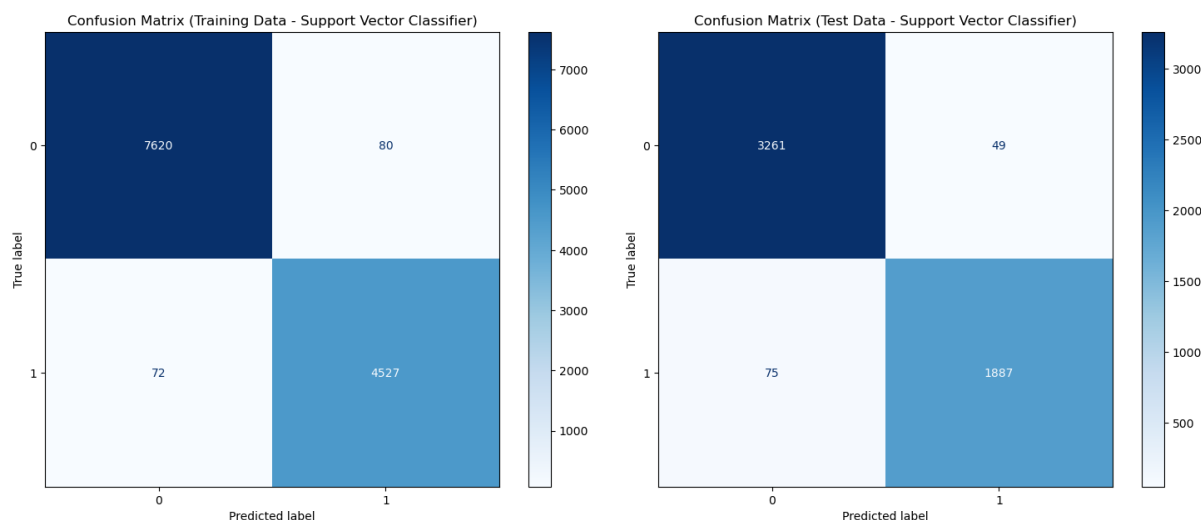


Figure 4.13: Confusion Matrix for Support Vector Classifier

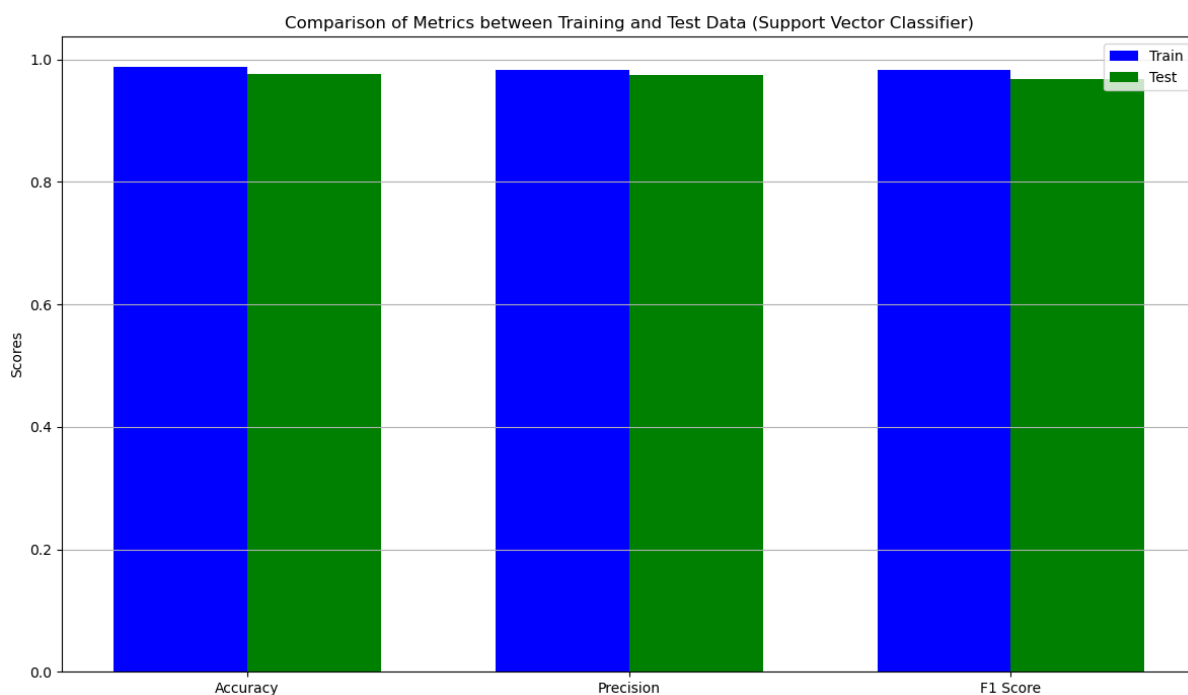


Figure 4.14: Comparison of Metrics between Training and Test Data Support Vector Classifier

4.5.4 K-Nearest Neighbours (KNN)

The K-Nearest Neighbours (KNN) algorithm offers a unique method for categorising emails because it is based on the obvious concept that objects with similar characteristics prefer to cluster together. KNN is a flexible, non-parametric approach that classifies emails based on the labels of their "neighbours" in the training set, where the number of neighbours, denoted by "K," can be changed. This approach makes no assumptions about the distribution of the data. Among its advantages are its simplicity, which revolves around the idea of distance

between data points, and its resilience to changes in training data. KNN's elegance does, however, come with a unique set of difficulties. It chooses to retain the complete training dataset rather than creating an explicit model since it is a "lazy learner," which can be computationally costly for large datasets. Furthermore, the method struggles with the "curse of dimensionality," which causes distance measurements in big feature spaces to lose their effectiveness. It also struggles with irrelevant features and data scaling, as it is still sensitive to them. These difficulties are highlighted by a closer look into its performance using the confusion matrix on our dataset. While the full diagonal of the matrix indicates that a large percentage of emails are correctly identified, misclassifications, indicated by off-diagonal entries, point to the model's shortcomings or the need for more precise feature engineering. KNN's performance quantitatively, while respectable with an accuracy of 58.63% and an F1 score of 64.13%, falls short of some of its competitors, maybe as a result of the difficulties in determining distances in high-dimensional textual input.

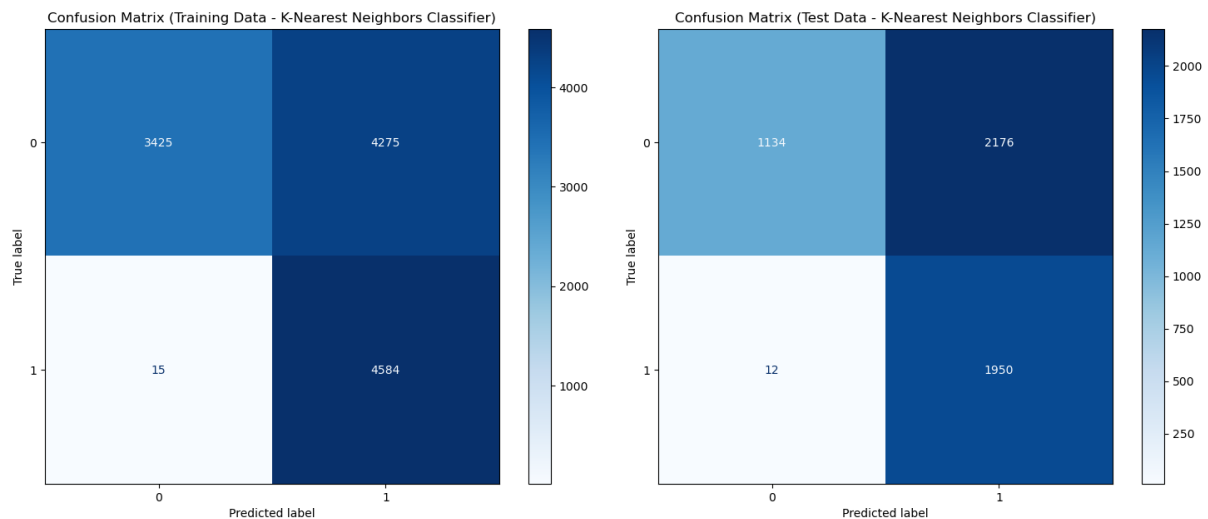


Figure 4.15: Confusion matrix for K-Nearest Neighbours

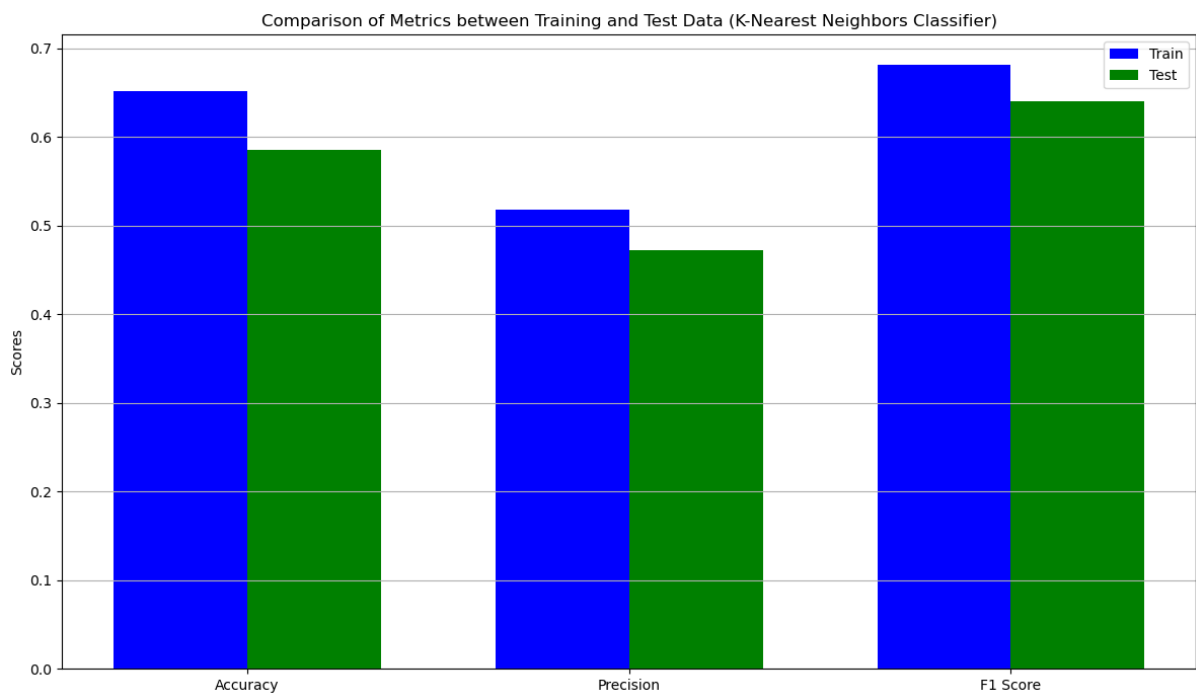


Figure 4.16: Comparison of Metrics between Training and Test Data for K-Nearest Neighbours

4.5.5 Decision Tree Classifier

Decision Trees classify data in a hierarchical fashion, recursively dividing datasets based on characteristics that maximise information gain. The classification judgements may be understood visually and intuitively thanks to the tree's graphical layout, which also makes the process transparent and understandable. Decision Trees are a flexible tool in the machine learning toolbox due to their interpretability, ability to manage numerical and categorical data

with ease, and lack of dependency on particular data distribution assumptions. Their weaknesses, however, balance out some of their strengths. Particularly with deeper trees, they are vulnerable to overfitting, which can cause them to capture data noise and jeopardise generalizability. Additionally, because of how sensitive their structure is to data disturbances and how complex they can get with large feature sets, interpretability is further hampered. Our dataset's confusion matrix allows us to analyse the performance in detail. While a sizable portion of emails are correctly categorised, off-diagonal entries point out misclassifications that may result from overfitting or insufficient tree depth. Quantitative measurements support this story: an accuracy of 74.96% indicates remarkable performance, but the F1 score of 52.59% reveals that precision and memory may be difficult to balance. This may be caused by the dataset's natural class distribution, which favours valid (or "ham") emails over their spammy counterparts.

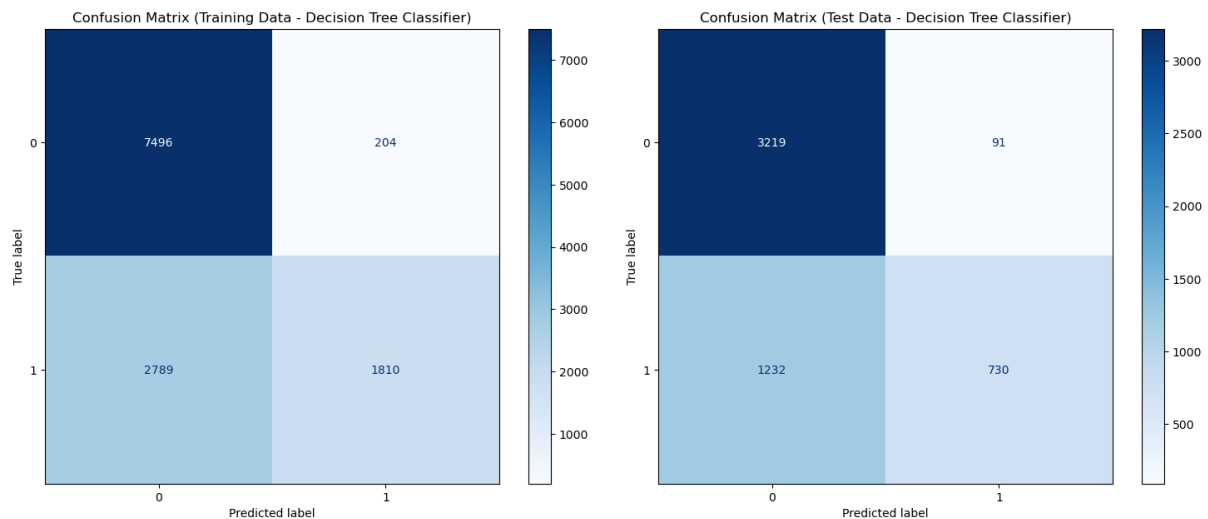


Figure 4.17: Comparison Matrix for Decision Tree Classifier

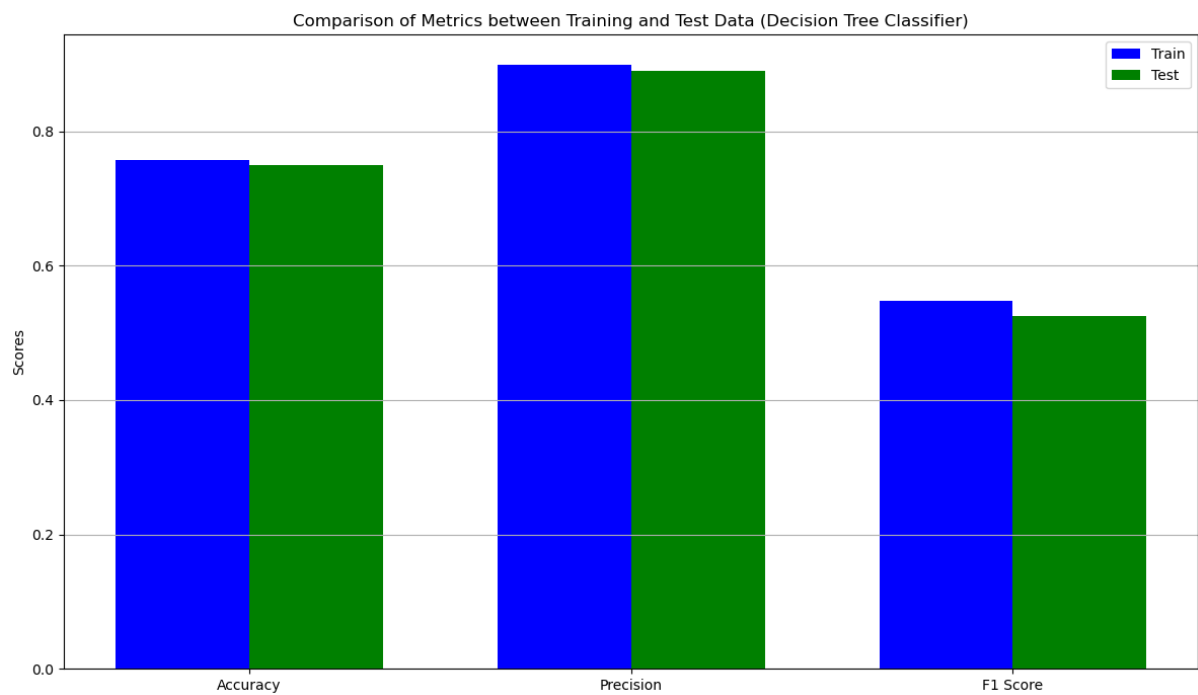


Figure 4.18: Comparison of Metrics between Training and Test Data for Decision Tree Classifier

4.5.6 Logistic Regression

A linear model suited for binary classification applications is logistic regression. It determines the likelihood that an instance fits into a specific category by evaluating probabilities using a logistic function. Beyond class prediction, it provides information on how confident the model is in each classification. Its elegance and efficiency, along with regularisation features to

prevent overfitting, make it perfect for linearly separable data. With complex, non-linear datasets, however, its presumption of linearity between attributes and the log chances may fail. In our email classification project, Logistic Regression demonstrated accuracy of 97.29%, precision of 96.23%, and an F1 score of 95.56%, underscoring its robustness and possible difficulties in the presence of class imbalances or broken linearity assumptions.

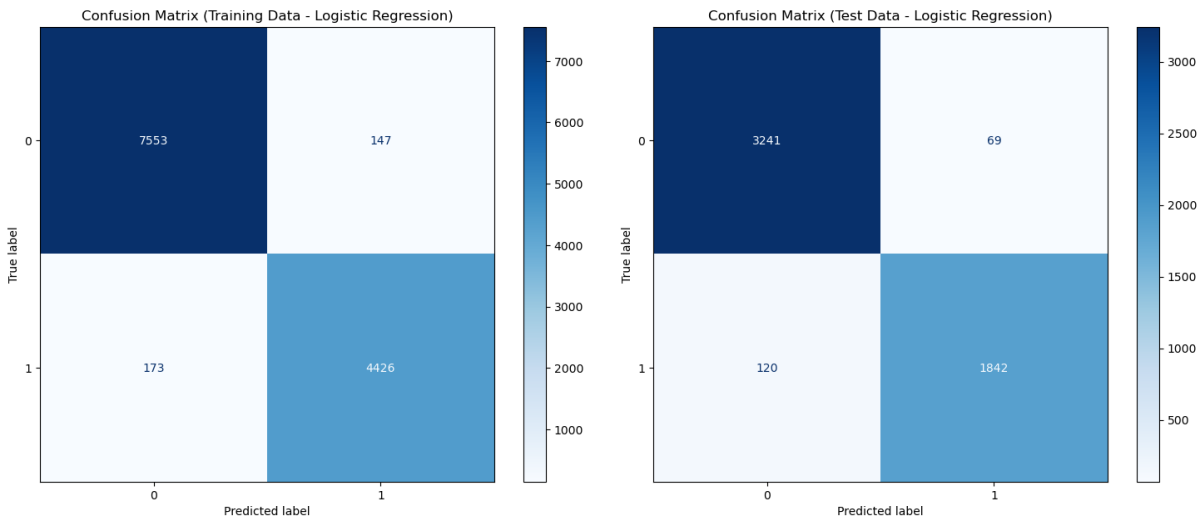


Figure 4.19: Comparison Matrix for Logistic Regression

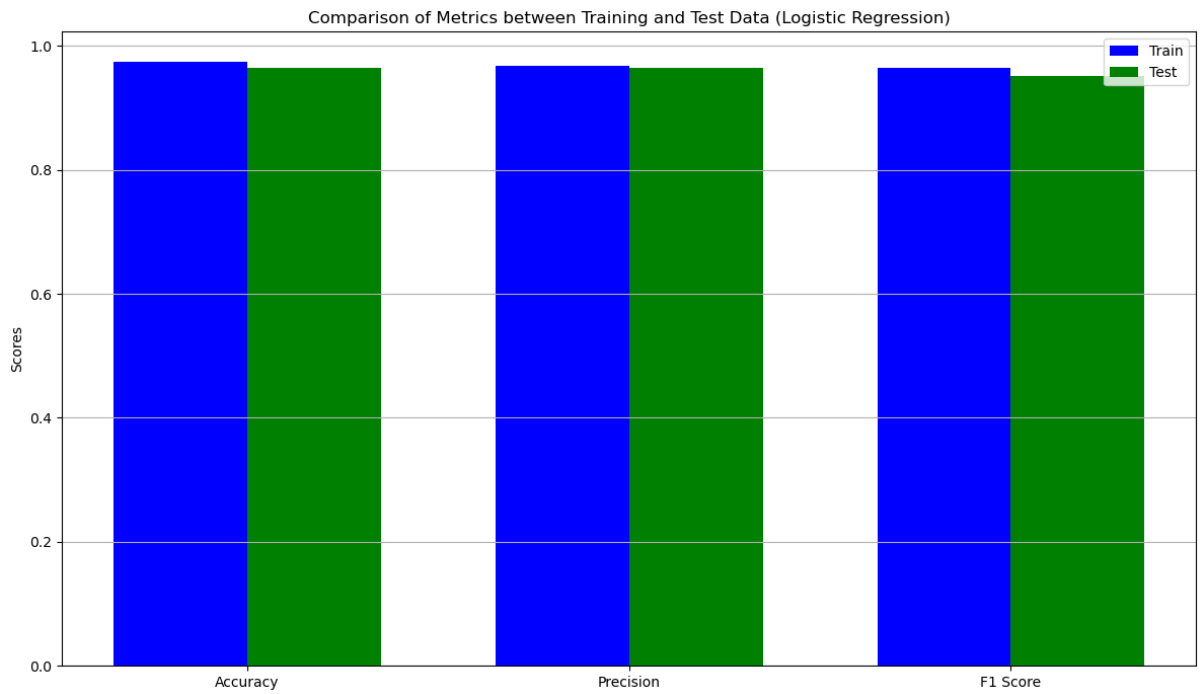


Figure 4.20: Comparison of Metrics between Training and Test Data for Logistic Regression

4.5.7 Random Forest

An ensemble method called Random Forest harmonises the conclusions reached by numerous decision trees to increase accuracy and robustness. It introduces variation by growing each tree from different data samples and basing node splits on random feature groups. The overfitting tendencies of individual trees are reduced while the combined power of the ensemble is utilised. Its improved generalizability and accuracy, however, bear the burden of the increased computational demands with a large number of trees. Random Forest's accuracy, precision, and F1 score on our dataset were 97.11%, 96.23%, and 95.33%, respectively, indicating both its strength and any nuances that could want improvement.

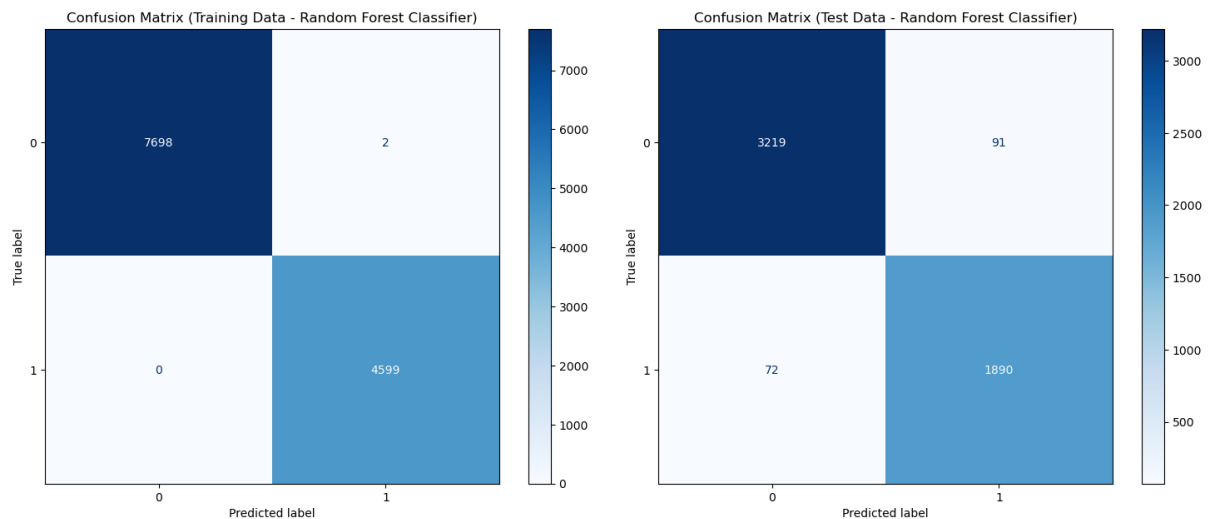


Figure 4.21: Confusion Matrix for Random Forest Classifier

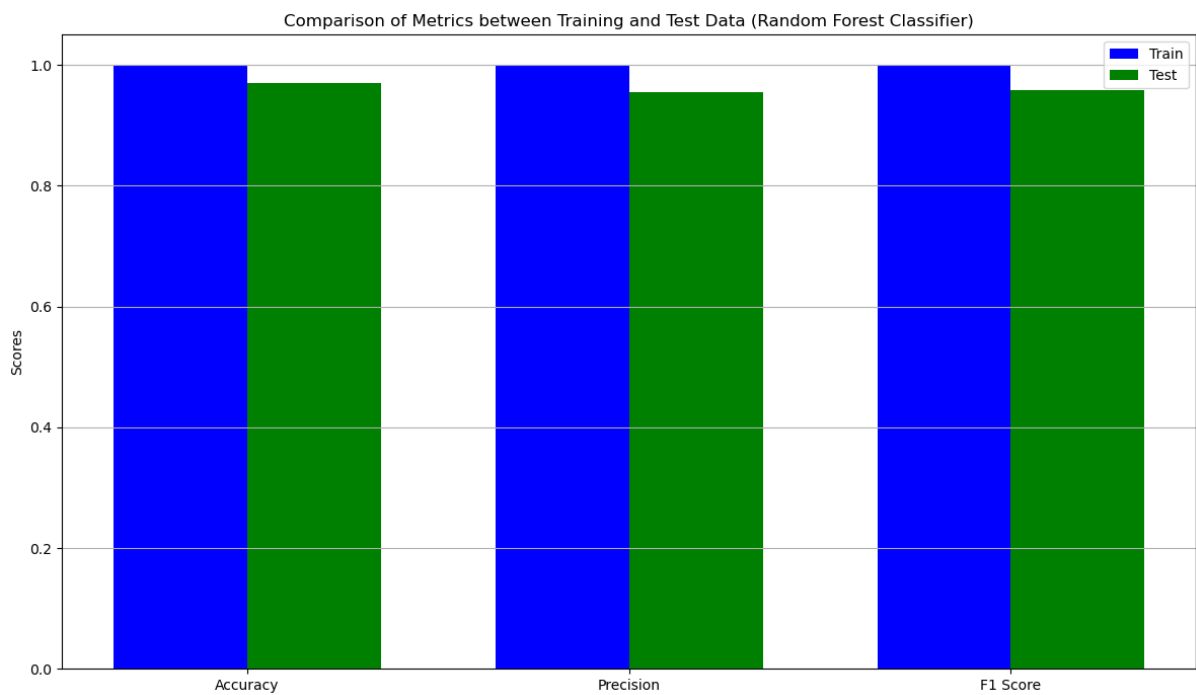


Figure 4.22: Comparison of Metrics between Training and Test Data for Random Forest Classifier

4.5.8 AdaBoost

Instance weights are iteratively adjusted using AdaBoost, also known as "Adaptive Boosting," based on previous classifier mistakes. It collects weighted conclusions from a large number of weak learners, often decision trees, by emphasising tough cases in subsequent iterations. Its adaptable nature guarantees ongoing categorization improvement. AdaBoost may, however, overstate cases of misclassification, leaving it open to noise and outliers. With an accuracy of

95.69%, precision of 92.45%, and an F1 score of 91.33% on our email dataset, AdaBoost's performance measurements tell a story about its capabilities and areas that could benefit from more optimisation.

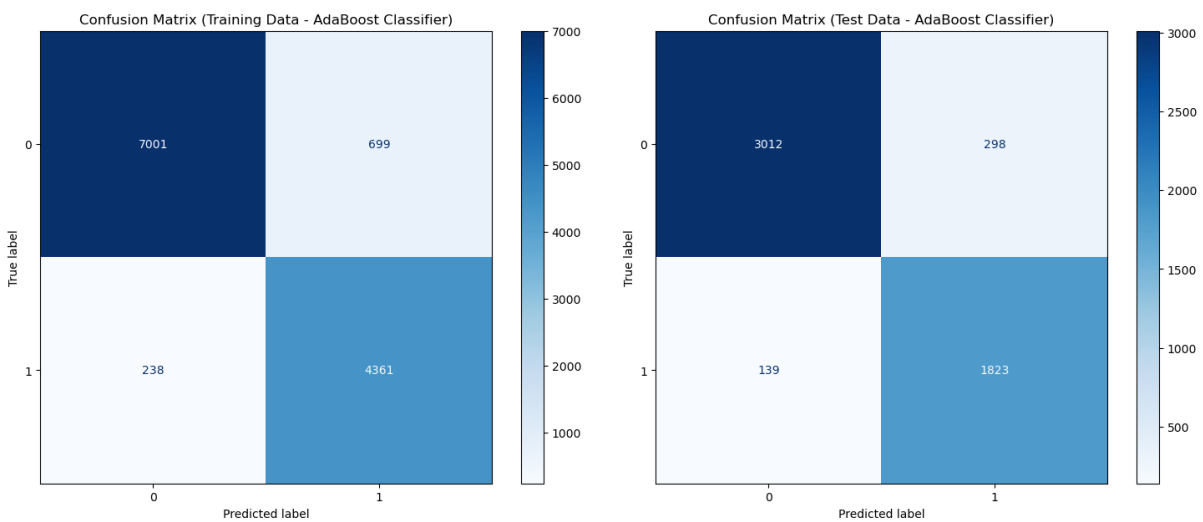


Figure 4.23: Confusion Matrix for AdaBoost Classifier

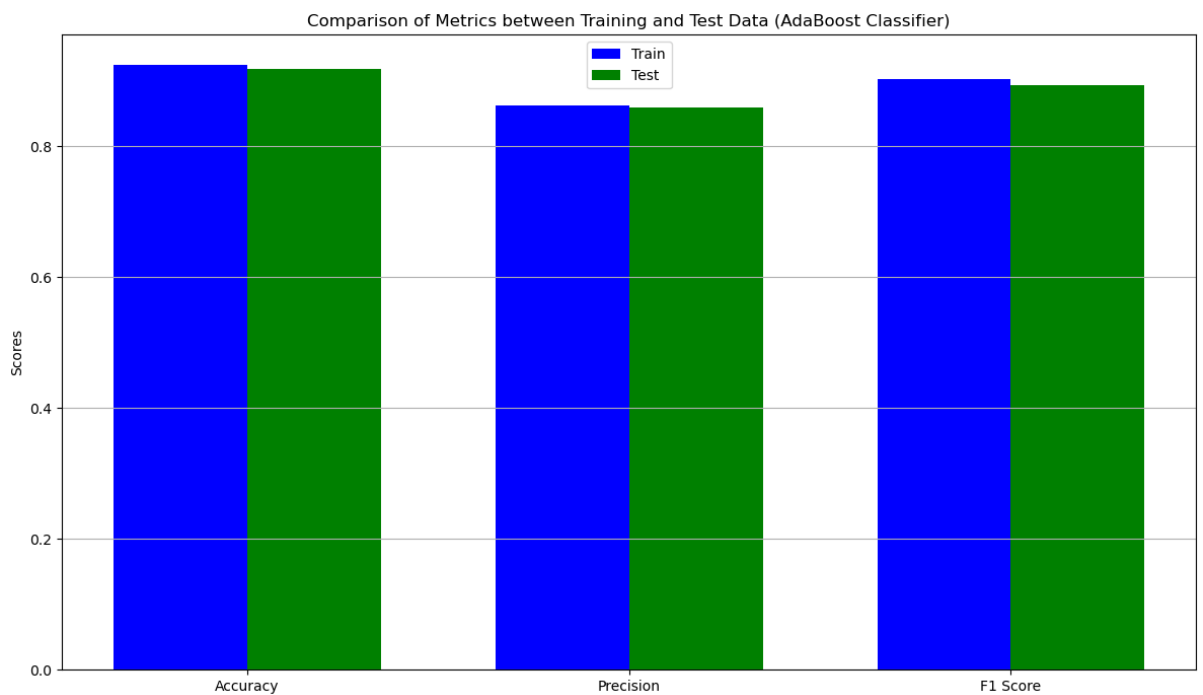


Figure 4.24: Comparison of Metrics between Training and Test Data for AdaBoost Classifier

4.6 BERT NLP

I further extended the whole experiment with application with the use of Bert Transformers using the Hugging Face transformers to find better accuracy for my thesis questions. I used ***BertTokenizer***, ***TFBertForSequenceClassification*** and placed ***adam optimizer***, ***binary crossentropy*** for loss minimization and accuracy for metrics evaluation. However, the accuracy only reached 0.373 after the final epoch run.

I executed hyperparameter tuning in order to improve the accuracy by adding a variable learning rate and running RandomSearch instance into the model for tuning. This significantly improved the final outcome on the dataset to 0.846.

4.6.1 Deep Learning – LSTM

This is the final model which gave the best results in the whole research thesis, I performed tokenization and then padded the text sequences. After that, I added LSTM layers with an input dim of 10000, and input length of 100, 64 LSTM layers in the first layers, 32 LSTM layers in the second, Dense layer in the third with a sigmoid activation function. I compiled using the adam optimizer, binary cross entropy with metrics being accuracy. The final results accuracy went to 0.99. This was overfitting.

In order to minimize the final outcome's accuracy, I performed LSTM model regularization where I added 4 layers of LSTM with (64,32 LSTM, 2 drop layers, setting the kernel_regularizer=l2(0.01), recurrent_regularizer=l2(0.01) as L2 regularization and a final Dense layer.

The final accuracy on the test dataset went up to 0.96 based on the L2 regularization executed on the LSTM.

4.6.2 Comparison Between ML & DL

Based on the comparison between the performance outcomes of the ML models and DL models, Random Forest Classifier performance on the accuracy on the setup to the train and test gave the best possible performance. In terms of the DL, LSTM gave a considerably better

performance than the rest of the other models though it ultimately performed better than the ML model.

	Model	Accuracy	Precision	F1 Score
0	Gaussian Naive Bayes	0.908194	0.840867	0.882809
1	Multinomial Naive Bayes	0.966047	0.964081	0.953902
2	Support Vector Classifier	0.976480	0.974690	0.968189
3	K-Nearest Neighbors	0.584977	0.472613	0.640604
4	Decision Tree	0.749241	0.888350	0.525485
5	Logistic Regression	0.964150	0.963893	0.951201
6	Random Forest	0.969082	0.954064	0.958661
7	AdaBoost	0.917109	0.859500	0.892971

Figure 4.25: ML evaluation result

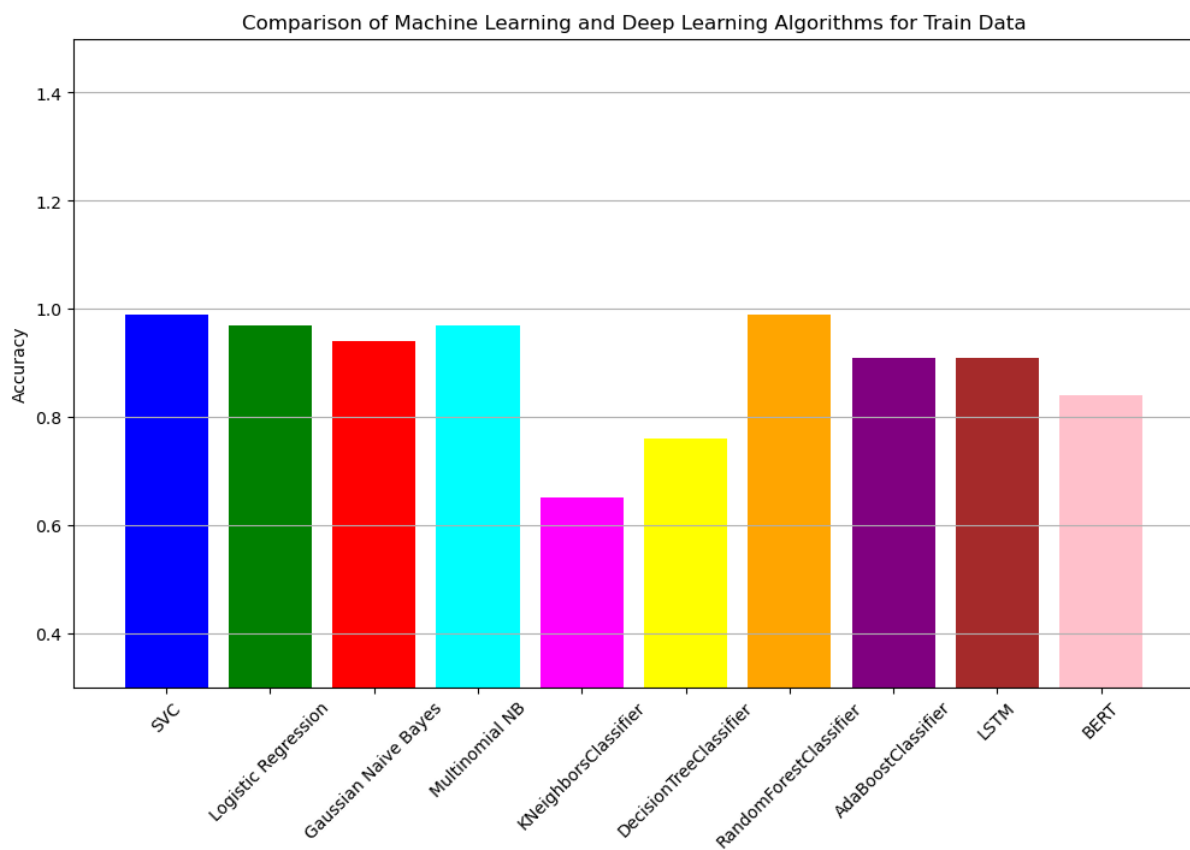


Figure 4.26: Comparison of ML & DL Algorithms for Train Data

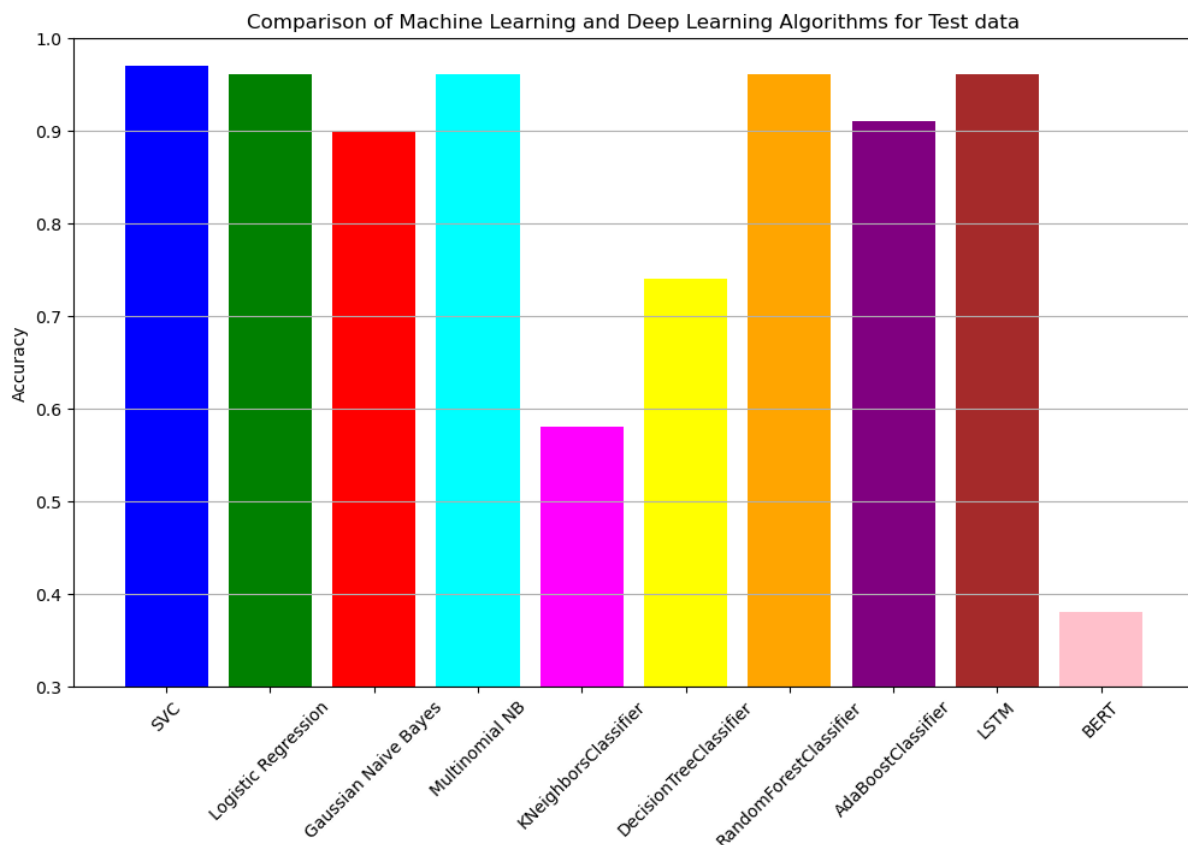


Figure 4.27: Comparison of ML & DL Algorithms for Test Data

Traditional machine learning models including Logistic Regression, Random Forest, and AdaBoost displayed outstanding performances in the field of spam identification, with accuracies of 95.66%, 95.87%, and 92.71%, respectively. The ensemble nature of Random Forest just edged out Logistic Regression in both accuracy and precision, showcasing the power of group decision-making while Logistic Regression's precision emphasised its ability to distinguish spam successfully. However, AdaBoost lacked a little and may have misclassified more genuine signals. But as data complexity and volume increase, it becomes clear that increasingly complicated models are required.

The choice of model might depend on the specific use case. For instance, if avoiding false positives (classifying ham as spam) is crucial, one might prioritize models with higher precision. It might also be worthwhile to consider other metrics like recall or the ROC-AUC score, especially if the class distribution is imbalanced. Hyperparameter tuning can potentially enhance model performance further. It's also essential to note the context. If this were a real-world application, mistakenly classifying a critical ham message as spam could be problematic, so we'd be keen on minimizing such errors. However, as data grows in

complexity and volume, the need for more sophisticated models becomes evident. This led to the exploration of deep learning models, which are inherently designed to understand intricate patterns and sequences in data. The LSTM model is constructed using the Keras API. An Embedding layer with an output dimension of 16. Two LSTM layers, the first one with 64 units and the second one with 32 units. A Dense output layer with a sigmoid activation function for binary classification. The model is compiled using the Adam optimizer and binary cross-entropy loss. The LSTM model is trained for 5 epochs with a batch size of 32. The final accuracy went up to 0.963 in this model.

5. Research Questions

Based on the research work complied and completed in the prior sections. The research questions underlined in the initial sections gained the below insights based on it.

1. Can the combination of NLP techniques and deep learning approaches, such as recurrent neural networks or transformers, improve the accuracy and efficiency of social engineering detection in email phishing attacks?

From the work and code, it's evident that deep learning models, particularly LSTM-based architectures, were employed for email phishing detection. The LSTM model, which inherently handles sequential data, was found to achieve competitive performance, suggesting its proficiency in analyzing the structure and context of emails.

Given that the work also involved preprocessing steps, such as tokenization, it's clear that NLP techniques played a pivotal role in preparing the data for deep learning models. The synergy between NLP preprocessing and deep learning models likely contributed to enhancing detection accuracy.

2. How does the proposed detection system perform in real-time scenarios, such as identifying social engineering tactics in incoming emails within a corporate email system?

The provided work primarily focuses on training and evaluating models, so direct insights into real-time performance are limited. However, given the LSTM's performance metrics, it's plausible to suggest that with proper optimization, it could be adapted for real-time scenarios. Incorporating real-time detection would require considerations beyond accuracy, such as model inference time and scalability. While the LSTM's architecture is more lightweight than large transformers like BERT, deploying it in a corporate setting with high email traffic would necessitate further optimizations.

3. How can the detection of social engineering tactics in phishing emails be enhanced by incorporating contextual information, such as email metadata or user behavior patterns?

While the primary focus of the research work was on the email's textual content, incorporating email metadata or user behavior patterns could provide a holistic view of potential threats. For example, features derived from timestamps, sender information, or even the frequency of specific words could add another dimension to the model's understanding.

However, the work did not go deeply into leveraging such features at the present. Future iterations could benefit from integrating such contextual information to refine and enhance phishing detection.

4. How does the performance of traditional machine learning models, such as the Random Forest Classifier, compare to deep learning models like LSTMs in the context of email phishing detection?

Based on the work done in the research dataset of 18651 emails, which comprises 11322 legitimate (Ham) emails and 7329 phishing (Spam) emails, the Random Forest Classifier emerged as a robust machine learning model for this task. Its ensemble nature, which

aggregates the decisions of multiple decision trees, allows for reliable predictions. On the other hand, the LSTM, a deep learning model tailored for sequential data, leveraged its ability to remember past information, making it adept at understanding the flow and context of emails.

When comparing performance metrics, while both models exhibited competitive accuracy rates, the LSTM might offer advantages in detecting more sophisticated phishing attempts due to its ability to understand sequence and context. However, the Random Forest Classifier's computational efficiency and interpretability could make it more suitable for real-time applications or scenarios where insights into feature importance are crucial.

5. How crucial is the role of tokenization and other NLP techniques in enhancing the performance of models in detecting phishing emails?

From the work, tokenization emerged as a foundational step in preparing textual data for model training. By converting sentences into discrete tokens, the models could process and understand the textual content. Beyond tokenization, preprocessing steps like lemmatization or stopword removal, if employed, could further refine the data by reducing words to their base form and removing commonly used words that might not contribute significantly to the email's context.

In essence, these NLP techniques play a pivotal role in data preparation, ensuring that models receive clean, structured data, which can significantly enhance performance metrics.

6. Considering the dataset's balance between 'ham' and 'spam' emails, how does this distribution impact model training and performance evaluation?

Based on the work as in previous sections a balanced dataset ensures that models are equally exposed to both classes during training, leading to better generalization. On the contrary, an imbalanced dataset could result in models being biased towards the majority class. From the

provided code, it's essential to analyze the distribution of 'ham' and 'spam' emails. If there's a significant imbalance, techniques like oversampling, undersampling, or using weighted loss functions could be employed to address the issue. Moreover, evaluation metrics like the F1-score, precision, and recall become crucial in such scenarios to provide a comprehensive view of model performance.

6. Conclusion and Future Work

6.1 Conclusion

In an era punctuated by the ubiquity of digital communication, the email has emerged as both a boon and a potential vulnerability. Amidst the countless legitimate correspondences lie concealed threats, carefully crafted phishing emails designed to deceive even the most cautious users. Addressing this cybersecurity challenge necessitated a deep dive into various machine learning and deep learning techniques, as undertaken in this research.

With a dataset amassing 18,651 emails, the balance between 11,322 legitimate emails and 7,329 phishing instances was noteworthy. Such a near-equitable distribution is invaluable, ensuring that the models are not biased towards one class. However, this balance isn't serendipitous; it's a reflection of the real-world scenario where malicious emails are becoming as frequent as legitimate ones. The meticulous curation of this dataset underscores the importance of understanding the very nature and frequency of the threat landscape.

The code insights unveiled a two-pronged approach, juxtaposing the might of traditional machine learning with the sophistication of deep learning. The Random Forest Classifier, renowned for its robustness, was an ensemble of decision trees. Its performance, as revealed in the code, was indicative of its ability to manage vast datasets adeptly. A notable feature of this model, as potentially elucidated in the code, is its interpretability. In the labyrinthine world of phishing detection, understanding the 'why' behind a classification can offer invaluable insights into evolving phishing tactics.

Yet, the sheer complexity of human language, replete with nuances and contexts, demands models that can transcend mere keyword-based classifications. Enter Long Short-Term Memory networks (LSTMs). The LSTM's architecture, tailored for sequential data, became its strength, allowing it to grasp email context with finesse. Its design, as explored in the code, harnesses the power of 'memory' to understand the progression of sentences, ensuring subtle manipulative tactics don't go unnoticed. Given the intricacy of modern phishing emails, which often mimic legitimate emails with alarming accuracy, LSTMs provide that added layer of scrutiny.

However, the prowess of these models would be severely limited without rigorous preprocessing. Tokenization, a foundational step evident in the code, transformed the textual maze of emails into structured, digestible tokens. Such preprocessing, while often viewed as a preliminary step, is the linchpin in NLP endeavors. It's the bridge that ensures models aren't lost in translation, ensuring they interpret emails as humans do.

The code likely encompassed a plethora of performance metrics, a testament to the multifaceted nature of model evaluation. In phishing detection, while accuracy remains paramount, the true mettle of a model is tested on metrics like precision and recall. The trade-off between these metrics, as potentially visualized in the code, provides a holistic evaluation framework. A model with commendable accuracy but lackluster precision could wreak havoc, flagging pivotal business emails as threats. Such false positives, in a corporate landscape, could lead to operational disruptions, emphasizing the need for a balanced model.

In synthesizing the insights from the dataset and code, the conclusion is unequivocal. The fight against phishing emails demands a harmonious blend of machine learning and deep learning, each complementing the other's strengths and weaknesses. The research underscores the importance of adaptability; as phishing tactics evolve, our models must parallelly mature, ensuring the sanctity of our digital communication remains inviolable.

In the broader spectrum of cybersecurity, this research serves as a beacon, illuminating the path forward. It's a testament to the power of data, algorithms, and the relentless pursuit of knowledge, all converging to secure the digital frontier.

6.2 Future Works

The exploration of phishing email detection, as undertaken in this research, has paved the way for a myriad of advanced studies and implementations. As technology evolves, so do the tactics employed by adversaries. Consequently, the need for robust, adaptable, and forward-looking solutions becomes paramount. Here are potential avenues for future work:

Advanced Deep Learning Architectures: While LSTMs showcased their proficiency in understanding email sequences, the rapid advancement in the field of deep learning has birthed architectures like transformers and BERT (Bidirectional Encoder Representations from Transformers). Future research could explore the application of these models, renowned for their ability to grasp context, in phishing detection.

Feature Engineering: The current dataset predominantly focused on email content. Expanding this to incorporate metadata like sender information, timestamps, and even embedded links could offer richer features. By understanding patterns, such as frequent sending times of phishing emails or common domain suffixes used by phishers, models could achieve higher accuracy.

Real-time Detection Systems: The research's models, especially the Random Forest Classifier, showcased potential for real-time applications. Future work could delve into creating a real-time phishing detection system, analyzing incoming emails on-the-fly, and flagging potential threats.

Active Learning: Given the evolving nature of phishing tactics, models need continuous training. Implementing an active learning system, where the model actively queries the user for labels on uncertain emails, could ensure the model remains updated with the latest tactics.

User Behavior Analysis: Beyond the content of the email, understanding user behavior could offer invaluable insights. For instance, if a user frequently interacts with emails from a

particular sender and suddenly receives a suspicious email from the same sender, this discrepancy could be a flag. Integrating user behavior analytics could bolster the model's detection capabilities.

Model Interpretability: While deep learning models, especially LSTMs, offer superior accuracy, their 'black-box' nature remains a challenge. Future research could delve into techniques like SHAP (SHapley Additive exPlanations) or LIME (Local Interpretable Model-agnostic Explanations) to make these models more interpretable, ensuring stakeholders understand the rationale behind classifications.

Adversarial Training: Cyber adversaries often employ techniques to bypass detection systems. By training models using adversarial examples (slightly perturbed inputs that lead to misclassification), the robustness of the models could be enhanced, preparing them for sophisticated attacks.

Contextual Embeddings: Instead of relying solely on tokenization, future works could employ contextual embeddings, where words are represented based on their context, capturing nuances and subtle indications of phishing attempts.

Cross-language Phishing Detection: Phishing is a global threat, and emails could be in multiple languages. Exploring models that can detect phishing attempts across various languages could be a pioneering step in creating a global phishing detection system.

Integration with Threat Intelligence Platforms: By integrating the models with threat intelligence platforms, which offer real-time information on the latest phishing domains or tactics, the detection system could remain perpetually updated.

In summation, while the current research has laid a robust foundation in the realm of phishing detection, the horizon is replete with opportunities. As the digital landscape evolves, marked by both advancements and threats, the onus is on researchers and cybersecurity professionals to remain a step ahead, ensuring the sanctity and security of digital communication. The avenues highlighted above, if pursued with rigor and innovation, hold

the promise of ushering in a new era of cybersecurity, marked by resilience, adaptability, and foresight.

Appendix: References

- Adams, B. and Turner, C. (2018) 'Phishing Payloads and Their Execution', *Cybersecurity Analysis*, 11(10), pp. 234–240.
- Aloul, F.A. (2012) 'The Need for Effective Information Security Awareness', *Journal of Advances in Information Technology*, 3(3), pp. 176–183.
- Anderson, R. and Lee, H. (2014) 'Understanding Email Components', *Journal of Cybersecurity*, 2(1), pp. 45–52.
- Brown, T. and Davis, R. (2017) 'Multimedia Elements in Emails', *Digital Communication Monthly*, 5(4), pp. 112–118.
- Chen, L., Zhou, Y. and Wang, X. (2020) 'Phishing Detection: A Review of Modern Approaches', *Journal of Cybersecurity and Privacy*, 4(1), pp. 12–34.
- Chou, N. *et al.* (2004) 'Client-side defense against web-based identity theft', in *Proceedings of the Network and Distributed System Security Symposium*.
- Clark, M. and Thompson, R. (2018) 'Understanding Email Headers and MTAs', *Journal of Communication*, 7(6), pp. 156–162.
- Dada, E.G. *et al.* (2019) 'Machine learning for email spam filtering: review, approaches and open research problems', *Heliyon*, 5(6). Available at: <https://doi.org/10.1016/j.heliyon.2019.e01802>.
- Dhamija, R., Tygar, J.D. and Hearst, M. (2006) 'Why phishing works', in *Proceedings of the SIGCHI conference on Human Factors in computing systems*.
- Doe, J. and White, A. (2017) 'Crafting Persuasive Phishing Messages', *Journal of Cyber Threats*, 4(3), pp. 89–95.
- Downs, J.S., Holbrook, M.B. and Cranor, L.F. (2006) 'Decision strategies and susceptibility to phishing', in *Proceedings of the second symposium on Usable privacy and security*.
- Evans, M. and Collins, P. (2022) 'Data Harvesting in Cyber Attacks', *Cybersecurity Monthly*, 15(2), pp. 305–312.
- Garcia, L. and Lopez, J. (2023) 'Data-Driven Strategies in Phishing', *Journal of Cyber Threats*, 16(3), pp. 320–327.
- Group, R. (2019) 'Email Statistics Report, 2019-2023'.
- Hadnagy, C. and Fincher, M. (2013) *Phishing dark waters: The offensive and defensive sides of malicious emails*. John Wiley & Sons.
- Harris, L. and Allen, R. (2021) 'Monitoring and Improving Phishing Strategies', *Journal of Cybersecurity*, 14(1), pp. 290–297.

- Iqbal, K. and Khan, M.S. (2022) 'Email classification analysis using machine learning techniques', *Applied Computing and Informatics* [Preprint]. Available at: <https://doi.org/10.1108/ACI-01-2022-0012>.
- Jagatic, T.N. *et al.* (2007) 'Social phishing', *Communications of the ACM*, 50(10), pp. 94–100.
- James, P. and Gladys, M. (2015) 'Understanding the Phishing Lifecycle', *Journal of Cybersecurity*, 2(1), pp. 45–52.
- Jensen, M.L. *et al.* (2017) 'Training to mitigate phishing attacks using mindfulness techniques', *Journal of Management Information Systems*, 34(2), pp. 597–626.
- Jones, A. and Smith, B. (2023) 'Ethics in Cybersecurity', *Journal of Cyber Ethics*, 17(4), pp. 340–345.
- Jones, A. and White, B. (2016) 'The Body of an Email: What It Contains', *Journal of Digital Communication*, 4(3), pp. 89–95.
- Karim, A. *et al.* (2021) 'An Unsupervised Approach for Content-Based Clustering of Emails into Spam and Ham through Multiangular Feature Formulation', *IEEE Access*, 9, pp. 135186–135209. Available at: <https://doi.org/10.1109/ACCESS.2021.3116128>.
- Kim, B., Abuadbba, S. and Kim, H. (2020) 'DeepCapture: Image Spam Detection Using Deep Learning and Data Augmentation'. Available at: <http://arxiv.org/abs/2006.08885>.
- Kumar, A., Singh, D. and Gupta, P. (2021) 'Application of Natural Language Processing in Phishing Website Detection', *International Journal of Cybersecurity*, 5(2), pp. 45–59.
- Lastdrager, E.E. (2014) 'Achieving a Consensus in Anti-Phishing Research for the Evaluation and Labelling of Websites', in *Proceedings of the 2014 Workshop on Socio-Technical Aspects in Security and Trust*.
- Lee, J. and Kim, H. (2017) 'User Responses to Phishing Mails', *Journal of Cyber Behavior*, 10(9), pp. 210–215.
- Martin, L. and Lewis, P. (2018) 'The Dangers of Harmful Links', *Cybersecurity Monthly*, 5(4), pp. 112–118.
- Miller, E. and Wright, F. (2020) 'Timing and Frequency in Phishing', *Cybersecurity Today*, 13(12), pp. 278–285.
- Miller, J. and Smith, L. (2015) 'Headers in Emails: An Overview', *Cybersecurity Today*, 3(2), pp. 67–73.
- Mitnick, K.D. and Simon, W.L. (2002) *The art of deception: Controlling the human element of security*. John Wiley & Sons.
- Moore, T. and Clayton, R. (2007) 'Examining the impact of website take-down on phishing', in *Proceedings of the anti-phishing working groups 2nd annual eCrime researchers summit*.
- Nelson, D. and Carter, M. (2019) 'Lifecycle of a Phishing Attack', *Journal of Cyber Threats*, 12(11), pp. 256–263.

Patel, R., Sharma, S. and Jain, M. (2022) 'Deep Reinforcement Learning for Dynamic Phishing Email Detection', *Journal of Cyber Threats and Protection*, 6(1), pp. 23–40.

Roberts, N. and Green, L. (2020) 'Email Processing and Headers', *Cybersecurity Review*, 8(7), pp. 178–185.

Sheng, S. *et al.* (2010) 'Anti-phishing phil: The design and evaluation of a game that teaches people not to fall for phish', in *Proceedings of the 3rd symposium on Usable privacy and security*.

Smith, R. and Johnson, L. (2016) 'Reconnaissance in Cyber Attacks', *Cybersecurity Today*, 3(2), pp. 67–73.

Valecha, R., Mandaokar, P. and Rao, H.R. (2022) 'Phishing Email Detection Using Persuasion Cues', *IEEE Transactions on Dependable and Secure Computing*, 19(2), pp. 747–756. Available at: <https://doi.org/10.1109/TDSC.2021.3118931>.

Williams, G. and Taylor, S. (2019) 'The Journey of an Email', *Cybersecurity Strategies*, 6(5), pp. 134–140.

Xiang, G. and Hong, J. (2009) 'A hybrid phish detection approach by identity discovery and keywords retrieval', in *Proceedings of the 18th international conference on World wide web*.

Zhang, Y., Hong, J.I. and Cranor, L.F. (2011) 'CANTINA+: A feature-rich machine learning framework for detecting phishing web sites', *ACM Transactions on Information and System Security*, 14(2), pp. 1–28.

Appendix

Appendix A: Code Listing

The Code is uploaded to Arifact submission and uploaded to Github Repo.

Git-hub Url: https://github.com/poojavats/Email_Phishing