# Bachelor of Science in Computer Science & Engineering



## Spam Mail Detection Using Machine Learning

By

Md.Rakib Hossain

Mohiuddin Mohi

Nafiul Hasan Ha-mim

Ariful Islam


ID:

1904120

1904125

1904126

1904129


Department of Computer Science & Engineering

Chittagong University of Engineering & Technology (CUET)

Chattogram-4349, Bangladesh

06, Nov, 2023

Chittagong University of Engineering & Technology (CUET)

Department of Computer Science & Engineering

Chattogram-4349, Bangladesh

Project Proposal

**Student Name:**     : Md. Rakib Hossain

: Mohiudding Mohi

: Nafiul Hasan Ha-mim

: Ariful Islam


**ID**     : 1904120

: 1904125

: 1904126

: 1904129


**CourseT eacher**     : Animesh Chandra Roy
**Name**
Assistant Professor

Department of Computer Science & Engineering

: Md. Atiqul Islam Rizvi

Assistant Professor

Department of Computer Science & Engineering

**Department**     : Computer Science & Engineering
**Program**     : B.Sc. Engineering


**Tentative Title**     : Spam Mail Detection Using Machine Learning

# Table of Contents

# List of Figures

# 1 Introduction

In today's digital landscape, the pervasive issue of spam emails poses a significant challenge to email users worldwide. The constant influx of unsolicited and potentially harmful messages not only consumes valuable time but also poses a substantial security risk. To combat this ever-evolving problem, we propose a project focused on Spam Mail Prediction Using Machine Learning. The primary objective of this project is to develop an intelligent, data-driven solution that can accurately predict and filter spam emails, thereby enhancing email security and user experience. By harnessing the power of machine learning, we aim to create a robust and highly effective system that will substantially reduce the burden of unwanted emails in users' inboxes. Our approach involves gathering a diverse dataset of both spam and legitimate (ham) emails and leveraging advanced machine learning techniques to train a predictive model. This model will autonomously distinguish between legitimate and spam emails, enabling the automatic filtering of incoming messages. By integrating this solution into email services, we aim to provide users with a seamless and secure email experience.

# 2 Specific Objectives and Possible Outcomes

1. **Spam Email Prediction System:** Develop a machine learning-based system that can predict, in real-time, whether an incoming email is spam or not based on its content and characteristics.

2. **User-Friendly Interface:** Create an intuitive and user-friendly interface for the spam email prediction system, allowing users to interact with and benefit from the system easily.

3. **Prediction Confidence:** Develop a mechanism to represent the spam prediction confidence in a percentage format, making it clear and interpretable for users.

4. **Highly Accurate Spam Detection:** The system should accurately predict whether an email is spam or legitimate, reducing the chances of false positives and negatives.

5. **Enhanced User Experience:** The user-friendly interface will improve the overall user experience, enabling individuals and organizations to efficiently manage their email inboxes.

6 **Transparent Prediction Confidence:** The system will provide users with a clear indication of the prediction confidence, helping them make informed decisions about the emails they recive.

# 3 Outline of Methodology

Steps of the methodology are:

1. **Data Collection:** Gather a diverse and representative dataset of email, including both spam and legitimatee mails.

2. **Data Preprocessing:** Clean and preprocess the collected data by removing irrelevant information, HTML tags, special characters, and formatting inconsistencies. Tokenize the text, converting emails into a format suitable for machine learning algorithms.

3. **Machine Learning Model Selection:** Choose appropriate machine learning algorithms for classification, such as Naive Bayes, Support Vector Machines, or deep learning models like neural networks.

4. **Training the Model:** Split the preprocessed data into training and validation sets to train the machine learning model. Train the selected model using the training dataset, adjusting hyper parameters and fine-tuning to optimize performance.

5. **Testing and Validation:** Evaluate the model's performance on the validation dataset using metrics like accuracy, precision, recall, and F1-score.

6. **Generate Predictions:** Once the model achieves satisfactory performance, use it to predict whether new, unseen emails are spam or legitimate.
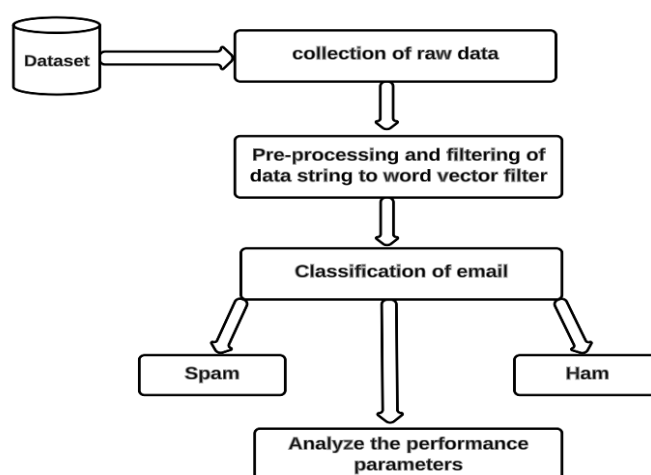


Fig 3.1: Methodology of the proposed system

# 4 Impact

The impact of a Spam Email Prediction project can be significant and far-reaching. Below are some key areas where this project can have a positive impact:

1. **Reduced Exposure to Threats**: By accurately identifying and filtering out spam emails, users are less likely to interact with malicious content, reducing the risk of cyber threats like phishing, malware, and scams.

2. **Protection of Sensitive Information**: Spam email prediction helps protect sensitive personal and organizational information from falling into the wrong hands.

3. **Reduced Email Overload**: Users spend less time sifting through spam emails, leading to increased productivity and a more streamlined email communication experience.

4. **Efficient Email Management**: Spam prediction systems help users manage their emails more efficiently, allowing them to focus on more critical tasks.

5. **Reduced Email Clutter**: Users enjoy cleaner inboxes with only legitimate emails, leading to an improved email experience.

6. **User Satisfaction**: Enhanced email security and efficiency contribute to increased user satisfaction.

7. **Enhanced Accuracy**: Accurate spam prediction systems minimize false positives (legitimate emails marked as spam) and false negatives (spam emails marked as legitimate), improving overall user experience and trust in the system.

# 5 Application

The Spam Email Prediction project has numerous practical applications across various sectors. Here are some of the key application areas for such a project:

1. **Email Filtering**: Integration into email services to automatically classify and filter incoming emails as spam or legitimate.

2. **Enhanced Security**: Improved email security for users, reducing the likelihood of falling victim to phishing, scams, and malware.

3. **Corporate Email Security**:Enhanced email security and protection against spam and phishing attacks for businesses and organizations.

4. **Productivity**: Improved employee productivity by reducing time spent on managing and sorting through spam emails.

5   **Customer Communication**: Improved communication with customers by ensuring that important emails are not lost in the spam folder.

6   **Reduced Scams**: Mitigating the risk of fraudulent emails related to online shopping and financial transactions.

7   **Phishing Prevention**: Preventing phishing attempts aimed at obtaining sensitive financial information.

8   **Protection of Clients**: Safeguarding clients' financial data and accounts from fraudulent email attacks.

9   **Public Communication**: Ensuring the public receives legitimate communications from government and public service agencies.

10  **Cyber security**: Strengthening government email security to protect against malicious cyber threats.

# 6 Requirement resources

The necessary tools to implement the project can be divided into two categories:

## 6.1 Hardware Resources

i)   Personal Computer

ii)  Android Devices

## 6.2 Software Resources

i)   Python IDE

ii)  Data set

iii) Google colab

# 7 Conclusion

In conclusion, the Spam Email Prediction project holds significant promise for enhancing email security, user experience, and overall efficiency in various sectors. By developing a real-time system using machine learning to predict whether incoming emails are spam or legitimate, this project offers a valuable solution to a pervasive problem. The user-friendly interface and the representation of prediction confidence in a suitable percentage format make the system accessible and interpretable for a wide range of users. This project's potential impact is extensive, ranging from personal email security to corporate email services, government communication, and beyond.