# E-MAIL SYSTEM (GMAIL CLONE) WITH CUSTOM SPAM CLASSIFIER

## UCS663: CONVERSATIONAL AI
## DATA SCIENCE
## MINI PROJECT

**Submitted By:**

Tushar Mittal

3CS10

101916042

**Submitted To:**

Dr. Sahil Sharma

Assistant Professor

**Computer Science and Engineering Department**

**Thapar Institute of Engineering and Technology, Patiala**

**May 2022**

# ABSTRACT

The aim of this project is to build an electronic mail system, i.e. basically a gmail clone with backend programmed in Django Framework, frontend in HTML, CSS and Javascript and database used being SQLite. The system has a mail client server that performs the function of sending and receiving mails, replying to received mails and archiving mails. Further it has been pre-trained to classify mails as spam or ham (using the Naïve-Bayes Classifier). Anyone can comfortably use the system by registering as a new user in a few steps or by logging in using the registered credentials. Viewing the mails as spam or ham is also very easy since it shows whether a received or sent mail in the respective inbox is spam or ham directly.

# TABLE OF CONTENTS

# INTRODUCTION

In today's globalized world, email is a primary source of communication. This communication can vary from personal, business, corporate to government. With the rapid increase in email usage, there has also been increase in the SPAM emails. SPAM emails, also known as junk email involves nearly identical messages sent to numerous recipients by email. Apart from being annoying, spam emails can also pose a security threat to computer system. It is estimated that spam cost businesses on the order of $100 billion in 2007.

In this project, I aim to perform automatic spam filtering in my self-created email system. I have tried to identify patterns using Naïve Bayes classification algorithm hence classifying the received emails as HAM or SPAM.

Naive Bayes classifiers are a popular statistical technique of e-mail filtering. They typically use bag-of-words features to identify spam e-mail, an approach commonly used in text classification. Naive Bayes spam filtering is a baseline technique for dealing with spam that can tailor itself to the email needs of individual users and give low false positive spam detection rates that are generally acceptable to users. It is one of the oldest ways of doing spam filtering, with roots in the 1990s.

The Naïve Bayes classification algorithm works very well in this scenario giving an accuracy of 97%.

# RELATED WORK

Email is such an integral part of our day-to-day life that it has become the primary source of communication for most. It is easy to use, safe and very fast. With billions of users worldwide and more than 300 billion emails being exchanged every day, email has come a long way since it was invented in 1971.

Today, people all over the world have a host of email clients to choose from. According to the latest data, as of January 2022, the two most used email clients in the world are Apple and Gmail. At 57.16 percent, Apple has the majority of the email client market share. It's used by more than half of the world's email users to send and receive emails. In the race to capture the majority of market share, companies like Apple and Google work continuously for the betterment of their respective email systems. Classifying Emails as spam or ham is a major task they need to deal.

For Example: Google's email client Gmail spam filter works in the following way:
1. First, it checks the email of the sender against Gmail's database of blacklisted domains.
2. If the email passes that (if the email or domain is unknown), Gmail will then check any links against its database of known malicious links and compare them to links in the incoming email.
3. After this, Gmail will also check for spelling and grammatical errors and go through its list of trigger words that are heavily featured in known spam emails.
4. Gmail also uses an in-house machine learning framework called **Tensorflow** – alongside some smart AI – to train new spam filters moving forward.

And Apple's Iphone works in the following way to  filter spam:
Traditionally, machine learning techniques formalize a problem of clustering of spam message collection through the objective function. The objective function is a maximization of similarity between messages in clusters, which is defined by k-nearest neighbor (kNN) algorithm.
Unfortunately, above approach do not provide good enough performance to filter spam e-mails for iPhone. Thus, Apple applies artificial bee-based decision tree (ABBDT) to filter spam e-mails for iPhone. In the proposed approach, decision tree is used to filter spam e-mails. In addition, artificial bee algorithm is used to ameliorate the testing accuracy of decision tree.

# DATASET

As data is the foundation of any model, selecting the correct dataset is the most critical stage in any machine learning task. I have chosen the following dataset for training the model by Naïve Bayes Classification Algorithm from the net and uploaded it is a part in my Git repository:

https://github.com/tush7301/Mail_system/blob/main/mail/spam_ham_dataset.csv



FIGURE 1: A SECTION OF THE DATASET USED FOR TRAINING THE CLASSIFIER



FIGURE 2: A SECTION OF THE DATASET USED FOR TRAINING THE CLASSIFIER

# METHODOLOGY

Naive Bayes methods are a set of supervised learning algorithms based on applying Bayes' theorem with the "naive" assumption of conditional independence between every pair of features given the value of the class variable. Bayes' theorem states the following relationship, given class variable $y$ and dependent feature vector $x_1$ through $x_n$, :

$$P(y \mid x_1, \ldots, x_n) = \frac{P(y)P(x_1, \ldots, x_n \mid y)}{P(x_1, \ldots, x_n)}$$

Using the naive conditional independence assumption that

$$P(x_i \mid y, x_1, \ldots, x_{i-1}, x_{i+1}, \ldots, x_n) = P(x_i \mid y),$$

for all i, this relationship is simplified to:

$$P(y \mid x_1, \ldots, x_n) = \frac{P(y) \prod_{i=1}^{n} P(x_i \mid y)}{P(x_1, \ldots, x_n)}$$

In spite of their apparently over-simplified assumptions, naive Bayes classifiers have worked quite well in many real-world situations, famously document classification and spam filtering. They require a small amount of training data to estimate the necessary parameters.

Naive Bayes learners and classifiers can be extremely fast compared to more sophisticated methods. The decoupling of the class conditional feature distributions means that each distribution can be independently estimated as a one dimensional distribution. This in turn helps to alleviate problems stemming from the curse of dimensionality.

**Multinomial Naïve Bayes:**

MultinomialNB implements the naive Bayes algorithm for multinomially distributed data, and is one of the two classic naive Bayes variants used in text classification (where the data are typically represented as word vector counts, although tf-idf vectors are also known to work well in practice). The distribution is parametrized by vectors $\theta_y = (\theta_{y1}, \ldots, \theta_{yn})$ for each class y, where n is the number of features (in text classification, the size of the vocabulary) and $\theta_{yi}$ is the probability $P(x_i \mid y)$ of feature i appearing in a sample belonging to class y.

The parameters $\theta_y$ is estimated by a smoothed version of maximum likelihood, i.e. relative frequency counting:

$$\hat{\theta}_{yi} = \frac{N_{yi} + \alpha}{N_y + \alpha n}$$

where $N_{yi} = \sum_{x \in T} x_i$ is the number of times feature i appears in a sample of class y in the training set T, and $N_y = \sum_i N_{yi}$ is the total count of all features for class y.

The smoothing priors $\alpha \geq 0$ accounts for features not present in the learning samples and prevents zero probabilities in further computations. Setting $\alpha = 1$ is called Laplace smoothing, while $\alpha < 1$ is called Lidstone smoothing.

The email system has been pretrained using Multinomial Naïve bayes classification algorithm of the sklearn library. The algorithm works quite well with the used dataset and gives an overall accuracy of more than 97% , which is clearly seen  on the mails it classifies as ham and spam in the inbox of users.
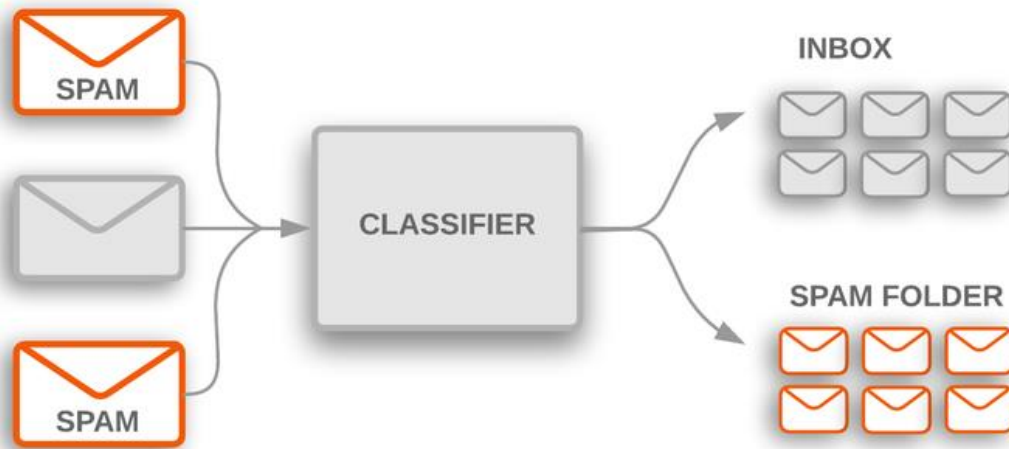
# PROPOSED SOLUTION ARCHITECTURE



FIGURE 3: BASIC ARCHITECTURE

When a message is received or sent via the proposed email system, it undergoes a series of steps on the backend.

1. The message is first preprocessed and tokenized.
2. Then the tokens undergo probability estimation.
3. Then the machine applies Multinomial Naïve bayes classifier to classify the messages as spam or ham.
4. If the message is classified as spam, then the message is tagged spam and moved to a separate inbox and if the message is classified as legitimate, it is tagged ham and is processed as usual.

(Naive Bayes classifiers work by correlating the use of tokens (typically words, or sometimes other things), with spam and non-spam e-mails and then using Bayes' theorem to calculate a probability that an email is or is not spam.)



FIGURE 4: GENERAL WORKING OF A CLASSIFIER

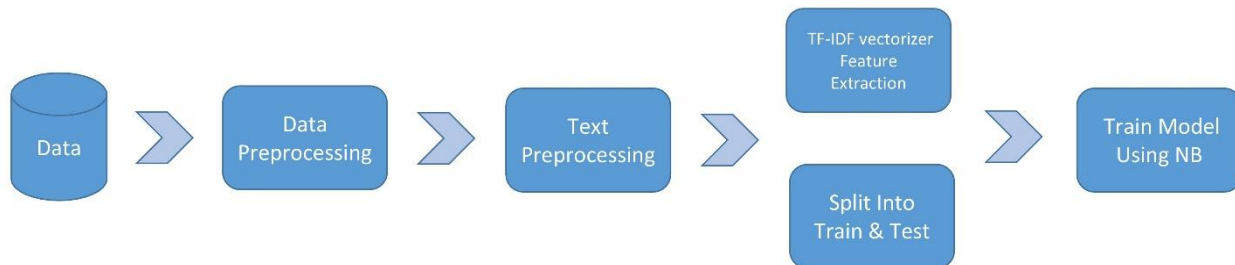The steps can be visualized in the following way:
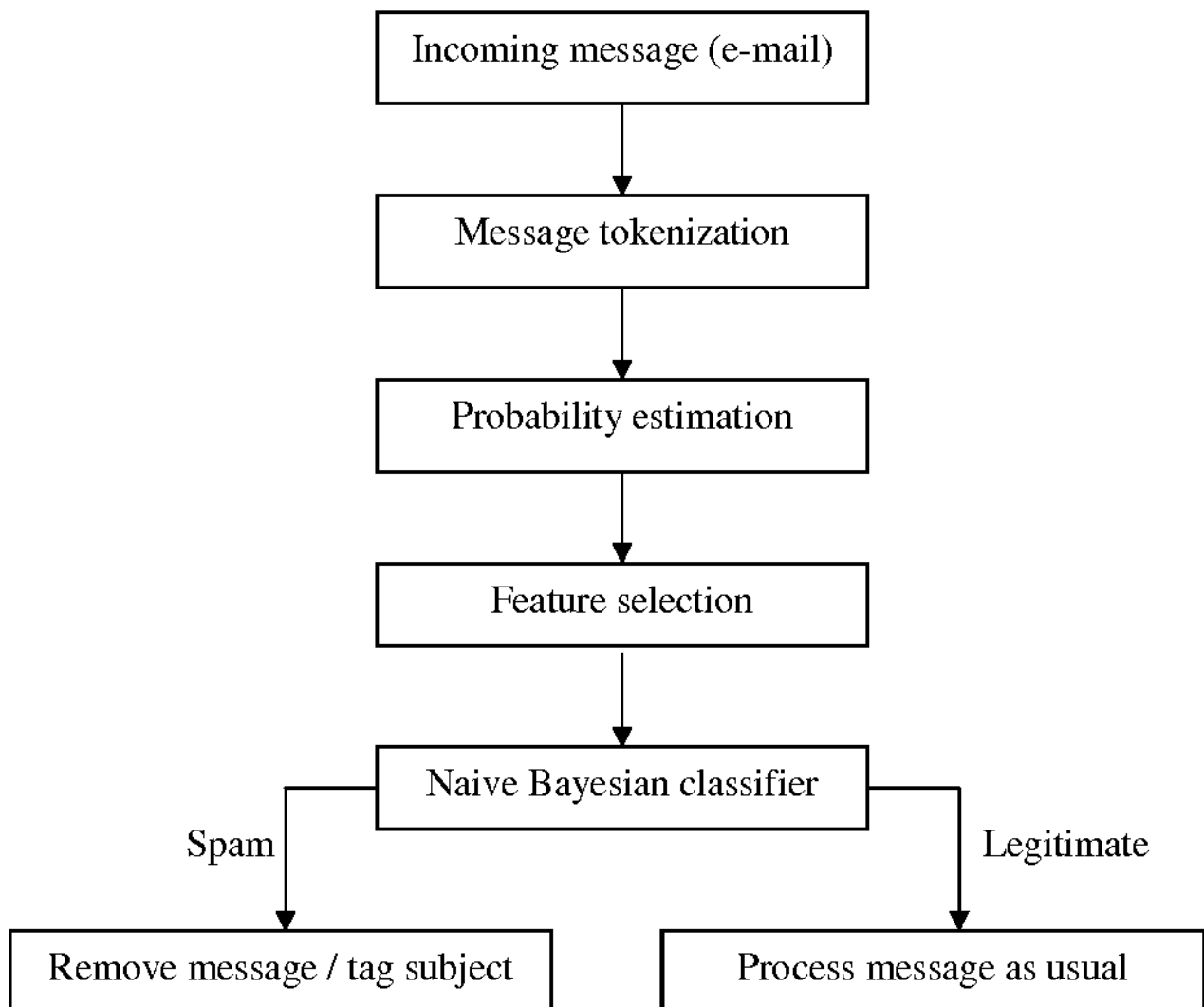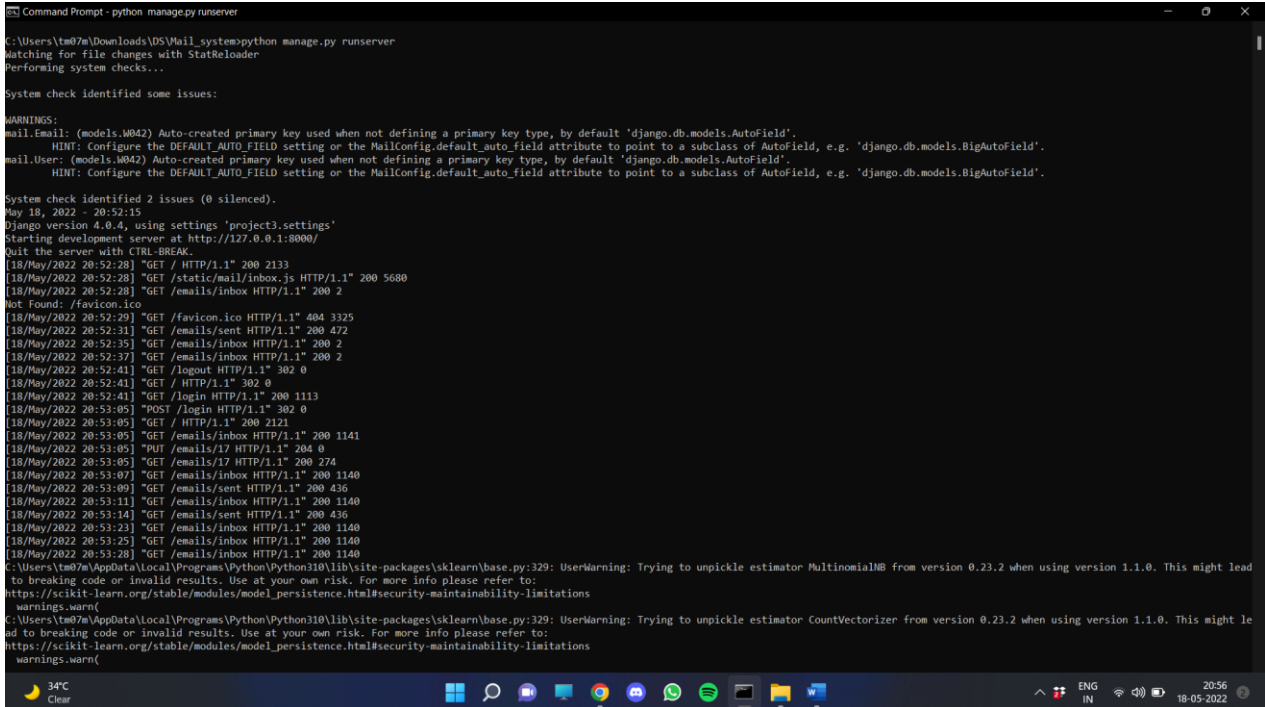


FIGURE 5: STEPS TO TRAIN THE MODEL



FIGURE 6: HOW AN INCOMING MAIL IS PROCESSED AS SPAM

# EXPERIMENTATION AND RESULTS



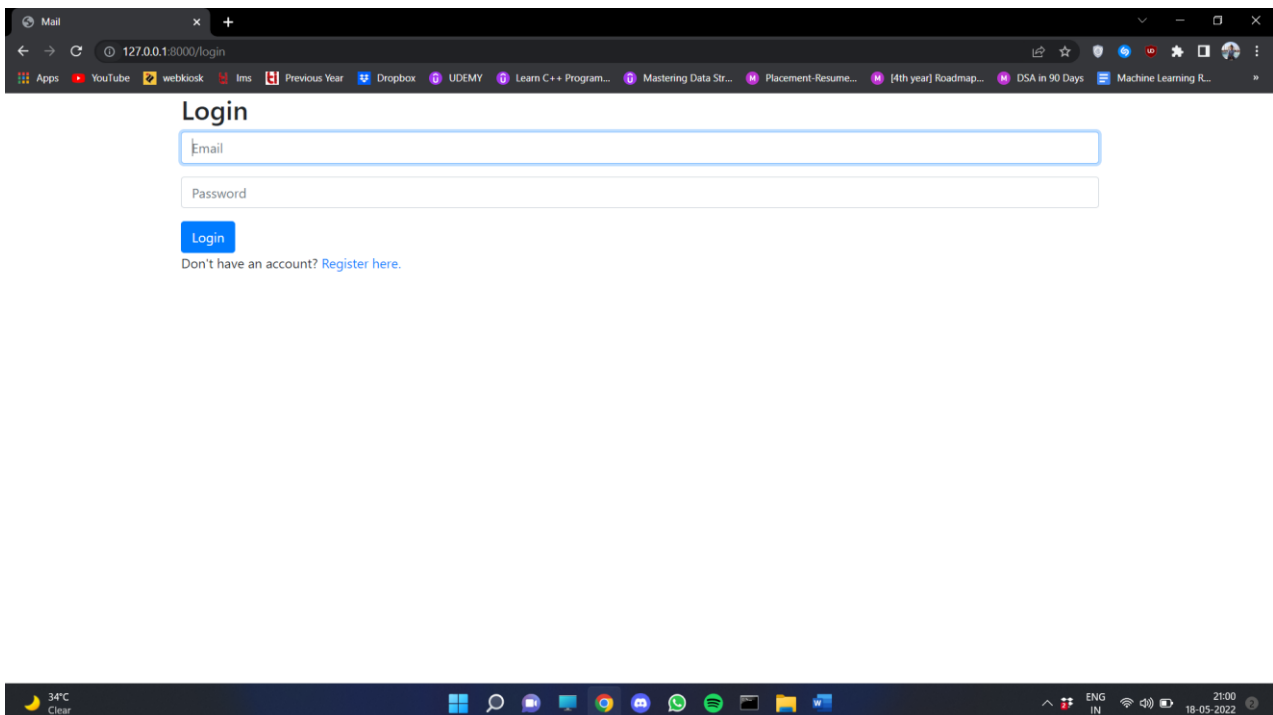FIGURE 7: RUNNING THE SERVER THROUGH COMMAND PROMPT TO ACCESS THE TERMINAL



FIGURE 8: LOGIN PAGE

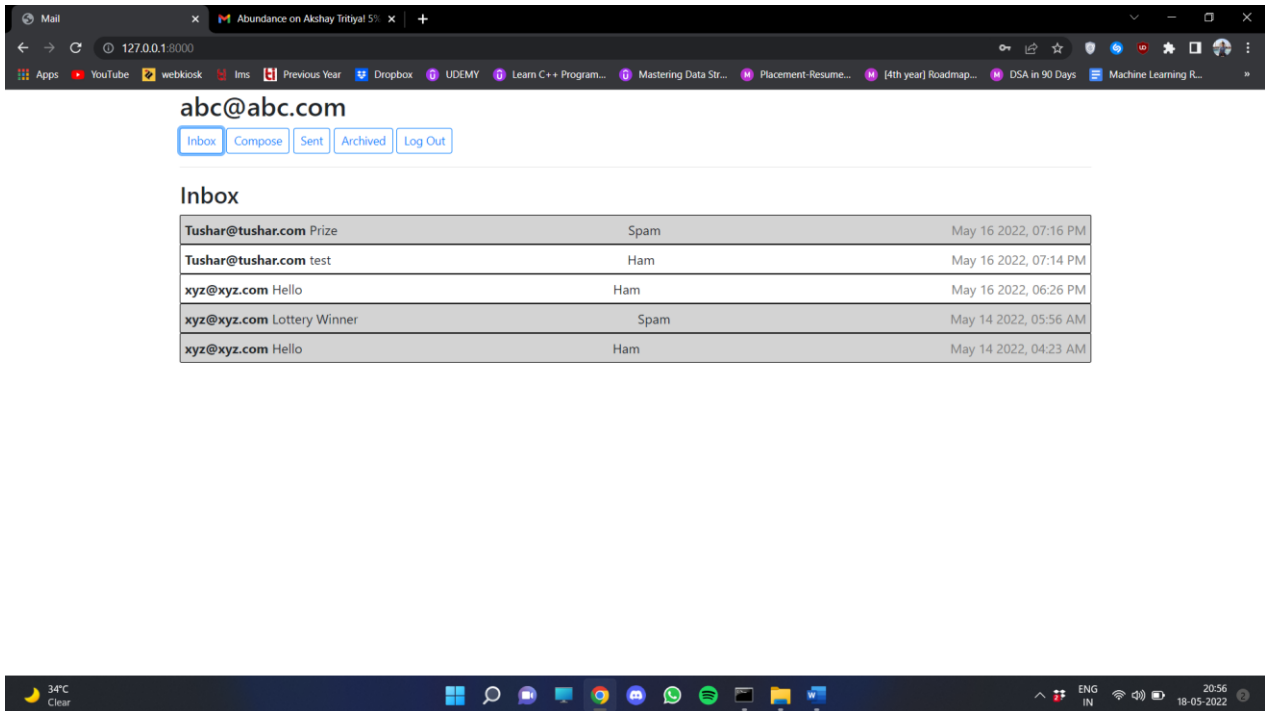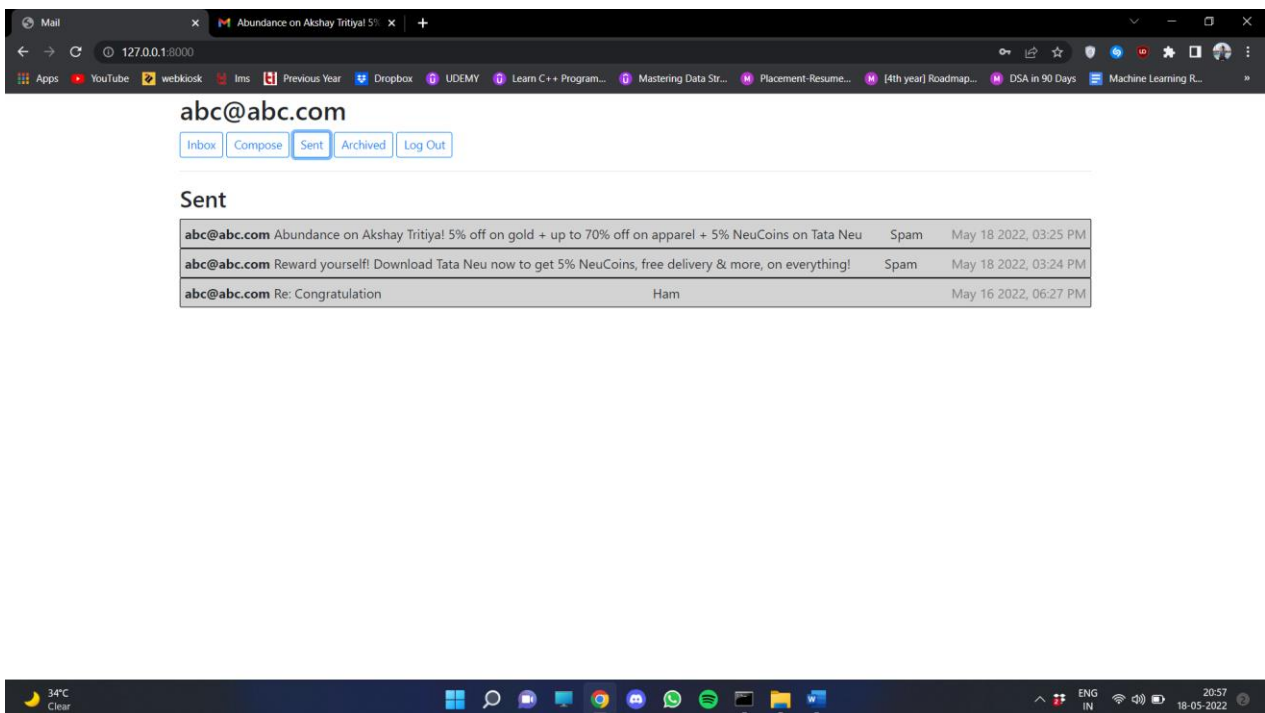FIGURE 9: SAMPLE USER 1-INBOX
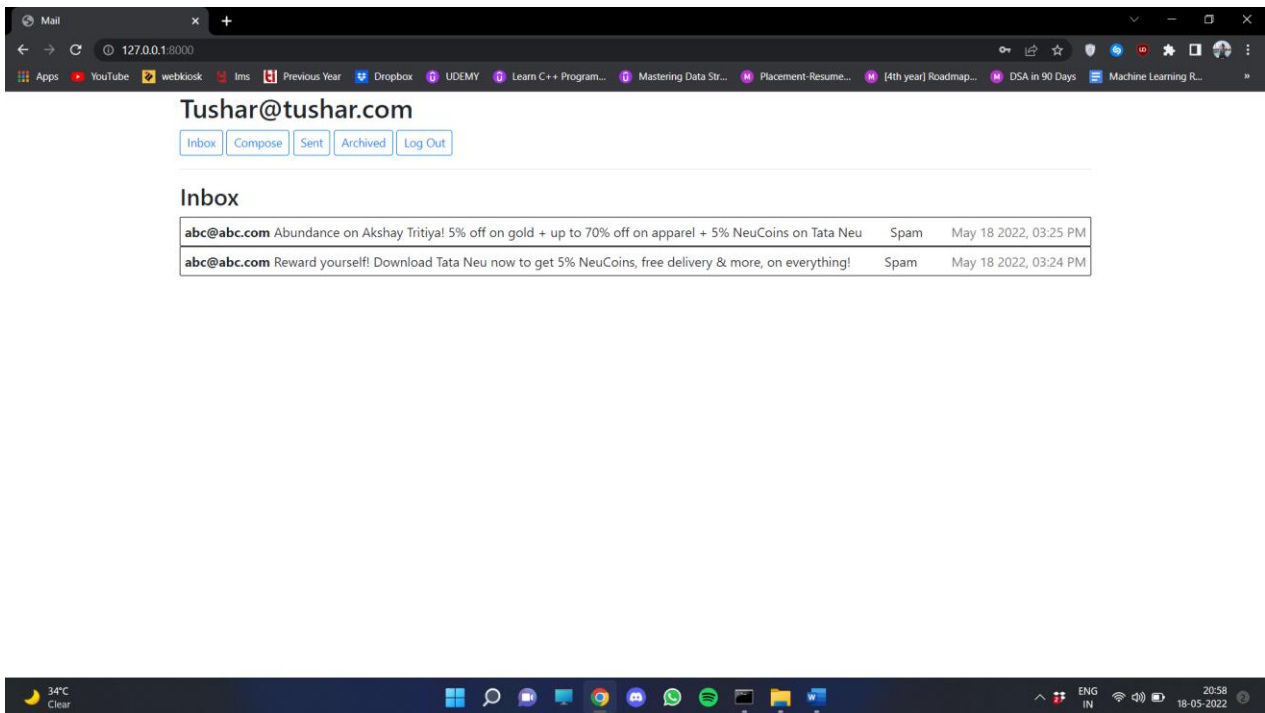


FIGURE 10: SAMPLE USER 1-SENT
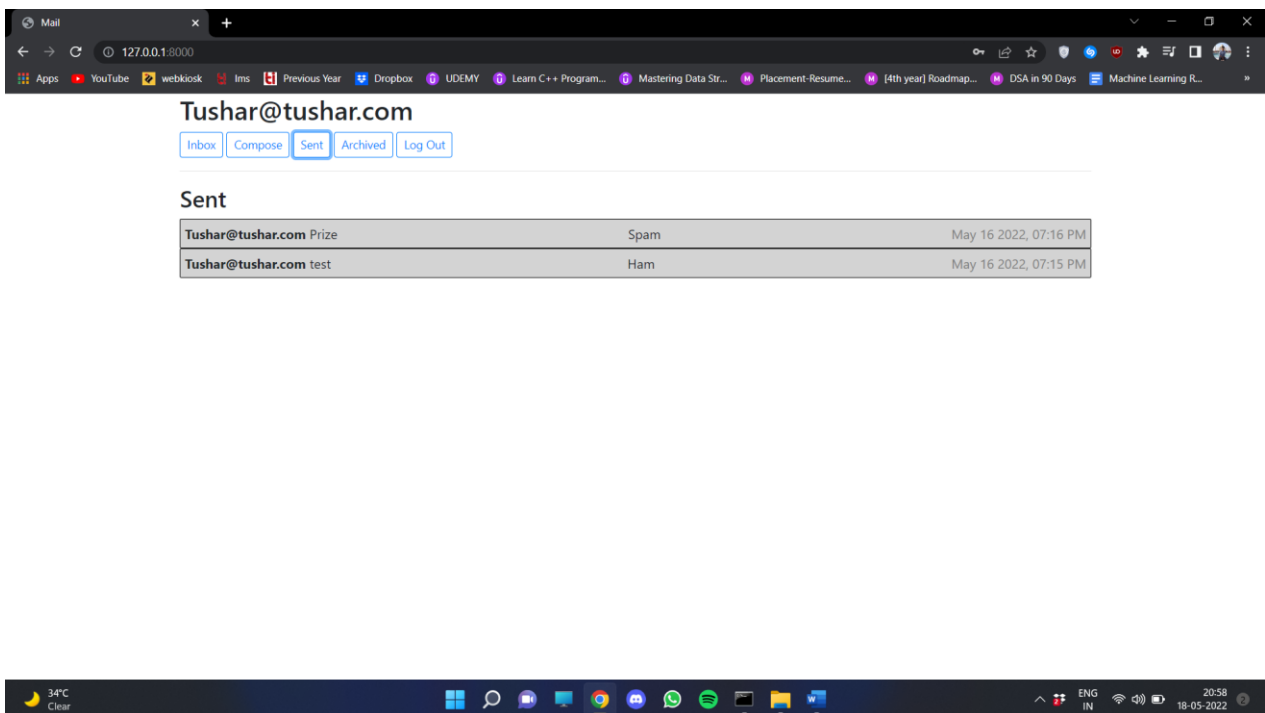
FIGURE 11: SAMPLE USER 2-INBOX



FIGURE 12: SAMPLE USER 2-SENT

# CONCLUSION AND FUTURE SCOPE

Spam filters detect unsolicited, unwanted, and virus-infested email (called spam) and stop it from getting into email inboxes. In today's world, spam filtering is a must to protect your business. Spam is not going away. It is estimated that 70 percent of all email sent globally is spam, and the volume of spam continues to grow.

In the proposed email system, spam filters are applied to both inbound email (email entering the network) and outbound email (email leaving the network).

Spam filters use "heuristics" methods, which means that each email message is subjected to a set of predefined rules (algorithms). Each rule assigns a numerical score to the probability of the message being spam, and if the score passes a certain threshold the email is flagged as spam and blocked from going further. There are different types of spam filters for different criteria:

- Content filters – parse the content of messages, scanning for words that are commonly used in spam emails.

- Header filters – examine the email header source to look for suspicious information (such as spammer email addresses).

- Blocklist filters – stop emails that come from a blocklist of suspicious IP addresses. Some filters go further and check the IP reputation of the IP address.

- Rules-based filters – apply customized rules designed by the organization to exclude emails from specific senders, or emails containing specific words in their subject line or body.

Any spam filtering solution cannot be 100 percent effective. However, a business email system without spam filtering is highly vulnerable, if not unusable. It is important to stop as much spam as you can, to protect your network from the many possible risks: viruses, phishing attacks, compromised web links and other malicious content. Spam filters also protect your servers from being overloaded with non-essential emails, and the worse problem of being infected with spam software that may turn them into spam servers themselves. By preventing spam email from reaching your employees' mailboxes, spam filters give an additional layer of protection to your users, your network, and your business.

The project can be accessed and view from my git repository at: https://github.com/tush7301/Mail_system

# REFERENCES

- https://scikit-learn.org/stable/modules/naive_bayes.html
- https://en.wikipedia.org/wiki/Naive_Bayes_spam_filtering
- https://en.wikipedia.org/wiki/Naive_Bayes_classifier
- https://www.geeksforgeeks.org/naive-bayes-classifiers/
- https://towardsdatascience.com/how-to-build-and-apply-naive-bayes-classification-for-spam-filtering-2b8d3308501
- https://github.com/ShubhamPy/Spam-Classifier