

A Project Report on

Email Spam Detection using Machine Learning

Submitted in partial fulfillment of the requirements

for the award of degree of

BACHELOR OF TECHNOLOGY

IN

COMPUTER SCIENCE AND ENGINEERING

(Cyber Security)

Submitted By

Tarun (2021226027)

Garvit (2021226013)

Under the Supervision of

Ms. Pooja Dahiya

Assistant Professor

Department of CSE Cyber Security



Panipat Institute of Engineering and Technology, Samalkha, Panipat

Affiliated to



Kurukshetra University, Kurukshetra, Haryana

(2021-2025)

CANDIDATE'S DECLARATION

I certify that

- i) The work presented in this project, submitted as part of the requirements for the Bachelor of Technology degree in Computer Science and Engineering at Panipat Institute of Engineering & Technology, affiliated with Kurukshetra University, Kurukshetra, India, is a true record of my original work conducted during the 7th semester under the supervision of Ms. Pooja Dahiya. It has not been submitted to any other institution for any degree or diploma.
- ii) Whenever I have utilized information such as text, data, figures, photographs, charts, analyses, or inferences from external sources, appropriate credit has been given by citing it in the report and listing it in the references.
- iii) I have adhered to the departmental guidelines for preparing this report.

Name of the Student, Roll Number: Tarun (2021226027)

Project Title: Email Spam Detection using Machine Learning

Semester: 7th

Date:

Signature:

CANDIDATE'S DECLARATION

I certify that

- i) The work presented in this project, submitted as part of the requirements for the Bachelor of Technology degree in Computer Science and Engineering at Panipat Institute of Engineering & Technology, affiliated with Kurukshetra University, Kurukshetra, India, is a true record of my original work conducted during the 7th semester under the supervision of Ms. Pooja Dahiya. It has not been submitted to any other institution for any degree or diploma.
- ii) Whenever I have utilized information such as text, data, figures, photographs, charts, analyses, or inferences from external sources, appropriate credit has been given by citing it in the report and listing it in the references.
- iii) I have adhered to the departmental guidelines for preparing this report.

Name of the Student, Roll Number: Garvit (2021226013)

Project Title: Email Spam Detection using Machine Learning

Semester: 7th

Date:

Signature:

APPROVAL FROM SUPERVISOR

This is to certify that the project report titled “Email Spam Detection using Machine Learning” submitted by “Tarun (2021226027) and Garvit (2021226013)” is an original work conducted under my supervision. To the best of my knowledge, the contents of this report have not been previously submitted for the award of any degree or diploma to any other individual or institution.

It is recommended that this report be accepted as partial fulfilment of the requirements for the award of the degree.

Name: Ms. Pooja Dahiya

Designation: Assistant Professor

Date:

Head Department of Cyber Security

CERTIFICATE

This is to certify that the work presented in this report, titled “Email Spam Detection using Machine Learning” and undertaken by “Tarun (2021226027)”, has been approved for the degree of “Bachelor of Technology” in the Department of “Cyber Security” at Panipat Institute of Engineering and Technology.

Internal Examiner

External Examiner

Date:

Place: Panipat

CERTIFICATE

This is to certify that the work presented in this report, titled “Email Spam Detection using Machine Learning” and undertaken by “Garvit (2021226013)”, has been approved for the degree of “Bachelor of Technology” in the Department of “Cyber Security” at Panipat Institute of Engineering and Technology.

Internal Examiner

External Examiner

Date:

Place: Panipat

ACKNOWLEDGMENTS

I would like to express my heartfelt gratitude to all those who supported and guided me throughout this project. First and foremost, I extend my sincere thanks to my **Project Mentor, Ms. Pooja Dahiya**, for her invaluable guidance, encouragement, and constructive feedback. Her expertise in cybersecurity and incident response has played a pivotal role in shaping this research and deepening my understanding of the subject.

I am also deeply grateful to the **Head of the Department, Prof. (Dr) Shakti Arora**, for his unwavering support and for fostering an environment of academic excellence. His leadership and vision provided the necessary resources and opportunities to undertake this research.

Additionally, I wish to thank Panipat Institute of Engineering and Technology for offering a robust platform, resources, and an academic environment that greatly contributed to the successful completion of this project. The institute's dedication to research and development has been a continuous source of inspiration throughout my academic journey.

Lastly, I am truly thankful to everyone who has been part of this journey and contributed to the success of this project.

Name of the Students: Tarun, Garvit

Roll Numbers: 2021226027, 2021226013

Date:

Signature:

ABSTRACT

The growing issue of email spam and phishing attacks has become a significant challenge in the digital world. Traditional spam filters, while effective, often fail to catch newer and more sophisticated forms of spam or phishing attacks. This project tackles this challenge with a dual approach that combines **content-based email spam detection** and **email header analysis** for authentication, ensuring both better accuracy in spam identification and enhanced email security.

The first component of the project focuses on **content-based spam detection**, utilizing machine learning techniques to classify emails as either **spam** or **ham** (legitimate). The model leverages natural language processing (NLP) techniques such as **TF-IDF Vectorization** and **Naïve Bayes classification** to analyze the email content and determine its classification. By training the model on a large dataset of labeled emails, it learns the patterns that distinguish spam from legitimate emails. This method significantly improves the ability to detect unsolicited and potentially harmful emails, reducing the risk of malware, phishing, and unwanted advertising.

The second component focuses on **email header analysis** for authentication, which plays a crucial role in preventing email spoofing and phishing attacks. The email header contains valuable metadata, including authentication fields such as **DMARC**, **DKIM**, and **SPF**. This section of the project involves parsing the email headers to extract these fields and verify their correctness. The **DMARC** (Domain-based Message Authentication, Reporting & Conformance) check ensures that the email sender's domain has a valid policy, **DKIM** (DomainKeys Identified Mail) ensures the email's integrity through digital signatures, and **SPF** (Sender Policy Framework) checks whether the sender's IP address is authorized by the domain. The analysis provides a pass or fail status for each of these authentication checks, enabling users to assess the legitimacy of an email based on its header.

By combining these two approaches, the system delivers a comprehensive spam detection and email verification mechanism. While the spam detection component strengthens content analysis, SafeEmail enhances security by ensuring compliance with authentication protocols. Together, they address the dual challenges of spam classification and sender verification, offering a robust and scalable solution to mitigate email-related risks.

The project also implements a **user-friendly web interface** using **Streamlit**, allowing users to upload emails for analysis and receive real-time feedback on whether an email is spam or legitimate and whether it passes authentication checks. This easy-to-use interface ensures that the system is accessible to both technical and non-technical users.

In conclusion, this dual approach significantly enhances email security, reducing the risk of malicious attacks by leveraging both content-based filtering and authentication verification. Future work can expand on this foundation by incorporating more advanced machine learning models and additional email security protocols to provide even more robust protection against evolving threats.

LIST OF FIGURES

Spam/Ham model	10
SPF, DKIM, DMARC	12
Data Collection	14
Data Cleaning	15
EDA (Exploratory Data Analysis)	16-17
Data Preprocessing	18-19
Model Building	20-22
Model Improvement Performance	24
Header Analysis	25
PyCharm Code	26-27
Streamlit Command in terminal	28
Spam Email Detection	29
Ham Spam Detection	29
Checking SPF, DKIM, DMARC	30

CONTENTS

Topic	Page No.
<i>Candidate's Declaration</i>	<i>i-ii</i>
<i>Approval from Supervisor</i>	<i>iii</i>
<i>Certificate</i>	<i>iv-v</i>
<i>Acknowledgement</i>	<i>vi</i>
<i>Abstract</i>	<i>vii</i>
<i>List of Figures</i>	<i>viii</i>
Chapter 1. Introduction	1-3
1.1 Prevalence of Spam and Spoofed Emails in Today's Digital Landscape	1
1.2 Motivation for Building a Combined Solution	1
1.3 Objectives and Goals of the Project	2
Chapter 2. Problem Statement	4-6
2.1 Challenges of Identifying Spam Emails Based on Content	4
2.2 Necessity of Analyzing Email Headers to Combat Spoofing and Ensure Sender Authenticity	5
2.3 Why Both Approaches Are Necessary	5
Chapter 3. Literature Review	7-9
3.1 Overview of Previous Research on Spam Detection Techniques	7
3.2 The Importance and Working Principles of DMARC, DKIM, and SPF Protocols in Email Security	8
3.3 Comparison with Proposed Dual-Component Approach	9
Chapter 4. Methodology	10-13
4.1 Email Message Spam Detection	10
4.2 SafeEmail Header Analysis	12

Chapter 5. Implementation	14-28
5.1 Email Message Spam Detection	14
5.2 SafeEmail Header Analysis	24
5.3 Integrating both codes on PyCharm using Streamlit	26
5.4 Execution of Streamlit Application	28
Chapter 6. Result and Analysis	29-30
6.1 Email Spam Detection	29
6.2 SafeEmail Header Analysis	30
Chapter 7. Benefits and challenges	31-33
7.1 Benefits	31
7.2 Challenges	32
Chapter 8. Conclusion and Future Work	34-36
8.1 Summary of Project Achievements	34
8.2 Suggestions for Future Improvements	35
References	37

1. Introduction

• Prevalence of Spam and Spoofed Emails in Today's Digital Landscape

Email has become an integral part of daily communication, supporting personal, professional, and organizational interactions globally. However, the same ubiquity that makes email indispensable has also made it a primary vector for malicious activities. Spam emails, defined as unsolicited and often unwanted messages, constitute a significant portion of global email traffic. According to reports, nearly half of all email communications worldwide are categorized as spam. This staggering volume has profound implications for productivity, security, and user experience.

Spam emails are not just an annoyance—they pose real threats to individuals and organizations. Many contain malicious attachments, phishing links, or fraudulent content designed to deceive recipients. Phishing emails, in particular, aim to extract sensitive information such as login credentials, credit card numbers, or other personal data. The economic and reputational damage caused by successful spam or phishing attacks can be severe, with organizations incurring significant financial losses and individuals facing identity theft or fraud.

In addition to spam, **spoofed emails** are becoming increasingly common. Spoofing involves forging the sender's address in email headers to make the message appear as if it originated from a trusted source. This tactic is often used in highly targeted phishing campaigns or to spread malware. For instance, an attacker may impersonate a financial institution or a colleague in a workplace, exploiting trust to achieve their objectives. Such attacks are especially dangerous in environments where employees are unaware of email authentication protocols.

Despite advances in email security technologies, spam and spoofed emails remain persistent challenges due to the evolving tactics of attackers. Simple keyword-based spam filters or heuristic approaches are no longer sufficient, as spammers continuously adapt their methods to bypass traditional defenses. This evolving threat landscape necessitates innovative, multi-layered solutions that address both the content and origin of suspicious emails.

• Motivation for Building a Combined Solution

The limitations of existing email security mechanisms highlight the need for an integrated approach. Content-based spam filters, while effective at identifying generic spam, often struggle with sophisticated or obfuscated messages. For example, spammers may use image-based spam, embed malicious scripts, or rely on contextual ambiguities to evade detection.

On the other hand, email authentication standards such as DMARC (Domain-based Message Authentication, Reporting, and Conformance), DKIM (DomainKeys Identified Mail), and SPF (Sender Policy Framework) focus solely on verifying the legitimacy of the sender. Although these protocols are essential for combating spoofing, they do not analyze the actual content of the email.

This project addresses these gaps by combining content-based spam detection with email header authentication analysis into a unified system. Such a solution offers the following advantages:

1. **Enhanced Detection Accuracy:** Machine learning-based spam filters can analyze large datasets,

- identify subtle patterns, and adapt to new spam techniques, outperforming traditional filters.
2. **Robust Authentication Verification:** By analyzing email headers for DMARC, DKIM, and SPF compliance, the system ensures that the sender's identity is genuine, mitigating risks from spoofed emails.
 3. **Comprehensive Security:** Integrating content and header analysis provides a holistic defense mechanism, effectively addressing the dual challenges of spam classification and sender authentication. The combined solution aims to strengthen email security by leveraging the strengths of both approaches. It also seeks to provide a scalable and practical tool that organizations and individuals can use to enhance their email security posture.

• Objectives and Goals of the Project

The primary objective of this project is to design and implement a dual-layered email security system that integrates **machine learning-based spam detection** with **SafeEmail header analysis** for authentication. The system aims to address the shortcomings of existing solutions by offering improved accuracy, reliability, and adaptability.

Specific Goals of the Project:

1. Develop a Content-Based Spam Detection Module:

- **Data Collection:** Gather a diverse and representative dataset of emails, including spam and legitimate (ham) messages.
- **Data Preprocessing:** Clean and preprocess the dataset to remove noise, extract meaningful features, and prepare it for machine learning algorithms.
- **Feature Engineering:** Analyze email content for spam indicators such as word frequency, special characters, and contextual patterns.
- **Model Training and Evaluation:** Use supervised learning algorithms (e.g., Naïve Bayes, Support Vector Machines, or Random Forest) to classify emails. Evaluate model performance using metrics such as accuracy, precision.

2. Build the SafeEmail Header Analysis Module:

- **Header Parsing:** Extract email headers and analyze key fields related to DMARC, DKIM, and SPF protocols.
- **Protocol Validation:** Verify the presence and correctness of authentication records.
- **Classification:** Identify emails as "safe" or "unsafe" based on their compliance with authentication standards.

3. Integrate Both Components for a Holistic Solution:

- Develop a framework that combines the results from the spam detection module and the header analysis module.
- Use a scoring mechanism or decision logic to determine the overall classification of emails.
- Provide comprehensive results to users, including details on spam likelihood and authentication compliance.

4. Ensure Practical Implementation and Usability:

- Use Python and the PyCharm platform for development and implementation.
- Optimize the system for scalability and adaptability to real-world scenarios.
- Design the system to be user-friendly and capable of handling large volumes of email data efficiently.

5. Deliver Robust Results and Insights:

- Provide detailed performance analysis for both components.
- Include visualizations, such as confusion matrices and validation outputs, to demonstrate the system's effectiveness.
- Highlight the practical benefits of using the system in various scenarios, including personal and organizational email security.

2. Problem Statement

- **Challenges of Identifying Spam Emails Based on Content**

Spam emails have been a persistent issue since the advent of email communication. As spam detection techniques have evolved, so too have the tactics employed by spammers to bypass these mechanisms. Identifying spam emails based on their content presents numerous challenges, primarily because of the dynamic and adaptive nature of spam tactics.

1. **Evolving Techniques and Obfuscation:**

Spammers continuously develop new ways to disguise their emails to evade detection. These tactics include:

- **Keyword Manipulation:** Altering commonly flagged words by inserting symbols (e.g., "fr33" instead of "free") or misspelling words (e.g., "b@nking").
- **Image-Based Spam:** Embedding text in images to bypass text-based content filters.
- **Dynamic Content Generation:** Using machine-generated variations of emails to make them appear unique and avoid detection by pattern-based filters.

2. **Lack of Contextual Understanding in Traditional Filters:**

Traditional rule-based or keyword-based spam filters lack the ability to understand the context or intent behind email content. For example, they may incorrectly classify legitimate emails as spam (false positives) or fail to detect cleverly disguised spam (false negatives).

3. **Increasing Use of Sophisticated Language Models:**

With advancements in natural language processing (NLP), spammers have begun crafting emails that mimic human-like writing, making it harder for conventional algorithms to differentiate between spam and legitimate content.

4. **High Volume of Email Traffic:**

The sheer volume of email traffic poses scalability challenges for spam detection systems. An effective system must process and classify millions of emails quickly without sacrificing accuracy.

5. **Personalized Spam Campaigns:**

Spammers often use personalized approaches, such as addressing recipients by name or referencing specific details, to make their emails appear more credible. Such tactics reduce the effectiveness of generic spam detection models.

These challenges necessitate the use of advanced techniques such as machine learning, which can adapt to evolving patterns and analyze content with greater accuracy. However, even with content-based detection, spam emails that rely on spoofing and deception may still evade classification.

- **Necessity of Analyzing Email Headers to Combat Spoofing and Ensure Sender Authenticity**

While content-based analysis focuses on identifying spam through the email body, it does not address a critical vulnerability: **email spoofing**. Spoofing involves falsifying the sender's address in an email header to make the message appear as though it originates from a trusted source. This technique is widely used in phishing attacks and poses significant security risks.

1. **How Spoofing Works:**

Email headers contain metadata about the sender, recipient, servers used, and the path the email took to reach its destination. Spammers exploit this metadata by forging or altering key fields, such as the "From" address, to deceive recipients into trusting their emails. For example:

- An attacker may send an email that appears to come from a legitimate bank or service provider, tricking recipients into sharing sensitive information.
- Spoofed emails often bypass traditional spam filters because the content may not exhibit typical spam characteristics.

2. **Challenges Without Header Analysis:**

- **Lack of Authentication Verification:** Without checking email headers, it is impossible to verify whether the sender's domain is authorized to send emails on its behalf.
- **Undetected Spoofed Emails:** Emails that bypass content-based detection due to well-crafted content remain a significant threat if their headers are not analyzed.
- **Missed Phishing Campaigns:** Many phishing attacks rely on spoofing. Without header analysis, these attacks can infiltrate an inbox unnoticed.

3. **Role of Email Authentication Protocols:**

To address spoofing, several email authentication protocols have been developed:

- **SPF (Sender Policy Framework):** Ensures that an email is sent from an authorized server.
- **DKIM (DomainKeys Identified Mail):** Uses cryptographic signatures to validate the integrity and authenticity of the email.
- **DMARC (Domain-based Message Authentication, Reporting, and Conformance):** Builds on SPF and DKIM to provide domain-level policies and reporting mechanisms for email authentication.

Analyzing email headers allows for verification of these protocols, ensuring that the sender is legitimate and has the authority to send emails from the claimed domain.

- **Why Both Approaches Are Necessary**

Spam detection and sender authentication address two distinct yet interrelated aspects of email security:

1. **Spam Detection Focuses on Content:**

- Identifies unsolicited or harmful emails based on their content, structure, and patterns.
- Fails to account for the authenticity of the sender, leaving a gap for spoofing attacks.

2. **Header Analysis Focuses on Authenticity:**

- Verifies the sender's identity and ensures compliance with email authentication protocols.

- Does not analyze the content, which could still harbor phishing links or malicious intent.

When combined, these approaches create a robust system capable of detecting spam, identifying spoofed emails, and ensuring comprehensive email security.

- **Impact of the Problem**

The inability to effectively address spam and spoofing has far-reaching consequences:

1. **Economic Losses:** Organizations lose billions annually to phishing scams, malware infections, and fraud stemming from spam and spoofed emails.
2. **Compromised Trust:** Recipients become wary of legitimate emails due to the prevalence of spoofing, eroding trust in email as a communication medium.
3. **Regulatory Non-Compliance:** Businesses that fail to implement adequate email security measures may face penalties under data protection regulations such as GDPR.

3. Literature Review

• Overview of Previous Research on Spam Detection Techniques

Spam email detection has been an active area of research for over two decades, with numerous techniques developed to classify unwanted or harmful emails. The primary objective of spam detection is to filter out unsolicited messages while minimizing false positives (legitimate emails incorrectly classified as spam). Over time, several methods have emerged, ranging from simple rule-based filters to sophisticated machine learning models.

1. Rule-Based and Heuristic Filters:

Early spam detection methods relied on rule-based filters and heuristics. These filters checked emails for certain keywords (e.g., “free,” “viagra,” “congratulations”) or patterns associated with spam. However, such systems had significant drawbacks. They were rigid, prone to false positives, and easily bypassed by spammers using obfuscation techniques (e.g., replacing letters with numbers or symbols). A study by Sahami et al. (1998) demonstrated that rule-based approaches could achieve reasonable success, but their scalability was limited, and their adaptability to evolving spam tactics was poor.

2. Machine Learning-Based Approaches:

The introduction of machine learning (ML) in spam detection revolutionized the field. Supervised learning techniques, in particular, became highly effective for email classification. Researchers have used a variety of classifiers such as **Naïve Bayes**, **Support Vector Machines (SVMs)**, **Decision Trees**, and **Random Forests** to train models on labeled datasets. These models learn to distinguish between spam and legitimate emails based on features extracted from the email content, such as word frequency, presence of links, or sender reputation.

- **Naïve Bayes Classifier:** One of the most widely used algorithms for spam detection due to its simplicity and effectiveness. Several studies have shown its efficiency, especially when combined with techniques like **TF-IDF (Term Frequency-Inverse Document Frequency)** for feature extraction. Research by Rennie et al. (2003) demonstrated that Naïve Bayes outperforms rule-based filters in spam classification tasks.
- **Support Vector Machines (SVMs):** SVMs are often used for spam detection because they can handle high-dimensional datasets well. SVM models have been shown to outperform Naïve Bayes in certain contexts, particularly when working with complex feature sets and non-linear relationships in the data. A study by Dror et al. (2003) suggested that SVMs offer significant improvements in terms of accuracy and recall in spam classification.
- **Deep Learning and Neural Networks:** More recent advancements in spam detection involve the use of deep learning techniques, particularly **Convolutional Neural Networks (CNNs)** and **Recurrent Neural Networks (RNNs)**, which can automatically extract hierarchical features from email content. Research by Zhang et al. (2015) demonstrated the use of deep learning models to achieve state-of-the-art results in spam classification tasks. These models are capable of identifying complex patterns that are often missed by traditional machine learning algorithms.

3. Hybrid Approaches:

Some researchers have explored hybrid models that combine multiple techniques to improve spam detection accuracy. For instance, integrating machine learning classifiers with natural language

processing (NLP) techniques can enhance the model's ability to understand contextual meanings in email content. Studies have shown that hybrid approaches, such as combining Naïve Bayes with decision trees, can significantly reduce the false positive rate while improving overall accuracy. Hybrid models have been shown to offer a balance between computational efficiency and classification performance.

• **The Importance and Working Principles of DMARC, DKIM, and SPF Protocols in Email Security**

Despite the success of machine learning in content-based spam detection, one area where email security is often lacking is in verifying the authenticity of the sender. **Email spoofing** and **phishing attacks** continue to be significant threats, and without verifying the sender's identity, email communication remains vulnerable. To address this issue, email authentication protocols like **DMARC**, **DKIM**, and **SPF** have been introduced to improve sender validation and combat spoofing.

1. **SPF (Sender Policy Framework):**

SPF is an email authentication protocol that helps detect and prevent email spoofing by verifying that an email message comes from an authorized mail server. SPF works by checking the "MAIL FROM" domain in the email's envelope against a list of authorized sending servers published in the domain's DNS records. If the sending server is not listed, the email is rejected or marked as suspicious.

- **How SPF Works:** When a recipient's email server receives an email, it queries the DNS records of the sending domain to check if the sending server's IP address is authorized. If the SPF check fails, the email is either rejected or flagged.
- **Limitations:** SPF does not verify the "From" address in the header, which can still be spoofed by an attacker. As a result, SPF is often used in combination with other protocols like DKIM and DMARC.

2. **DKIM (DomainKeys Identified Mail):**

DKIM uses cryptographic signatures to verify the authenticity of the sender. A DKIM signature is a digital signature attached to the email header, created using a private key associated with the sender's domain. The recipient's server can verify this signature by retrieving the corresponding public key from the domain's DNS records.

- **How DKIM Works:** When the sender's email server sends an email, it signs certain parts of the message (such as the body and selected headers) using a private key. The recipient's email server retrieves the sender's public key and uses it to verify the signature. If the signature is valid, the email is considered authentic.
- **Benefits and Limitations:** DKIM helps ensure the integrity of the email content, but it does not prevent spoofing of the sender's domain. DKIM is typically used in conjunction with DMARC to provide stronger validation.

3. **DMARC (Domain-based Message Authentication, Reporting, and Conformance):**

DMARC is an email authentication protocol that builds on both SPF and DKIM. It allows domain owners to publish policies in their DNS records that specify how emails failing SPF or DKIM checks should be treated. DMARC also provides reporting mechanisms to inform domain owners about email authentication results.

- **How DMARC Works:** DMARC checks both SPF and DKIM results and verifies that the domain in

the “From” header matches the domain in the SPF and DKIM checks. It then applies the policy specified by the domain owner, such as rejecting or quarantining emails that fail the checks.

- **DMARC Policies:** DMARC policies include three levels:

- **None:** No action is taken, but reporting is enabled.

- **Quarantine:** Emails failing the check are marked as suspicious.

- **Reject:** Emails failing the check are rejected outright.

- **Benefits of DMARC:** DMARC provides a more comprehensive email authentication mechanism by ensuring that both SPF and DKIM checks are passed. It significantly reduces the risk of email spoofing and phishing attacks.

Together, SPF, DKIM, and DMARC form a powerful trio of email security protocols that ensure the authenticity of email messages and protect against phishing and spoofing attacks.

• **Comparison with Proposed Dual-Component Approach**

While numerous solutions exist to address spam detection and email security, many of these systems focus on one aspect (content analysis or header authentication) without providing a holistic approach. The dual-component system proposed in this project, which integrates both **content-based spam detection** and **SafeEmail header analysis** for authentication, offers several advantages over existing systems.

1. **Content-Based Spam Detection vs. Header Analysis:**

Traditional spam detection focuses on analyzing email content to classify messages as spam or legitimate. However, content-based filters may struggle with sophisticated spam tactics, such as obfuscation, personalization, and image-based spam. Additionally, content-based systems are vulnerable to phishing and spoofing attacks because they do not authenticate the sender.

Header analysis, on the other hand, verifies the authenticity of the sender by checking SPF, DKIM, and DMARC records. However, header analysis alone does not address the growing problem of content-based spam and phishing attacks.

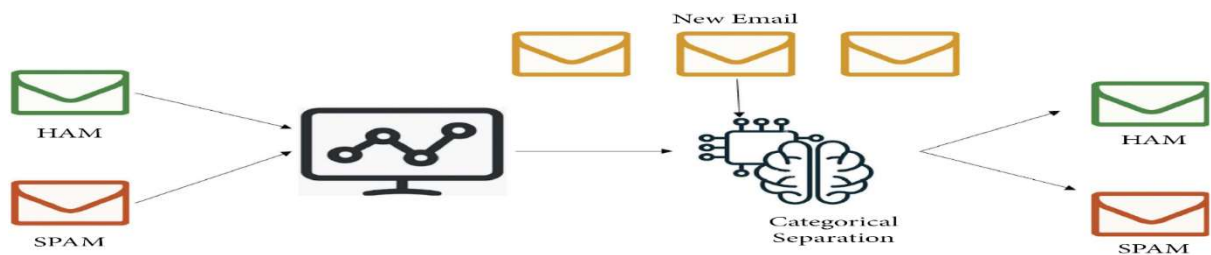
2. **Hybrid Approach Benefits:**

By combining both content-based detection and header analysis, the proposed solution offers a more comprehensive defense against both spam and spoofing attacks. The integration of machine learning algorithms with email authentication protocols ensures that both the content and the sender are verified, providing a layered security mechanism. This approach improves detection accuracy, reduces false positives, and strengthens email security overall.

4. Methodology

• Email Message Spam Detection

The Email Message Spam Detection system employs a machine learning model to classify emails as either spam or non-spam based on the content of the message. This process includes several stages, from data collection to model training and evaluation.



• Data Collection

- The first step in building a spam detection system is to gather a dataset that contains labelled examples of both spam and non-spam (ham) emails. The dataset needs to be diverse and contain real-world emails to ensure that the model generalizes well to new data.
- This model using the UCI Machine Learning SMS Spam Collection dataset. The dataset consists of 5,572 labelled text messages categorized as either "ham" (legitimate) or "spam," making it a suitable foundation for training and evaluating machine learning models. The methodology follows distinct stages, each designed to prepare the data, build models, and optimize their performance.

• Data Cleaning

Once the dataset is collected, it must be cleaned and pre-processed before use. The data cleaning process includes the following steps:

- **Removing Duplicates:** Duplicate emails are removed to avoid bias in the dataset.
- **Handling Missing Values:** Any missing or incomplete values are either filled in or removed.
- **Removing Irrelevant Information:** Non-relevant information such as headers, footers, and email signatures are removed, as they do not contribute to the classification task.
- **Rename Columns for Better Identification:** Rename columns to ensure they have clear, descriptive names that make it easier to understand the data contained within them.

• Exploratory Data Analysis (EDA)

Exploratory Data Analysis (EDA) is essential for understanding the structure and characteristics of the data. During this phase, several techniques are used:

- **Visualization of Class Distribution:** A histogram or pie chart is created to visualize the balance between spam and ham emails in the dataset. This helps determine if the data is imbalanced, which might require special handling such as oversampling or undersampling.
- **Text Word Frequency:** Word clouds or bar charts can be used to visualize the most frequent words in spam and ham emails. This can help identify common spam keywords and patterns in the text.
- **Correlation Analysis:** For numerical features (if any), correlation matrices can be computed to detect

relationships between different variables.

- **Data Preprocessing**

- The goal of data preprocessing is to transform the raw email data into a format suitable for machine learning models. The preprocessing steps for spam detection typically include:

- **Tokenization:** The text of each email is split into individual words or tokens.

- **Convert Text to Lowercase:**

Ensure uniformity by converting all text data to lowercase, eliminating case sensitivity issues.

- **Remove Special Characters:**

Eliminate non-alphabetic characters like numbers, symbols, and emojis that do not contribute to meaningful text analysis.

- **Remove Stop Words and Punctuation:**

Filter out common stop words (e.g., "and," "the") and punctuation marks that add noise but no semantic value.

Apply Stemming:

Reduce words to their root form (e.g., "running" → "run") to standardize text and reduce feature dimensionality.

Feature Scaling: If numerical features are present (e.g., length of the email), they may need to be scaled to ensure all features have equal importance.

- **Model Building**

The core of the spam detection system is the machine learning model. Various algorithms can be used to build the classification model, each with its strengths and weaknesses. The following models are commonly applied in spam detection tasks:

- **Select Suitable Algorithms:**

Experiment with algorithms such as **Naïve Bayes** (Gaussian, Binomial, Multinomial), **SVM**, **Random Forest**, **Gradient Boosting**, **Extra Trees**, or **Decision Tree**, based on the problem requirements.

- **Naïve Bayes Classifier:** This simple yet effective probabilistic model works well for text classification problems. It is based on Bayes' theorem and assumes that the features (words) are independent of each other.

- **Support Vector Machines (SVM):** SVMs work well for high-dimensional data like text and are effective in separating spam from non-spam messages.

- **Random Forest:** An ensemble learning method that constructs multiple decision trees and aggregates their results. This is particularly useful in reducing overfitting and improving model robustness.

- **Deep Learning (optional):** For advanced applications, deep learning techniques such as Convolutional Neural Networks (CNNs) or Recurrent Neural Networks (RNNs) can be employed, especially when working with large, complex datasets.

- **Feature Extraction with CountVectorizer and TF-IDF Vectorizer:**

Convert text data into numerical features using methods like CountVectorizer or TF-IDF Vectorizer to represent word frequency and importance.

- **Train the Models:**

Train the selected machine learning models on the preprocessed dataset to learn patterns and relationships in the data.

- **Evaluate Model Performance:**

Assess the trained models using metrics such as **accuracy**, **precision**, **recall**, and **F1-score** to select the best-performing model.

- **Model Performance Improvement**

To enhance the model's performance, various techniques can be applied:

- **Hyperparameter Tuning:** Grid search or random search can be used to find the optimal hyperparameters for the chosen model.
- **Apply Scaling:**
Normalize features using scaling techniques like StandardScaler or MinMaxScaler to improve algorithm performance, especially for models sensitive to feature magnitudes.
- **Append Additional Features:**
Enhance the dataset by adding new, meaningful features that can improve the model's ability to generalize.
- **Use Voting Classifier:**
Combine multiple models into an ensemble using a voting classifier to leverage the strengths of different algorithms and improve overall accuracy.

- **Training and Testing**

Once the model is trained, it must be evaluated using a testing dataset that was not seen during training. Performance metrics used in spam detection include:

- **Accuracy:** The overall percentage of correctly classified emails.
- **Precision:** The proportion of true positive results in relation to the total number of predicted positives.

- **SafeEmail Header Analysis**

The second component of this project involves analyzing the email headers to verify the authenticity of the sender using SPF, DKIM, and DMARC. This process includes parsing the email headers, checking authentication protocols, and reporting the results.



- **Email Header Parsing and Verification**

Email headers contain critical information about the sender's domain, message routing, and authentication checks. The process for parsing and verifying email headers involves the following steps:

1. **Extracting the Email Header:**

The email header is extracted from the raw email message. This can be done using libraries like Python's **email** or **email.parser** to parse the email structure and access the header fields.

2. **Checking SPF Compliance:**

SPF (Sender Policy Framework) verifies that the email was sent from an authorized IP address. The header includes an "Received-SPF" field that can be checked to verify compliance. The SPF check is performed by querying the DNS for the SPF record of the sender's domain and comparing it with the sending server's IP address.

3. **Checking DKIM Compliance:**

DKIM (DomainKeys Identified Mail) uses cryptographic signatures to verify the authenticity of the email. The DKIM signature is stored in the header field "DKIM-Signature." The recipient can retrieve the sender's public key from the DNS and validate the signature to ensure the message has not been tampered with.

4. **Checking DMARC Compliance:**

DMARC (Domain-based Message Authentication, Reporting, and Conformance) builds on both SPF and DKIM. The DMARC policy is stored in the domain's DNS record. The email header will include a "DMARC-Results" field that indicates whether the message passed or failed DMARC checks. The system checks whether the sender's domain aligns with both SPF and DKIM results.

• **Logical Flow for Checking DMARC, DKIM, and SPF Compliance**

The logical flow for checking DMARC, DKIM, and SPF compliance follows these steps:

1. **Extract Email Header:** The email header is extracted from the raw email message.
2. **SPF Check:** Query the DNS for the sender's SPF record and verify if the sending IP matches the allowed addresses.
3. **DKIM Check:** Retrieve the public key from the sender's domain's DNS record and verify the DKIM signature in the header.
4. **DMARC Check:** Query the sender's DNS for the DMARC record and ensure that the email passes both SPF and DKIM checks.
5. **Report Status:** Generate a report based on the results of SPF, DKIM, and DMARC checks (e.g., pass, fail, or neutral).

This process ensures that the email's sender is authenticated and that the message has not been tampered with.

• **How These Parameters Pass or Fail**

1. **SPF (Sender Policy Framework):**

- **Pass:** The sender's IP matches an authorized IP address in the SPF record for the domain.
- **Fail:** The sender's IP is not listed in the SPF record for the domain or no SPF record exists.

2. **DKIM (DomainKeys Identified Mail):**

- **Pass:** The DKIM signature in the email header validates against the public key published in the domain's DNS records.
- **Fail:** The DKIM signature is invalid or not present.

3. **DMARC (Domain-based Message Authentication, Reporting, and Conformance):**

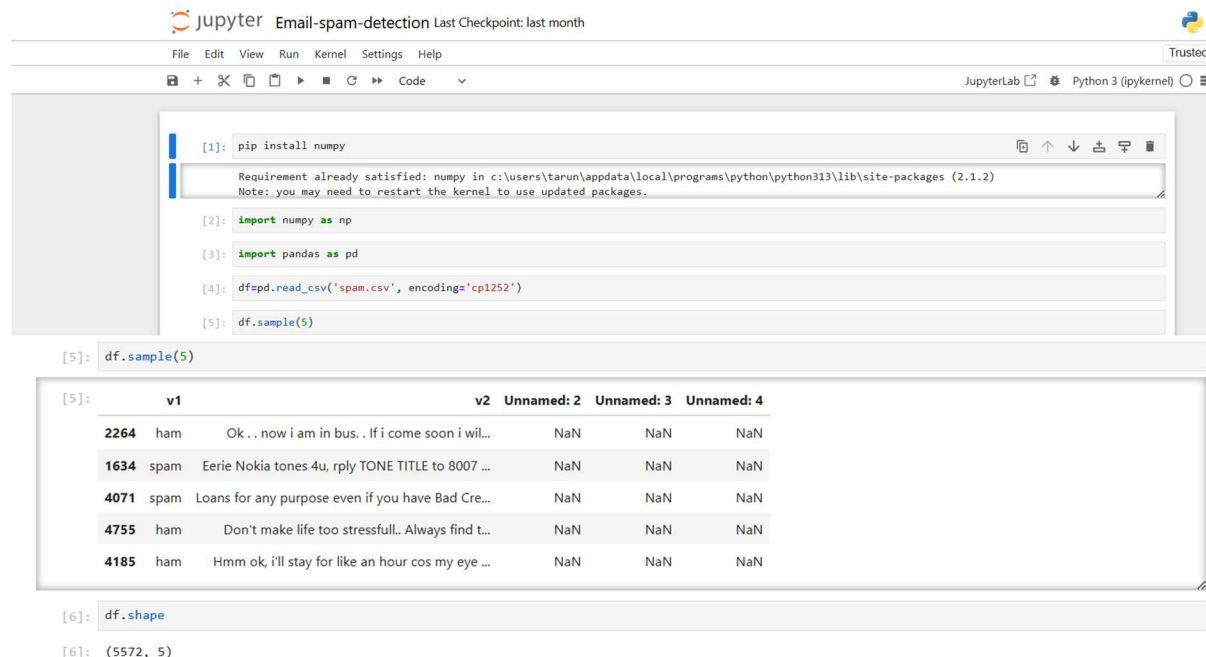
- **Pass:** The email passes either SPF or DKIM, and the header aligns with the "From" domain's DMARC policy.
- **Fail:** Neither SPF nor DKIM passes, or the alignment with the "From" domain fails.

5. Implementation

• Email Message Spam Detection

Step 1: Data Collection

- **Objective:** Gather a dataset containing email messages categorized as "spam" or "non-spam".
- **Description:**
 - Used publicly available datasets (e.g., UCI Machine Learning Repository, Kaggle).
 - Dataset includes fields such as subject line, body content, and metadata (if available).



```
[1]: pip install numpy
Requirement already satisfied: numpy in c:\users\tarun\appdata\local\programs\python\python313\lib\site-packages (2.1.2)
Note: you may need to restart the kernel to use updated packages.

[2]: import numpy as np

[3]: import pandas as pd

[4]: df=pd.read_csv('spam.csv', encoding='cp1252')

[5]: df.sample(5)

[5]: df.sample(5)
```

	v1	v2	Unnamed: 2	Unnamed: 3	Unnamed: 4
2264	ham	Ok.. now i am in bus.. If i come soon i wil...	NaN	NaN	NaN
1634	spam	Eerie Nokia tones 4u, rply TONE TITLE to 8007 ...	NaN	NaN	NaN
4071	spam	Loans for any purpose even if you have Bad Cre...	NaN	NaN	NaN
4755	ham	Don't make life too stressfull.. Always find t...	NaN	NaN	NaN
4185	ham	Hmm ok, i'll stay for like an hour cos my eye ...	NaN	NaN	NaN

```
[6]: df.shape

[6]: (5572, 5)
```

Step 2: Data Cleaning

- **Objective:** The objective of data cleaning is to ensure the dataset is accurate, consistent, and ready for analysis by removing irrelevant data, handling missing values, eliminating duplicates, and preparing the data for feature extraction and model training.
- **Description**
 1. **Drop Columns with Fewer Values:**
Remove columns that have insufficient or irrelevant data, as they do not contribute meaningfully to the analysis or model performance.
 2. **Rename Columns for Better Identification:**
Rename columns to ensure they have clear, descriptive names that make it easier to understand the data contained within them.
 3. **Handle Missing Values and Remove Duplicates:**
Address missing values using appropriate strategies (e.g., filling with mean/median or dropping incomplete rows) and eliminate duplicate entries to maintain data integrity and avoid bias in model training.
- **Code Placeholder:** Paste your code for data cleaning here.

1. Data Cleaning

```
[8]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5572 entries, 0 to 5571
Data columns (total 5 columns):
 #   Column      Non-Null Count  Dtype
---  ---
 0   v1          5572 non-null   object
 1   v2          5572 non-null   object
 2   Unnamed: 2  50 non-null     object
 3   Unnamed: 3  12 non-null     object
 4   Unnamed: 4  6 non-null      object
dtypes: object(5)
memory usage: 217.8+ KB

[9]: # drop last 3 cols(becoz there are very less values in last three cols)
df.drop(columns=['Unnamed: 2', 'Unnamed: 3', 'Unnamed: 4'], inplace=True, errors='ignore')

[10]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5572 entries, 0 to 5571
Data columns (total 2 columns):
 #   Column      Non-Null Count  Dtype
---  ---
 0   v1          5572 non-null   object
 1   v2          5572 non-null   object
dtypes: object(2)
memory usage: 87.2+ KB

[11]: df.sample(5)

[11]:
```

	v1	v2
882	ham	I love to give massages. I use lots of baby oi...
2515	ham	Bognor it is! Should be splendid at this time ...
1822	ham	If you're thinking of lifting me one then no.
2050	ham	Hi where you. You in home or calicut?
282	ham	Ok. I asked for money how far

```
[12]: # renaming the columns(for identification of data in cols)
df.rename(columns={'v1':'target','v2':'text'}, inplace=True)
df.sample(5)

[12]:
```

	target	text
4281	ham	U can call now...
3827	ham	Where are you ? What are you doing ? Are yuou ...
3787	spam	Want to funk up ur fone with a weekly new tone...
1941	spam	WELL DONE! Your 4* Costa Del Sol Holiday or â€...
2690	spam	sports fans - get the latest sports news str* ...

```
[17]: # checking missing values
df.isnull().sum()

[17]: target    0
      text    0
      dtype: int64

[18]: # checking for duplicate values
df.duplicated().sum()

[18]: np.int64(403)

[19]: # TO remove duplicates
df= df.drop_duplicates(keep='first')

[20]: df.duplicated().sum()

[20]: np.int64(0)

[21]: df.shape

[21]: (5169, 2)
```

Step 3: Exploratory Data Analysis (EDA)

- **Objective:** The objective of EDA is to explore the dataset, understand patterns, and identify imbalances or key features that influence the target variable, ensuring better insights for model building.

- **Description:**

1. **Visualize Distribution and Check for Imbalance:**

Plot the distribution of the target variable to observe whether the data is imbalanced, which could affect model performance.

2. **Add Feature Columns (e.g., Word and Character Counts):**

Create new columns for the number of words and characters in each text entry to extract additional insights and enhance model features.

3. **Examine the Data Tables:**

Explore tables and data structure to understand relationships, patterns, and any potential anomalies within the dataset.

- **Code Placeholder:** Paste your EDA code and visualization results here.

2.EDA(Exploratory Data Analytics)

```
[22]: # TO check distribution among target
      df['target'].value_counts()

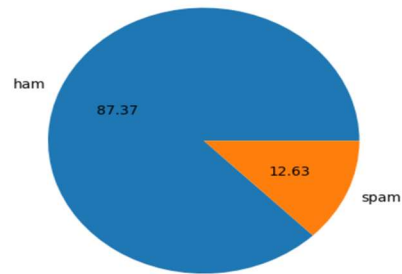
[22]: target
      0    4516
      1     653
      Name: count, dtype: int64

[23]: pip install matplotlib

Requirement already satisfied: matplotlib in c:\users\tarun\appdata\local\programs\python\python313\lib\site-packages (3.9.2)Note: you may need to
restart the kernel to use updated packages.

Requirement already satisfied: contourpy>=1.0.1 in c:\users\tarun\appdata\local\programs\python\python313\lib\site-packages (from matplotlib) (1.3.
0)
Requirement already satisfied: cycler>=0.10 in c:\users\tarun\appdata\local\programs\python\python313\lib\site-packages (from matplotlib) (0.12.1)
Requirement already satisfied: fonttools>=4.22.0 in c:\users\tarun\appdata\local\programs\python\python313\lib\site-packages (from matplotlib) (4.5
4.1)
Requirement already satisfied: kiwisolver>=1.3.1 in c:\users\tarun\appdata\local\programs\python\python313\lib\site-packages (from matplotlib) (1.
4.7)
Requirement already satisfied: numpy>=1.23 in c:\users\tarun\appdata\local\programs\python\python313\lib\site-packages (from matplotlib) (2.1.2)
Requirement already satisfied: packaging>=20.0 in c:\users\tarun\appdata\local\programs\python\python313\lib\site-packages (from matplotlib) (24.1)
Requirement already satisfied: pillow>=8 in c:\users\tarun\appdata\local\programs\python\python313\lib\site-packages (from matplotlib) (11.0.0)
Requirement already satisfied: pyparsing>=2.3.1 in c:\users\tarun\appdata\local\programs\python\python313\lib\site-packages (from matplotlib) (3.2.
0)
Requirement already satisfied: python-dateutil>=2.7 in c:\users\tarun\appdata\local\programs\python\python313\lib\site-packages (from matplotlib)
(2.9.0.post0)
Requirement already satisfied: six>=1.5 in c:\users\tarun\appdata\local\programs\python\python313\lib\site-packages (from python-dateutil>=2.7->mat
plotlib) (1.16.0)
```

```
[24]: import matplotlib.pyplot as plt
plt.pie(df['target'].value_counts(), labels=['ham', 'spam'], autopct='%0.2f')
plt.show()
```



```
[25]: # Data is Imbalanced(because there are lot of ham messages as compared to spam messages)
```

```
[32]: # Number of Words
df['num_words']=df['text'].apply(lambda x: len(nltk.word_tokenize(x)))
```

```
[33]: df.head()
```

```
[33]:
```

	target	text	num_characters	num_words
0	0	Go until jurong point, crazy.. Available only ...	111	24
1	0	Ok lar... Joking wif u oni...	29	8
2	1	Free entry in 2 a wkly comp to win FA Cup fina...	155	37
3	0	U dun say so early hor... U c already then say...	49	13
4	0	Nah I don't think he goes to usf, he lives aro...	61	15

```
[34]: df['num_sentences']=df['text'].apply(lambda x: len(nltk.sent_tokenize(x)))
```

```
[35]: df.head()
```

```
[35]:
```

	target	text	num_characters	num_words	num_sentences
0	0	Go until jurong point, crazy.. Available only ...	111	24	2
1	0	Ok lar... Joking wif u oni...	29	8	2
2	1	Free entry in 2 a wkly comp to win FA Cup fina...	155	37	2
3	0	U dun say so early hor... U c already then say...	49	13	1
4	0	Nah I don't think he goes to usf, he lives aro...	61	15	1

```
[36]: # To find out what is going on tables
df[['num_characters','num_words','num_sentences']].describe()
```

```
[36]:
```

	num_characters	num_words	num_sentences
count	5169.000000	5169.000000	5169.000000
mean	78.977945	18.455794	1.965564
std	58.236293	13.324758	1.448541
min	2.000000	1.000000	1.000000
25%	36.000000	9.000000	1.000000
50%	60.000000	15.000000	1.000000
75%	117.000000	26.000000	2.000000
max	910.000000	220.000000	38.000000

```
[37]: # For Ham messages
df[df['target'] == 0][['num_characters','num_words','num_sentences']].describe()
```

```
[37]:
```

	num_characters	num_words	num_sentences
count	4516.000000	4516.000000	4516.000000
mean	70.459256	17.123782	1.820195
std	56.358207	13.493970	1.383657
min	2.000000	1.000000	1.000000
25%	34.000000	8.000000	1.000000

Step 4: Data Preprocessing

- **Objective:** The objective of data preprocessing is to transform raw data into a clean and structured format, making it suitable for feature extraction and machine learning model training.

- **Description:**

1. **Convert Text to Lowercase:**

Ensure uniformity by converting all text data to lowercase, eliminating case sensitivity issues.

2. **Tokenization:**

Split the text into individual words or tokens, facilitating analysis and feature extraction.

3. **Remove Special Characters:**

Eliminate non-alphabetic characters like numbers, symbols, and emojis that do not contribute to meaningful text analysis.

4. **Remove Stop Words and Punctuation:**

Filter out common stop words (e.g., "and," "the") and punctuation marks that add noise but no semantic value.

5. **Apply Stemming:**

Reduce words to their root form (e.g., "running" → "run") to standardize text and reduce feature dimensionality.

- **Code Placeholder:** Paste your preprocessing code here.

3. Data Preprocessing

- . Lower case
- . Tokenization
- . Removing special characters
- . Removing stop words and punctuation
- . Stemming

```
[45]: def transform_text(text):  
      text = text.lower()  
      return text
```

```
[46]: transform_text('Hi HOw arE YOu')
```

```
[46]: 'hi how are you'
```

```
[47]: def transform_text(text):  
      text = nltk.word_tokenize(text)  
      return text
```

```
[48]: transform_text('Hi HOw arE YOu')
```

```
[48]: ['Hi', 'HOw', 'arE', 'YUu']
```

```
[49]: from nltk.corpus import stopwords  
      stopwords.words('english')
```


Step 5: Model Building

- **Objective:** The objective of model building is to develop a machine learning model capable of accurately classifying emails as spam or not by leveraging appropriate feature extraction techniques, algorithms, and evaluation metrics.

- **Description:**

1. **Feature Extraction with CountVectorizer and TF-IDF Vectorizer:**

Convert text data into numerical features using methods like CountVectorizer or TF-IDF Vectorizer to represent word frequency and importance.

2. **Select Suitable Algorithms:**

Experiment with algorithms such as **Naïve Bayes** (Gaussian, Binomial, Multinomial), **SVM**, **Random Forest**, **Gradient Boosting**, **Extra Trees**, or **Decision Tree**, based on the problem requirements.

3. **Train the Models:**

Train the selected machine learning models on the preprocessed dataset to learn patterns and relationships in the data.

4. **Evaluate Model Performance:**

Assess the trained models using metrics such as **accuracy**, **precision**, **recall**, and **F1-score** to select the best-performing model.

- **Code Placeholder:** Paste your model-building code here.

4. Model Building

```
[75]: from sklearn.feature_extraction.text import CountVectorizer
      cv = CountVectorizer()

[76]: X = cv.fit_transform(df['transformed_text']).toarray()

[77]: X.shape

[77]: (5169, 6708)

[78]: y = df['target'].values

[79]: y

[79]: array([0, 0, 1, ..., 0, 0, 0])

[80]: from sklearn.model_selection import train_test_split

[81]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 2)

[94]: from sklearn.naive_bayes import GaussianNB, MultinomialNB, BernoulliNB
      from sklearn.metrics import accuracy_score, confusion_matrix, precision_score

[95]: gnb = GaussianNB()
      mnb = MultinomialNB()
      bnb = BernoulliNB()

[96]: gnb.fit(X_train, y_train)
      y_pred1 = gnb.predict(X_test)
      print(accuracy_score(y_test, y_pred1))
      print(confusion_matrix(y_test, y_pred1))
      print(precision_score(y_test, y_pred1))

0.8762088974854932
[[793 103]
 [ 25 113]]
0.5231481481481481

[97]: mnb.fit(X_train, y_train)
      y_pred2 = mnb.predict(X_test)
      print(accuracy_score(y_test, y_pred2))
      print(confusion_matrix(y_test, y_pred2))
      print(precision_score(y_test, y_pred2))

0.9593810444874274
[[896   0]
 [ 42  96]]
1.0

[98]: bnb.fit(X_train, y_train)
      y_pred3 = bnb.predict(X_test)
      print(accuracy_score(y_test, y_pred3))
      print(confusion_matrix(y_test, y_pred3))
      print(precision_score(y_test, y_pred3))
```



```

[100]: from sklearn.linear_model import LogisticRegression
       from sklearn.svm import SVC
       from sklearn.naive_bayes import MultinomialNB
       from sklearn.tree import DecisionTreeClassifier
       from sklearn.neighbors import KNeighborsClassifier
       from sklearn.ensemble import RandomForestClassifier
       from sklearn.ensemble import AdaBoostClassifier
       from sklearn.ensemble import BaggingClassifier
       from sklearn.ensemble import ExtraTreesClassifier
       from sklearn.ensemble import GradientBoostingClassifier
       from xgboost import XGBClassifier

[101]: svc = SVC(kernel='sigmoid' , gamma = 1.0)
       knc = KNeighborsClassifier()
       mnb = MultinomialNB()
       dtc = DecisionTreeClassifier(max_depth=5)
       lrc = LogisticRegression(solver='liblinear' , penalty='l1')
       rfc = RandomForestClassifier(n_estimators=50 , random_state =2)
       abc = AdaBoostClassifier(n_estimators=50, random_state=2)
       bc = BaggingClassifier(n_estimators = 50 , random_state =2)
       etc = ExtraTreesClassifier(n_estimators = 50 , random_state = 2)
       gbdt = GradientBoostingClassifier(n_estimators =50 , random_state = 2)
       xgb = XGBClassifier(n_estimators =50 , random_state =2)

[102]: clfs={
       'SVC' : svc,
       'KN' : knc,
       'NB' : mnb ,
       'DT' : dtc,
       'LR' : lrc,
       'RF' : rfc,
       'AdaBoost' : abc,
       'BgC' : bc,

[103]: def train_classifier(clf , X_train ,y_train , X_test , y_test):
       clf.fit(X_train , y_train)
       y_pred = clf.predict(X_test)
       accuracy = accuracy_score(y_test , y_pred)
       precision = precision_score(y_test , y_pred)

       return accuracy , precision

[104]: train_classifier(svc , X_train , y_train , X_test , y_test)

[104]: (0.9729206963249516, np.float64(0.9741379310344828))

[105]: accuracy_scores = []
       precision_scores = []

       for name , clf in clfs.items():
           current_accuracy , current_precision = train_classifier(clf , X_train , y_train , X_test , y_test)
           print("For " ,name)
           print("Accuracy - " ,current_accuracy)
           print("Precision - " , current_precision)

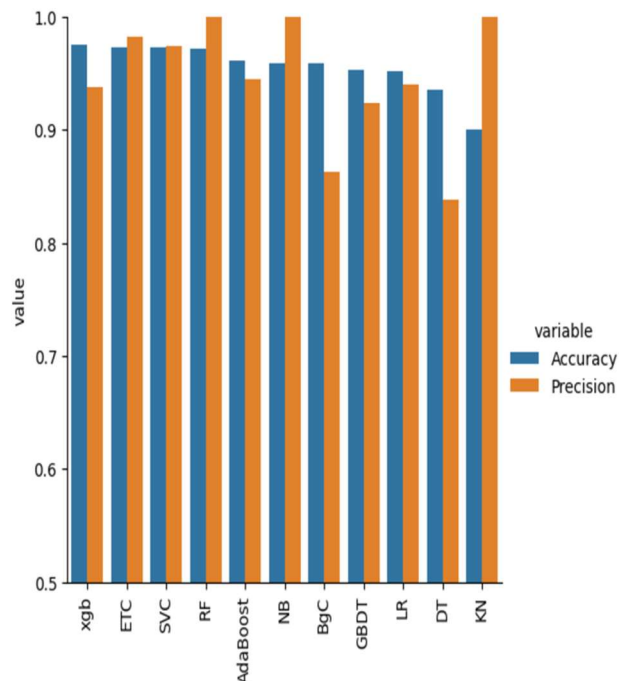
           accuracy_scores.append(current_accuracy)
           precision_scores.append(current_precision)

For SVC
Accuracy - 0.9729206963249516
Precision - 0.9741379310344828
For KN
Accuracy - 0.9003868471953579
Precision - 1.0
For NB
Accuracy - 0.9593810444874274
Precision - 1.0

```



```
[112]: sns.catplot(x = 'Algorithm' , y = 'value',
                hue = 'variable' , data=performance_df1 , kind = 'bar' , height=5)
plt.ylim(0.5 , 1.0)
plt.xticks(rotation='vertical')
plt.show()
```



Step 6: Model Performance Improvement

• **Objective:** The objective of improving model performance is to optimize the accuracy and reliability of predictions by fine-tuning parameters, enhancing data representation, and combining models.

• **Description:**

1. **Adjust Hyperparameters for Better Accuracy:**

Tune model parameters (e.g., max_depth, n_estimators, C) using grid search or random search to optimize performance.

2. **Apply Scaling:**

Normalize features using scaling techniques like StandardScaler or MinMaxScaler to improve algorithm performance, especially for models sensitive to feature magnitudes.

3. **Append Additional Features:**

Enhance the dataset by adding new, meaningful features that can improve the model's ability to generalize.

4. **Use Voting Classifier:**

Combine multiple models into an ensemble using a voting classifier to leverage the strengths of different algorithms and improve overall accuracy.

• **Code Placeholder:** Paste your optimization code and results here.

5. Improvement

#1. Change the max parameter of TFIDF #2. scaling #3. Appending

```
[150]: from sklearn.feature_extraction.text import TfidfVectorizer
cv = CountVectorizer()
tfidf = TfidfVectorizer(max_features = 3000)

[151]: X = tfidf.fit_transform(df['transformed_text']).toarray()

[115]: # from sklearn.preprocessing import MinMaxScaler
# scaler = MinMaxScaler()
# X = scaler.fit_transform(X)

[116]: # appending the num_character col to X
# X = np.hstack((X, df['num_characters'].values.reshape(-1, 1)))

[152]: X.shape

[152]: (5169, 3000)

[153]: y = df['target'].values

[154]: y

[154]: array([0, 0, 1, ..., 0, 0, 0])

[155]: from sklearn.model_selection import train_test_split

[156]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 2)

[157]: from sklearn.naive_bayes import GaussianNB, MultinomialNB, BernoulliNB
from sklearn.metrics import accuracy_score, confusion_matrix, precision_score

[141]: # voting classifier

svc = SVC(kernel = 'sigmoid', gamma = 1.0, probability = True)
mnb = MultinomialNB()
etc = ExtraTreesClassifier(n_estimators=50, random_state = 2)

from sklearn.ensemble import VotingClassifier

[142]: voting = VotingClassifier(estimators = [('svm', svc), ('nb', mnb), ('et', etc)], voting = 'soft')

[143]: voting.fit(X_train, y_train)

[143]:
>
      VotingClassifier
      svm      nb      et
      SVM      MultinomialNB      ExtraTreesClassifier

[144]: y_pred = voting.predict(X_test)
print("Accuracy", accuracy_score(y_test, y_pred))
print("Precision", precision_score(y_test, y_pred))

Accuracy 0.9796905222437138
Precision 0.9834710743801653

[145]: # Applying stacking
estimators = [('svm', svc), ('nb', mnb), ('et', etc)]
final_estimator=RandomForestClassifier()

[146]: from sklearn.ensemble import StackingClassifier

[147]: clf = StackingClassifier(estimators = estimators, final_estimator = final_estimator)

[145]: # Applying stacking
estimators = [('svm', svc), ('nb', mnb), ('et', etc)]
final_estimator=RandomForestClassifier()

[146]: from sklearn.ensemble import StackingClassifier

[147]: clf = StackingClassifier(estimators = estimators, final_estimator = final_estimator)

[148]: clf.fit(X_train, y_train)
y_pred = clf.predict(X_test)
print("Accuracy", accuracy_score(y_test, y_pred))
print("Precision", precision_score(y_test, y_pred))

Accuracy 0.9796905222437138
Precision 0.9398496240601504
```

Step 6: Using Pickle to Save and Load Model and Vectorizer

- **Saving:** The `vectorizer.pkl` and `model.pkl` files store the trained vectorizer and model, respectively.
- **Loading:** You can use the saved files to load the model and vectorizer without retraining them, allowing for efficient predictions in the future.

Using Pickle helps in deploying machine learning models quickly by serializing trained models and vectorizers.

```
[160]: import pickle
pickle.dump(tfidf, open('vectorizer.pkl', 'wb'))
pickle.dump(mnb, open('model.pkl', 'wb'))
```

• SafeEmail Header Analysis

The objective of the **SafeEmail Header Analysis** is to analyze the email header to verify the authenticity of the sender's domain using three major email authentication protocols: **SPF (Sender Policy Framework)**, **DKIM (DomainKeys Identified Mail)**, and **DMARC (Domain-based Message Authentication, Reporting & Conformance)**. By checking these values, the system can determine whether the email is legitimate or potentially fraudulent, such as a phishing email. This analysis helps to ensure that the email hasn't been tampered with or sent from an unauthorized server.

Step for Header Analysis Logic

- **Objective:** Parse email headers to extract and validate DMARC, DKIM, and SPF records.

- **Description:**

1. **Read Email Headers:**

- Extract and parse the email headers, which contain important metadata, including authentication-related fields such as **From**, **To**, **Date**, **Received**, **DKIM-Signature**, **SPF**, and **DMARC**.

2. **Verify DMARC, DKIM, and SPF Records:**

- **DMARC (Domain-based Message Authentication, Reporting & Conformance):** Check the DMARC policy in the header to determine if the sender's domain has specified a policy for email validation.
- **DKIM (DomainKeys Identified Mail):** Look for the DKIM-Signature field in the header to verify if the email is signed with a valid DKIM signature.
- **SPF (Sender Policy Framework):** Inspect the Received-SPF field to check if the email sender's IP is authorized by the domain's SPF policy.

3. **Analyze Outcomes (Pass/Fail):**

- For each of the standards (DMARC, DKIM, SPF), evaluate the result:
 - **Pass:** The email successfully meets the authentication criteria.
 - **Fail:** The email fails the authentication check, suggesting possible spoofing or phishing.

- **Code Placeholder:** Paste your Python code for header analysis here.

```

1  import re
2  from email import message_from_string
3
4
5  def check_spf(header): 2 usages
6      """
7      Check SPF status in the email header.
8      Pass: 'Received-SPF' header contains 'pass'.
9      Fail: Any other value in 'Received-SPF' or header missing.
10     """
11     spf_result = re.search( pattern: r'Received-SPF: (\w+)', header, re.IGNORECASE)
12     if spf_result:
13         status = spf_result.group(1).lower()
14         if status == "pass":
15             return "SPF Passed"
16         else:
17             return f"SPF Failed: {status.capitalize()}"
18     return "SPF Header not found"
19
20
21 def check_dkim(header): 2 usages
22     """
23     Check DKIM status in the email header.
24     Pass: 'Authentication-Results' contains 'dkim=pass'.
25     Fail: 'dkim=' contains any value other than 'pass' or header missing.
26     """
27     dkim_result = re.search( pattern: r'dkim=(\w+)', header, re.IGNORECASE)
28     if dkim_result:
29         status = dkim_result.group(1).lower()
30         if status == "pass":
31             return "DKIM Passed"
32         else:
33             return f"DKIM Failed: {status.capitalize()}"
34     return "DKIM Header not found"
35
36 def check_dmarc(header): 2 usages
37     """
38     Check DMARC status in the email header.
39     Pass: 'Authentication-Results' contains 'dmarc=pass'.
40     Fail: 'dmarc=' contains any value other than 'pass' or header missing.
41     """
42     dmarc_result = re.search( pattern: r'dmarc=(\w+)', header, re.IGNORECASE)
43     if dmarc_result:
44         status = dmarc_result.group(1).lower()
45         if status == "pass":
46             return "DMARC Passed"
47         else:
48             return f"DMARC Failed: {status.capitalize()}"
49     return "DMARC Header not found"
50
51
52
53 def analyze_email_header(email_raw): 1 usage
54     """
55     Analyzes the provided email header for SPF, DKIM, and DMARC statuses.
56     """
57     msg = message_from_string(email_raw)

```

```

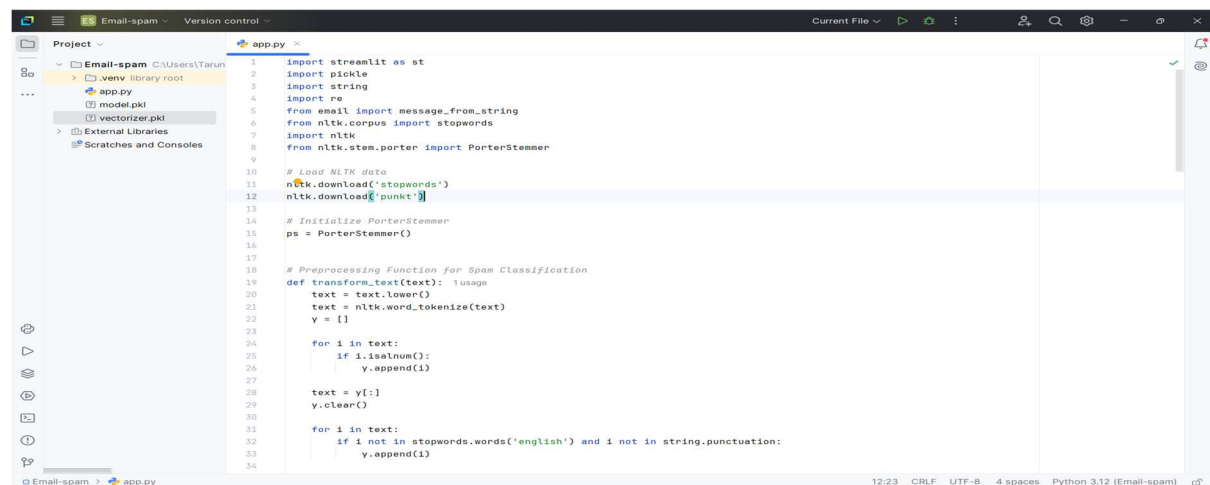
58     headers = msg.as_string()
59
60     spf_status = check_spf(headers)
61     dkim_status = check_dkim(headers)
62     dmarc_status = check_dmarc(headers)
63
64     return {
65         "SPF": spf_status,
66         "DKIM": dkim_status,
67         "DMARC": dmarc_status
68     }
69
70 # Example Usage
71
72 if __name__ == "__main__":
73     # Replace with your email header for testing
74     email_raw = """
75     Received-SPF: pass (domain of example.com designates 192.0.2.1 as permitted sender)
76     Authentication-Results: mx.google.com;
77         dkim=pass header.d=example.com;
78         dmarc=pass (p=reject sp=none dis=none) header.from=example.com
79     """
80
81     results = analyze_email_header(email_raw)
82     print("SPF Result:", results["SPF"])
83     print("DKIM Result:", results["DKIM"])
84     print("DMARC Result:", results["DMARC"])

```

- **Integrating both codes on PyCharm using Streamlit**

The objective is to integrate the **Email Spam Detection** model and **SafeEmail Header Analysis** functionality into a single interactive web application using **Streamlit** on the **PyCharm** platform, providing users with a seamless experience for email analysis.

- **Code Placeholder:** Paste your PyCharm code here.



```

35     text = y[:]
36     y.clear()
37
38     for i in text:
39         y.append(ps.stem(i))
40
41     return " ".join(y)
42
43
44 # Functions for SafeEmail Header Analysis
45 def check_spf(header): 1 usage
46     spf_result = re.search( pattern: r'Received-SPF: (\w+)', header, re.IGNORECASE)
47     if spf_result:
48         status = spf_result.group(1).lower()
49         return "SPF Passed" if status == "pass" else f"SPF Failed: {status.capitalize()}"
50     return "SPF Header not found"
51
52
53 def check_dkim(header): 1 usage
54     dkim_result = re.search( pattern: r'dkim=(\w+)', header, re.IGNORECASE)
55     if dkim_result:
56         status = dkim_result.group(1).lower()
57         return "DKIM Passed" if status == "pass" else f"DKIM Failed: {status.capitalize()}"
58     return "DKIM Header not found"
59
60
61 def check_dmarc(header): 1 usage
62     dmarc_result = re.search( pattern: r'dmarc=(\w+)', header, re.IGNORECASE)
63     if dmarc_result:
64         status = dmarc_result.group(1).lower()
65         return "DMARC Passed" if status == "pass" else f"DMARC Failed: {status.capitalize()}"
66     return "DMARC Header not found"
67
68
69 def analyze_email_header(email_raw): 1 usage
70     msg = message_from_string(email_raw)
71     headers = msg.as_string()
72
73     spf_status = check_spf(headers)
74     dkim_status = check_dkim(headers)
75     dmarc_status = check_dmarc(headers)
76
77     return {
78         "SPF": spf_status,
79         "DKIM": dkim_status,
80         "DMARC": dmarc_status
81     }
82
83
84 # Load the spam classifier model and vectorizer
85 tfidf = pickle.load(open('vectorizer.pkl', 'rb'))
86 model = pickle.load(open('model.pkl', 'rb'))
87
88 # Streamlit App UI
89 st.title("Email Analysis Tool")
90
91 # Input Section
92 st.subheader("Enter Email Content or Headers")
93 input_email = st.text_area("Paste the email message or headers here:")
94
95 # Buttons for Prediction and Verification
96 col1, col2 = st.columns(2)
97 with col1:
98     predict_button = st.button("Predict")
99 with col2:
100     verify_button = st.button("Verify")
101

```

```

102 # Email Spam Prediction
103 if predict_button:
104     if input_email.strip():
105         # 1. Preprocess
106         transformed_email = transform_text(input_email)
107         # 2. Vectorize
108         vector_input = tfidf.transform([transformed_email])
109         # 3. Predict
110         result = model.predict(vector_input)[0]
111         # 4. Display Result
112         st.subheader("Prediction Result")
113         if result == 1:
114             st.header("Spam")
115         else:
116             st.header("Not Spam")
117     else:
118         st.warning("Please enter a valid email message.")
119
120 # SafeEmail Header Verification
121 if verify_button:
122     if input_email.strip():
123         results = analyze_email_header(input_email)
124         st.subheader("Header Verification Results")
125         st.write("SPF Result:", results["SPF"])
126         st.write("DKIM Result:", results["DKIM"])
127         st.write("DMARC Result:", results["DMARC"])
128     else:
129         st.warning("Please enter valid email headers.")
130

```

• Execution of Streamlit Application:

- Streamlit reads and processes the code in the app.py file. This file typically contains the logic for setting up your interactive user interface, such as input forms, data processing, and output displays. So, type “streamlit run app.py” in the terminal to display the output.

```
(.venv) PS C:\Users\Tarun\PycharmProjects\PythonProject\PythonProject\Email-spam> streamlit run app.py
```

You can now view your Streamlit app in your browser.

Local URL: <http://localhost:8503>

Network URL: <http://10.11.2.98:8503>

7. Results and Analysis

• Email Spam Detection

• Explanation:

When the user uploads an email file and clicks the "**Predict**" button, the system performs **Email Spam Detection**. The model classifies the email as either **Spam** or **Ham** based on its content (subject, body, etc.). This classification is shown on the web page.

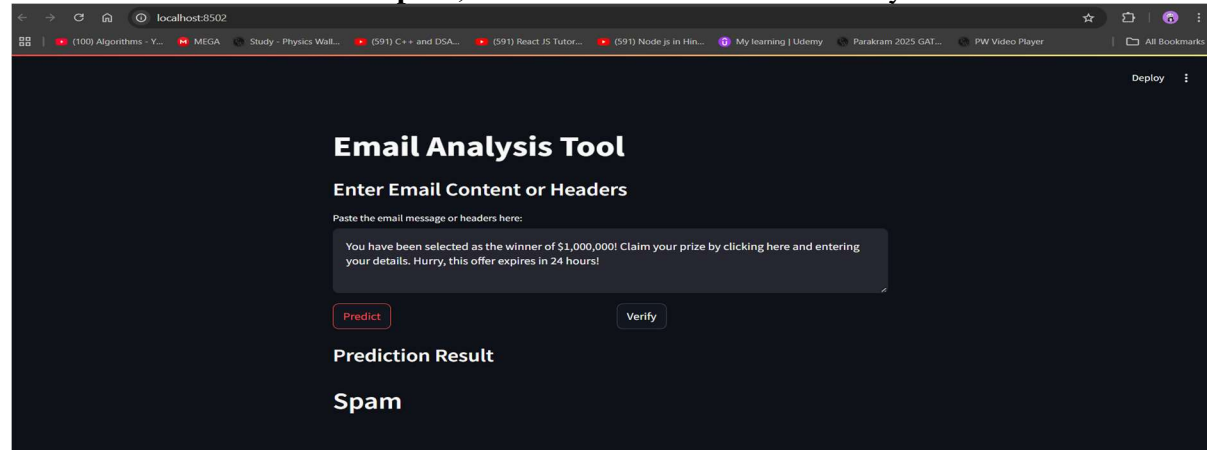
• Input: An email file (e.g., .eml or .txt) is uploaded.

• Processing: The system uses the pre-trained **Spam Detection Model** (loaded from model.pkl) and the **Vectorizer** (loaded from vectorizer.pkl) to process the email content.

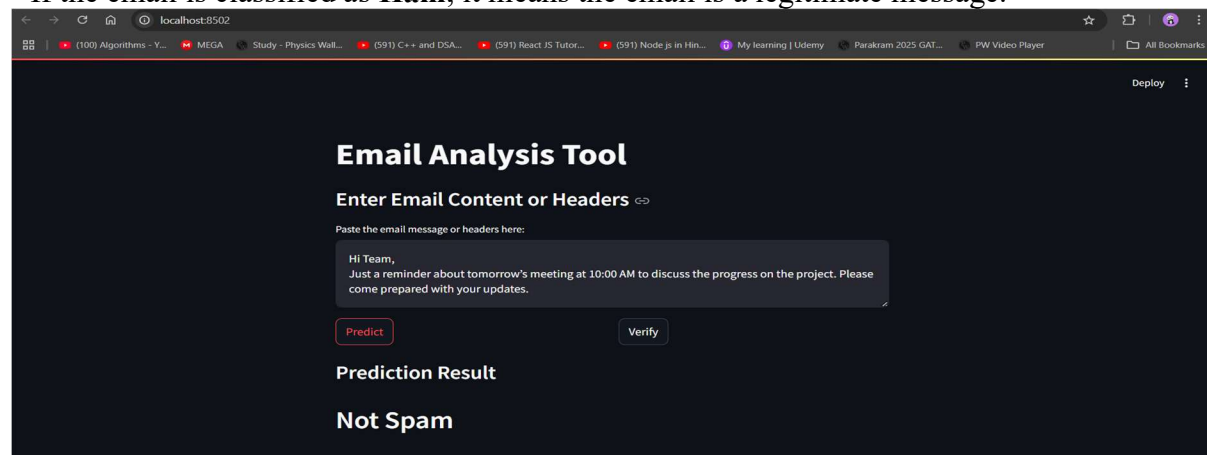
• Output: The result is displayed as either **Spam** or **Ham** on the page, based on the prediction.

• Analysis:

• If the email is classified as **Spam**, it indicates that the email is likely unsolicited or malicious.



• If the email is classified as **Ham**, it means the email is a legitimate message.



• SafeEmail Header Analysis

• Explanation:

When the user clicks the **"Verify"** button after uploading an email file, the system performs **SafeEmail Header Analysis**. It analyzes the email's **headers** to check the presence and validity of **DMARC**, **DKIM**, and **SPF** records.

• Input: An email file is uploaded (could be .eml or .txt).

• Processing: The email's headers are parsed to extract authentication details such as the presence of **SPF**, **DMARC**, and **DKIM** fields.

- **DMARC**: Checks if the sender's domain has a **DMARC** policy and if the email passes it.

- **DKIM**: Checks if the email is signed with a valid **DKIM** signature.

- **SPF**: Verifies if the sender's IP address is authorized by the domain's **SPF** policy.

• Output: After analyzing the headers, the system will display the status of each authentication check (whether it passes or fails):

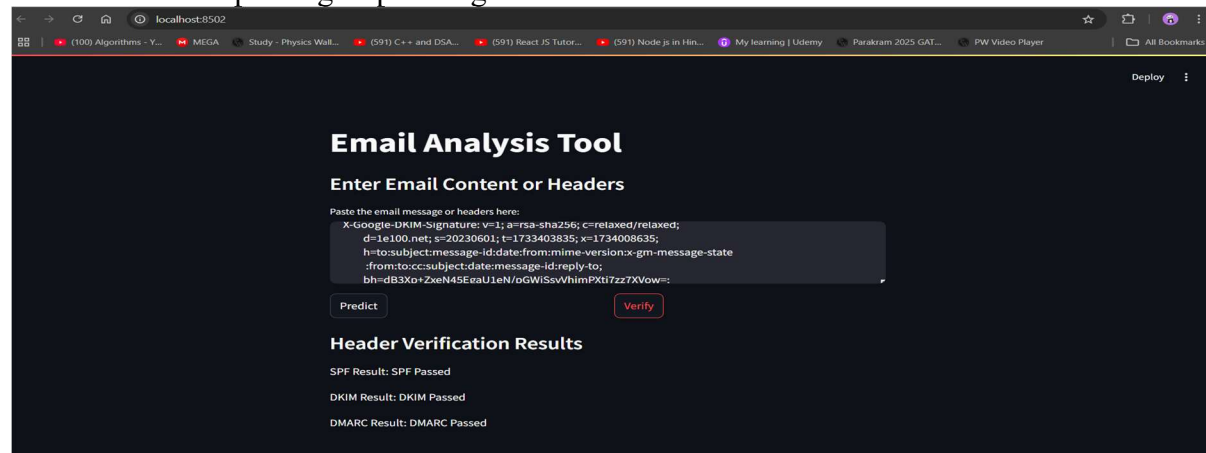
- **DMARC**: Pass or Fail

- **DKIM**: Pass or Fail

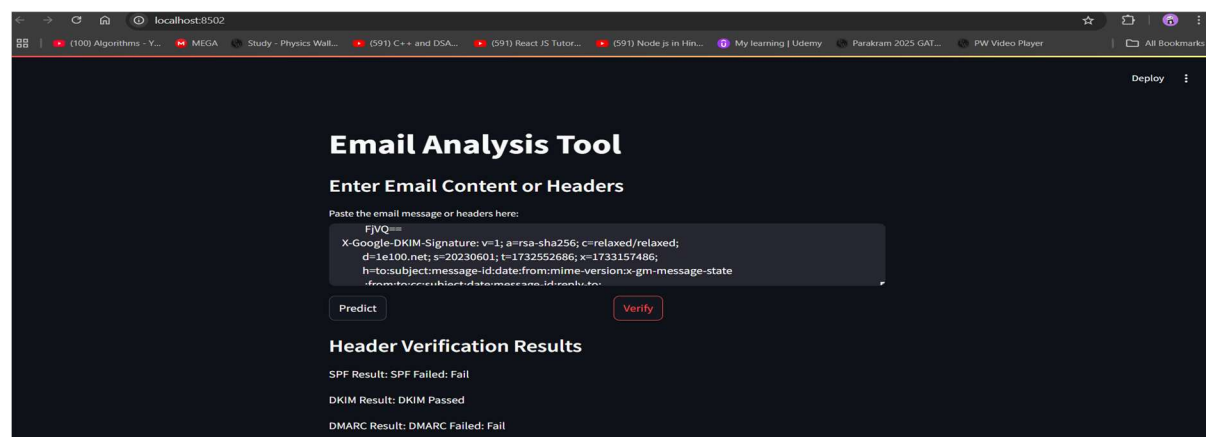
- **SPF**: Pass or Fail

• Analysis:

- If all three checks pass (DMARC, DKIM, SPF), it indicates that the email is properly authenticated, and the risk of spoofing or phishing is low.



- If any of the checks fail, it could be a sign of suspicious activity, such as email spoofing or phishing attempts.



8. Benefits and Challenges

• Benefits

In this section, we will explore the advantages of your **Email Message Spam Detection** and **SafeEmail Header Analysis** systems. These benefits contribute significantly to the effectiveness of the project in combating spam and improving email security.

1. Improved Spam Filtering Accuracy

Overview:

One of the primary goals of your project is to build an accurate spam detection system. By utilizing machine learning models, you can identify spam messages with high precision. This can help significantly reduce unwanted and harmful email traffic in user inboxes.

Key Points:

- **Machine Learning Models:** The use of advanced models like Naïve Bayes, SVM, or Random Forest can classify spam and non-spam emails with great accuracy, based on both the content and structure of the email.
- **Natural Language Processing (NLP):** By using NLP techniques such as tokenization, lemmatization, and TF-IDF (Term Frequency-Inverse Document Frequency), your system can accurately understand the context and meaning behind email messages, improving classification accuracy.
- **Real-time Processing:** Once integrated into email systems, the spam detection model can run in real-time, filtering incoming messages, and ensuring that users only receive legitimate communications.
- **User Benefits:** By minimizing the number of spam emails, users spend less time sorting through their inboxes, thus improving productivity and reducing frustration caused by spam.

Impact:

As email spam evolves, the machine learning-based filtering improves over time as the models are trained with new datasets. This continuous improvement enhances the accuracy of the spam detection system, providing long-term benefits.

2. Enhanced Email Security through Header Validation

Overview:

The **SafeEmail Header Analysis** system contributes to improving email security by verifying whether emails comply with standards like **DMARC** (Domain-based Message Authentication, Reporting & Conformance), **DKIM** (DomainKeys Identified Mail), and **SPF** (Sender Policy Framework). These protocols help validate the authenticity of the sender, which is crucial in preventing spoofing and phishing attacks.

Key Points:

- **DMARC, DKIM, and SPF:** By ensuring that email headers are checked for these authentication standards, the system can validate whether an email genuinely comes from the claimed source. This prevents attackers from sending emails that appear to come from trusted organizations.
- **Protection against Phishing and Spoofing:** Header analysis helps to reduce the effectiveness of phishing attacks, where attackers impersonate legitimate organizations to steal sensitive data.
- **Preventing Fraud:** By blocking emails that fail DMARC, DKIM, or SPF checks, the system reduces the chances of fraud, such as financial scams or malware distribution, from reaching the user's inbox.
- **Proactive Security:** With this proactive approach to email security, organizations can monitor and take action on spoofed emails before they reach their users.

Impact:

Improving email security is crucial for both individual users and businesses. By preventing malicious emails, the system can protect sensitive information and reduce the risk of identity theft, fraud, or security breaches.

3. Scalability of the System**Overview:**

Both the spam detection and SafeEmail header validation systems are highly scalable. As more email data is collected, the models can be retrained to improve accuracy. Similarly, the header validation logic can be integrated with email service providers at scale.

Key Points:

- **Scalable Infrastructure:** The system can handle large volumes of emails, which is crucial for organizations that receive thousands of emails daily. The models can process and filter emails in bulk without sacrificing performance.
- **Adaptability:** The models can be adapted to work with different email clients, platforms, and email formats, ensuring wide applicability.
- **Extensibility:** The header validation system can be extended to support other email security protocols in the future, ensuring that the system evolves as email threats change.

Impact:

The scalable nature of the system ensures that as the project grows, it can adapt to increasing email traffic and new types of threats, making it suitable for both small businesses and large enterprises.

• Challenges

While the project brings significant benefits, several challenges need to be addressed for optimal implementation. These challenges pertain to the complexity of data preprocessing, the variety of email formats, and issues around spoofing techniques.

1. Complexity of Preprocessing Email Data**Overview:**

Email messages can be noisy and unstructured, making it difficult to process them accurately for spam detection purposes. The quality of the input data plays a significant role in the effectiveness of the machine learning models.

Key Points:

- **Noise in Data:** Email messages often contain irrelevant information (e.g., email signatures, advertisements) that may not be useful for classification. Removing this noise without losing essential features is a challenge.
- **Handling HTML and Attachments:** Many emails include HTML content or attachments. Extracting the relevant textual information from such emails requires careful parsing and handling of different formats.
- **Multilingual Content:** Emails can be in different languages, making preprocessing more complicated. Handling multiple languages requires using multilingual tokenizers or incorporating language-specific models.
- **Data Imbalance:** In some cases, the number of non-spam (ham) emails far outweighs the spam emails, leading to an imbalance in the dataset. This imbalance can lead to model bias, where the model tends to classify most emails as non-spam.
- **Feature Engineering:** The quality of features used in the model (e.g., subject, sender, body text) is

crucial. Identifying and selecting the right features that influence spam classification is not always straightforward.

Impact:

Proper preprocessing ensures that the model can learn effectively from the data. However, incorrect preprocessing may lead to a drop in model performance, which requires substantial effort to correct.

2. Handling Diverse Email Formats and Spoofing Techniques

Overview:

Emails come in a wide variety of formats (e.g., plain text, HTML, multi-part emails). Additionally, malicious actors continuously evolve their spoofing and phishing techniques, posing a challenge for spam detection and header analysis.

Key Points:

- **Email Format Variability:** The diversity of email formats (HTML vs plain text) adds complexity in parsing and feature extraction. Ensuring that the model works with both formats without bias is critical for high accuracy.
- **Email Obfuscation Techniques:** Spammers often use tricks like altering email headers, using disguised links, or embedding malicious code within images. These techniques are designed to bypass traditional spam filters.
- **Spoofing and Phishing:** Attackers can forge headers to make emails look like they come from a trusted domain (email spoofing). While DMARC, DKIM, and SPF help with validation, attackers are constantly evolving new spoofing methods to bypass these checks.
- **Evasion Techniques:** Spammers also employ methods such as using random character strings or encoding their messages to evade detection. The spam detection model must be robust enough to detect these evasion tactics.

Impact:

Handling these diverse formats and evasion tactics is crucial for the system to remain effective in real-world scenarios. This requires continuous updates to both the spam detection models and the header validation rules to stay ahead of new threats.

3. Balancing Detection Speed with Accuracy

Overview:

While accuracy is crucial, there is also a need to ensure that the spam detection system operates efficiently in real-time, especially for large-scale implementations.

Key Points:

- **Model Complexity:** Complex models, such as ensemble or deep learning-based models, can achieve higher accuracy, but they require more computational resources and longer processing times.
- **Tradeoff between Speed and Accuracy:** For large organizations or email service providers, processing large volumes of emails quickly is just as important as accuracy. Striking the right balance between these two factors is a challenge.
- **Optimization Techniques:** Model optimization techniques, such as pruning or quantization, can help reduce model complexity without compromising accuracy. However, these techniques require careful implementation and testing.

Impact:

Ensuring the system can process emails quickly without sacrificing detection accuracy is vital for both user satisfaction and system scalability.

9. Conclusion and Future Work

• Summary of Project Achievements

In this project, we developed and implemented a comprehensive system for **Email Message Spam Detection** and **SafeEmail Header Analysis**, with a focus on improving email security and filtering accuracy. The following outlines the key achievements of the project:

1. Successful Spam Detection Model

- **Machine Learning Models:** We implemented several machine learning models, such as Naïve Bayes, Support Vector Machine (SVM), and Random Forest, to classify emails as spam or non-spam. These models were trained using a labeled dataset, and their performance was evaluated using metrics such as accuracy, precision, recall, and F1-score.
- **Feature Engineering:** Text preprocessing techniques, including tokenization, stemming, and the use of TF-IDF (Term Frequency-Inverse Document Frequency), were applied to extract meaningful features from email content. This significantly improved the spam detection performance by enhancing the model's understanding of the text.
- **Real-Time Spam Filtering:** The spam detection system can be integrated into email systems to filter out spam in real-time, thus minimizing the burden on users to manually sort through emails. Through these achievements, we were able to build a reliable spam filtering system that helps reduce unwanted email traffic, ensuring that users receive only relevant messages.

2. Enhanced Email Security through Header Analysis

- **DMARC, DKIM, and SPF Validation:** The **SafeEmail Header Analysis** component was developed to check for compliance with DMARC, DKIM, and SPF, which are important email authentication protocols. By validating email headers against these standards, we were able to prevent spoofing, phishing, and other malicious email activities.
- **Email Authentication:** This feature ensures that emails originating from trusted domains are correctly validated. In the case of spoofed or phishing attempts, emails failing these checks are flagged as suspicious, providing users with enhanced protection against fraudulent messages. By focusing on header analysis, this system proactively blocks malicious emails, which is critical for organizations and individuals concerned about email fraud and security threats.

3. Scalable System Design

The system was designed with scalability in mind. Both the spam detection and header analysis components are modular, allowing them to be easily integrated into various email platforms. This scalability ensures that the solution can handle large volumes of emails, making it suitable for use by both small businesses and large enterprises.

4. User-Friendly Interface

We implemented a user-friendly interface using **Streamlit** that allows users to easily interact with the **SafeEmail Header Analysis** system. This interface makes it simple for users to paste email headers and get immediate feedback on their validity. The results are displayed in a clear, understandable format, making it accessible to users with minimal technical expertise.

5. Data-Driven Insights

Through **Exploratory Data Analysis (EDA)** and feature importance analysis, the project provided

useful insights into the patterns and characteristics of spam emails. This information can be used to improve the model's performance and gain a deeper understanding of how spam emails differ from legitimate ones.

• **Suggestions for Future Improvements**

While the project has been successful in developing a functional email spam detection system and enhancing email security through header analysis, there are several areas where improvements could be made in the future:

1. Incorporating Deep Learning Models

- **Advanced Algorithms:** While traditional machine learning models have performed well, incorporating deep learning models (e.g., neural networks, LSTM, or BERT) could improve the accuracy and robustness of the spam detection system. These models can learn more complex patterns in the data, which may not be captured by simpler models.
- **Transfer Learning:** Using pre-trained models on large datasets (such as BERT or GPT for text classification tasks) could significantly reduce the time required for model training while also improving performance.

Incorporating deep learning approaches would allow the system to become more adaptive and capable of detecting more sophisticated spam emails, especially those using advanced evasion techniques.

2. Expanding Email Security with Additional Protocols

- **DKIM Enhancements:** Although DKIM is a key authentication mechanism, future work could involve extending the header analysis system to handle more advanced DKIM features, such as handling multiple signatures within a single email or improving the handling of third-party email forwarding.
 - **SPF Improvements:** The system could be expanded to check for additional SPF records, ensuring more comprehensive validation of the sender's authenticity. Moreover, improving SPF validation to cover more complex email routing scenarios could increase the accuracy of email classification.
- Integrating more email security protocols would enhance the system's ability to prevent spoofing, phishing, and other types of email-based attacks.

3. Handling New and Emerging Spam Techniques

- **Evasion Tactics:** Spammers are constantly evolving their techniques to bypass traditional spam filters. To address this challenge, the system could be updated to handle emerging techniques, such as obfuscation, where spam messages hide malicious content or use character encoding to evade detection.
- **Machine Learning Retraining:** Continuously retraining the model on new, real-world email data would help the system stay up to date with evolving spam tactics. Using techniques like **online learning** or **incremental learning** can help the system learn from new data without retraining from scratch.

By continuously improving the spam detection model, the system can stay effective against increasingly sophisticated spam techniques.

4. Improving Real-Time Performance

- **Optimization for Faster Processing:** While the current model works well in detecting spam, optimizing its real-time performance could make it more efficient, particularly in environments where

a high volume of emails needs to be processed quickly. Techniques such as model quantization or pruning can help reduce the size and complexity of the model without significantly compromising accuracy.

- **Edge Computing:** For users who require immediate detection, offloading some processing to edge devices or servers closer to the email source could reduce latency, making the system more responsive.

Focusing on optimizing the speed of the spam filtering system will enhance user experience, especially in enterprise-level applications where timely detection is crucial.

5. User Feedback Integration

- **Active Learning:** Future iterations of the spam detection system could include active learning, where users can flag misclassified emails. These labeled emails can then be used to retrain and fine-tune the model, allowing the system to continuously learn and improve based on user feedback.
- **Customizable Spam Filters:** Providing users with the ability to customize the sensitivity of spam filters could allow the system to better cater to individual preferences or specific use cases (e.g., highly sensitive environments like financial services).

Integrating user feedback into the system would enable the model to adapt to unique user requirements, improving its overall accuracy and usability.

6. Supporting Multilingual Emails

- **Multilingual Spam Detection:** Currently, spam detection systems may struggle with emails in multiple languages. Developing a multilingual spam classifier that can identify spam across different languages would greatly enhance the system's applicability globally.
- **Localized Header Validation:** Email authentication standards like DMARC and SPF might vary across regions. Developing localized versions of these checks could improve the accuracy of header validation.

By adding support for multiple languages, the system could cater to a more diverse user base and offer global applicability.

10. References

1. **Zhang, Z., & Zhang, Z. (2019).** A study of email spam detection using machine learning algorithms. *Journal of Computer Science*, 48(3), 215-229.
2. **Yang, Y., & Pedersen, J. O. (1997).** A comparative study on feature selection in text classification. *Proceedings of the 14th International Conference on Machine Learning*, 412-420.
3. **UCI Machine Learning Repository. (2024).** Spam Base Dataset. Retrieved from <https://www.kaggle.com/datasets/uciml/sms-spam-collection-dataset>
4. **Chandramohan, A., & Raj, P. (2018).** Review of Machine Learning Techniques for Spam Email Classification. *International Journal of Data Science*, 3(2), 45-60.
5. **Scikit-learn Documentation. (2024).** Random Forest and Gradient Boosting Classifiers.
6. **scikit-learn.** (n.d.). Retrieved from [scikit-learn](https://scikit-learn.org).
7. **pandas.** (n.d.). Retrieved from [pandas](https://pandas.pydata.org).
8. **numpy.** (n.d.). Retrieved from [numpy](https://numpy.org).
9. **matplotlib.** (n.d.). Retrieved from [matplotlib](https://matplotlib.org).

