# CHAPTER 1
# INTRODUCTION

Email is a popular medium for communication, but it also poses a significant risk to users' security and privacy due to the high volume of spam emails. Spam emails are unwanted emails that are sent to a large number of recipients for advertising or other malicious purposes. Email spam can lead to several problems, including phishing attacks, malware distribution, and identity theft. Therefore, email spam detection is crucial in preventing these risks. Machine learning algorithms, especially Naive Bayes, have been widely used in email spam detection.

Naive Bayes is a probabilistic algorithm that uses Bayes' theorem to calculate the probability of a given email being spam or not. The Naive Bayes algorithm is simple, fast, and effective for text classification tasks, making it a popular choice for email spam detection.

In this paper, we propose a Naive Bayes-based approach for email spam detection. The proposed approach involves preprocessing the email data, feature extraction, and then applying the Naive Bayes algorithm for classification. We use a publicly available dataset to evaluate the performance of the proposed approach. The results show that our approach achieves a high accuracy rate in detecting spam emails. Therefore, the proposed approach can be used as an effective method for email spam detection, which can help in improving the overall user experience and reducing the risk of security breaches.

I hope this email finds you well. I am writing to provide you with an introduction to the concept of Naive Bayes and its application in email spam detection.Once trained, the Naive Bayes model can be used to classify new, unseen emails as spam or ham. The model calculates the probabilities of an email being spam or ham based on the observed features. It then assigns the label with the higher probability to the email.

One of the advantages of Naive Bayes is its simplicity and efficiency. It can quickly process large volumes of data and make predictions in real-time. Additionally, Naive Bayes performs well even with limited training data, making it suitable for spam detection applications where new types of spam continuously emerge.

In conclusion, Naive Bayes is a powerful and widely used algorithm for email spam detection. By leveraging Bayesian probability and statistical independence assumptions, it provides an effective and efficient means of filtering out unwanted spam emails from our inboxes. I hope this introduction to Naive Bayes for email spam detection has been informative. Should you have any further questions or require more detailed information, please do not hesitate to reach out.

## 1.1    BACKGROUND

Email spam detection has been an active research area in the field of natural language processing for several years. With the increasing volume of emails being sent every day, the need for an effective method to detect and filter out spam emails has become more critical.Several approaches have been proposed for email spam detection, including rule-based methods, content-based methods, and machine learning-based methods.

Machine learning-based methods have gained popularity in recent years due to their ability to learn from large datasets and identify patterns and features that are characteristic of spam emails. Naive Bayes algorithm is one of the most  commonly used machine learning algorithms for email spam detection.

It is a probabilistic algorithm that uses Bayes' theorem to calculate the probability of a given email being spam or not. In the Naive Bayes algorithm, each email is represented as a bag of words, and the probability of the email being spam is calculated based on the occurrence of words in the email. The algorithm assumes that the occurrence of words in an email is independent of each other, hence the name "Naive Bayes."

One of the advantages of the Naive Bayes algorithm is its simplicity and efficiency. It can be trained on a large corpus of emails and can quickly classify new emails as spam or not. Additionally, the algorithm can handle a   large number of features, making it suitable for high-dimensional data like text.

In summary, email spam detection using Naive Bayes algorithm is a popular and effective method that has been widely used in research and industry. The algorithm's simplicity and efficiency make it a popular choice for email spam detection, and it has been shown to achieve high accuracy rates in identifying spam emails.

## 1.2  MOTIVATION

The motivation for email spam detection using Naive Bayes algorithm stems from the need to protect users from the risks associated with spam emails. Email is one of the most commonly used communication channels, and the volume of emails sent every day is enormous. A significant portion of these emails are spam emails, which are unwanted and often malicious. Spam emails can contain phishing links, malware, and fraudulent messages, which can lead to several problems, including identity theft, financial loss, and damage to reputation. Email service providers and users need an effective method to detect and filter out spam emails to improve the user experience and prevent security breaches. Machine learning algorithms, especially Naive Bayes, have been shown to be effective in email spam detection. Naive Bayes is a probabilistic algorithm that can be trained on a large corpus of emails to identify patterns and features that are characteristic of spam emails. The motivation for email spam detection using Naive Bayes algorithm is to provide an effective and efficient method for detecting and filtering out spam emails. The proposed approach can help in improving the overall user experience and reducing the risk of security breaches, which is crucial in today's digital world.

## 1.3  OBJECTIVE

The objective of email spam detection using Naive Bayes is to develop a machine learning model that can accurately classify email messages as either spam or not spam (ham). Naive Bayes is a probabilistic algorithm that calculates the probability of a new email message belonging to either the spam or ham category based on the occurrence of certain keywords or features.

The goal is to use a set of labeled training data to teach the model to recognize patterns in the text of email messages that are associated with either spam or ham. Once the model has been trained, it can be used to classify new, unlabeled email messages with a high degree of accuracy. This can be particularly useful for individuals and organizations who receive a large volume of email and need to quickly and efficiently sort through their inbox to identify and prioritize important messages.

The objectives for spam email detection are typically to identify and classify incoming emails as either spam or legitimate (also known as "ham"). Here are some of the main objectives of a spam email detection system:

**Accurate classification**: The system should accurately identify and classify emails as spam or ham with high precision and recall.

**Minimizing false negatives**: False negatives are legitimate emails that are incorrectly classified as spam. Minimizing false negatives is important because it ensures that important messages are not lost.

**Minimizing false positives**: False positives are spam emails that are incorrectly classified as legitimate. Minimizing false positives is important because it reduces the amount of spam that users receive.

**Efficient processing:** The system should be able to process a large number of incoming emails quickly and efficiently.

**Easy to use:** The system should be easy for users to set up and use.

**Robust to changes:** The system should be robust to changes in spam tactics and be able to adapt to new forms of spam.

**Low cost:** The system should have a low cost, both in terms of time and resources, for both users and system administrators.

**Privacy:** The system should respect users' privacy and not collect or share any sensitive information.


## 1.4   SCOPE OF THE PROJECT

The scope of the project for email spam detection using Naive Bayes algorithm involves designing and implementing a system that can effectively detect and filter out spam emails. The project will involve the following tasks:

**1. Data Collection**: Collecting a dataset of emails that contains both spam and non-spam emails.

**2. Data Preprocessing:** Preprocessing the email data to remove unnecessary information, such as email headers, and converting the emails into a suitable format for further analysis.

**3. Feature Extraction**: Extracting relevant features from the preprocessed email data that can be used to classify emails as spam or non-spam. The features could include word frequency, presence of specific words or phrases, and other characteristics.

**4. Training the Naive Bayes Model**: Training a Naive Bayes model on the extracted features to learn to classify emails as spam or non-spam.

**5. Model Evaluation**: Evaluating the performance of the trained model using metrics such as accuracy, precision, recall, and F1 score.

**6. Deployment:** Deploying the trained model in a system that can effectively filter out spam emails in real-time.

The project's scope also includes testing the system on a real-world dataset of emails and comparing its performance with existing email spam detection methods.

# CHAPTER 2
# LITERATURE SURVEY

## 2.1 LITERATURE SURVEY

Several studies have been conducted on email spam detection using Naive Bayes algorithm. Here are some of the relevant literature surveys on the topic:

Nikhil Govil, Kunal Agarwal, Aashi Bansal "Email Spam Filtering: A Review", the authors provide a comprehensive review of different techniques used for email spam filtering, including Naive Bayes algorithm. The study highlights the effectiveness of Naive Bayes algorithm in email spam detection and discusses various features that can be used to improve its performance [1]

Nikhil Kumar,Sankat Sonawal,Nishant "A Comparative Study of Machine Learning Techniques for Email Spam Filtering" compares the performance of different machine learning algorithms, including Naive Bayes, in email spam filtering. The authors conclude that Naive Bayes algorithm is one of the most effective algorithms for email spam filtering, with high accuracy and low false positive rate [2]

Simran Gibson, Biju Isaac, Li Zhang "An Effective Email Spam Filtering System Using Machine Learning Techniques", the authors propose an email spam filtering system that uses Naive Bayes algorithm and other machine learning techniques. The study shows that Naive Bayes algorithm outperforms other machine learning algorithms in terms of accuracy and computational efficiency [3]

Gaiev Salim,Karimovich,Khamidov,Olimov Iskhandar"An Overview of Email Spam Filtering Techniques: Naive Bayes and Beyond" provides an overview of different techniques used for email spam filtering, with a focus on Naive Bayes algorithm. The authors discuss various preprocessing techniques and feature selection methods that can be used to improve the performance of Naive Bayes algorithm in email spam filtering [4]

Mansoor RAZA,Muhuna Ali Muslam,authors have highlighted several features contained in the email header which will be used to identify and classify spam messages efficiently .Those features are selected based on their performance in detecting spam messages. This paper also communalize each features contains in Yahoo mail, Gmail and Hotmail so a generic spam messages detection mechanism could be proposed for all major email providers [5]

Nandhini S, DR.Jeen, a new approach based on the strategy that how frequently words are repeated was used. The key sentences, those with the keywords, of the incoming emails have to be tagged and thereafter the grammatical roles of the entire words in the sentence need to be determined, finally they will be put together in a vector in order to take the similarity between received emails. K-Mean algorithm is used to classify the received e-mail. Vector determination is the method used to determine to which category the e-mail belongs to [6]

Mohammed Abudhalik,Abdul Malik , authors described about cyber attacks .Phishers and malicious attackers are frequently using email services to send false kinds of messages by which target user can lose their money and social reputations. These results into gaining personal credentials such as credit card number, passwords and some confidential data .In This paper ,authors have used Bayesian Classifiers .Consider every single word in the mail. Constantly adapts to new forms of spam [7]

Priya S,Annie Uthra R, proposed system attempts to use machine learning techniques to detect a pattern of repetitive keywords which are classified as spam. The system also proposes the classification of emails based on other various parameters contained in their structure such as Cc/Bcc, domain and header. Each parameter would be considered as a feature when applying it to the machine learning algorithm. The machine learning model will be a pre-trained model with a feedback mechanism to distinguish between a proper output and an ambiguous output. This method provides an alternative architecture by which a spam filter can be implemented. This paper also takes into consideration the email body with commonly used keywords and punctuations [8]

Fahima Hossain,Rajib Kuamr,Halder They describe a focused literature survey of Artificial Intelligence Revised (AI) and Machine learning methods for email spam detection. They have used the "image and textual dataset for the e-mail spam detection with the employment of various methods [9]

Olusanya Olatuniji,They have used methods of KNN algorithm, Reverse DBSCAN algorithm with experimentation on dataset. For the text recognition, OCR library" is employed but this OCR doesn't perform well. [10]

Overall, the literature survey suggests that Naive Bayes algorithm is an effective method for email spam detection.The algorithm is simple, efficient, and can achieve high accuracy rates in identifying spam emails. Additionally, researchers have proposed various techniques to improve the performance of Naive Bayes algorithm in email spam filtering, such as feature selection and preprocessing techniques.

## 2.2  INFERENCE OF THE LITERATURE REVIEW

The literature review for email spam detection using Naive Bayes reveals several key findings and insights:

1. Naive Bayes Algorithm: The Naive Bayes algorithm has been widely used for email spam detection due to its simplicity, efficiency, and effectiveness. It is based on the principle of Bayes' theorem, which calculates the probability of an event based on prior knowledge. Despite its "naive" assumption of independence between features, Naive Bayes has shown competitive performance in spam classification tasks.

2. Feature Selection: The choice of features plays a crucial role in email spam detection. Literature suggests that using both content-based features (such as keywords, frequency of certain words, and text patterns) and header-based features (sender information, subject line, etc.) can improve the accuracy of spam detection. Feature selection techniques, such as information gain, chi-square, and mutual information, have been applied to select the most informative features for classification.

3. Preprocessing Techniques: Preprocessing techniques are employed to enhance the quality of data before feeding it into the Naive Bayes classifier. Common preprocessing steps include tokenization, stemming or lemmatization, stop-word removal, and handling of special characters and HTML tags. These steps help to reduce noise, normalize the data, and improve the generalization of the classifier.

4. Training and Evaluation: To train a Naive Bayes spam classifier, a labeled dataset consisting of both spam and non-spam emails is required. The dataset is typically divided into training and testing sets. Cross-validation techniques, such as k-fold cross-validation, are often used to evaluate the performance of the classifier and estimate its generalization ability. Evaluation metrics such as accuracy, precision, recall, and F1 score are commonly used to measure the effectiveness of the classifier.

5. Challenges and Limitations: While Naive Bayes is a popular algorithm for email spam detection, it does have limitations. The assumption of feature independence may not hold true for all datasets, and the algorithm can be sensitive to the presence of rare or unseen features. Moreover, Naive Bayes may struggle with handling imbalanced datasets where the number of spam and non-spam emails is significantly different.

# CHAPTER 3
# SYSTEM ANALYSIS

## 3.1 PROBLEM DEFINATION

Problem definition for email spam detection using Naive Bayes algorithm is to design and implement a system that can accurately distinguish between spam and non-spam emails. The system should be able to analyze incoming emails and classify them as either spam or non-spam based on their content.

The goal is to create a system that can effectively filter out unwanted emails and improve the overall user experience. Spam emails can be a major source of annoyance and frustration for email users, and can also pose security risks if they contain malicious content.

The system should be able to learn from a dataset of emails that contains both spam and non-spam emails. The Naive Bayes algorithm should be used to train a model on the dataset, which can then be used to classify new incoming emails. The model should be able to identify relevant features from the email data, such as the frequency of specific words or phrases, that can help distinguish between spam and non-spam emails.

The accuracy of the system should be evaluated using appropriate metrics, such as precision, recall, and F1 score. The goal is to design and implement a system that can achieve high accuracy rates in identifying spam emails while minimizing false positives and false negatives.

Overall, the problem definition for email spam detection using Naive Bayes algorithm is to create an effective system that can filter out unwanted emails and improve the overall email experience for users.

## 3.2 PROPOSED SYSTEM

Naive Bayes is a popular and simple machine learning algorithm that is often used for spam email detection. The algorithm is based on Bayes' Theorem, which states that the probability of a given event (in this case, that an email is spam) is based on prior knowledge of related events (in this case, the words and phrases in the email).

Naive Bayes algorithms are simple and computationally efficient, making them well suited for large datasets. There are two main types of Naive Bayes algorithms used for

spam email detection: Multinomial Naive Bayes and Bernoulli Naive Bayes.

Multinomial Naive Bayes is based on word frequencies, counting the number of times each word appears in a message. The algorithm calculates the probability that a message is spam based on the presence of specific words or phrases that are commonly associated with spam.

Bernoulli Naive Bayes is based on the presence or absence of words, rather than the frequency of their occurrence. This approach is often used when the size of the data is limited and the focus is on detecting specific keywords or phrases.

In both cases, the algorithm trains on a dataset of labeled messages (spam and ham) and uses this training data to make predictions about new incoming messages. The prediction is made by calculating the probability that a given message is spam based on the presence or absence of certain words and phrases.

Overall, Naive Bayes is a fast and effective algorithm for spam email detection, and it is widely used in many commercial and open-source spam filtering systems

## 3.3 ADVANTAGES

**Simple and fast**: Naive Bayes algorithms are simple and computationally efficient, making them well suited for large datasets and real-time applications.

**Good accuracy**: Naive Bayes algorithms have been shown to perform well in various classification tasks, including spam email detection. This is due to their ability to make accurate predictions based on a relatively small amount of training data.

**Robust to irrelevant features:** Naive Bayes algorithms are not sensitive to irrelevant features, which can be a problem for other machine learning algorithms. This means that Naive Bayes algorithms can still make accurate predictions even if the data contains irrelevant or redundant features.

**Handle missing data**: Naive Bayes algorithms can handle missing data well, since they calculate probabilities based on the presence or absence of individual features, rather than their absolute values.

**Easy to implement:** Naive Bayes algorithms are easy to implement and understand, making them accessible to practitioners with a limited background in machine learning.

# 3.4 SYSTEM SPECIFICATION

Here are some of the system specifications that may be considered for email spam detection using Naive Bayes algorithm:

1. **Programming Language**: The system can be implemented using programming languages such as Python or Java.

2. Dataset: A suitable dataset of emails that contains both spam and non-spam emails will be required for training and testing the Naive Bayes model.

3. **Preprocessing Techniques**: The emails will need to be preprocessed to remove any unnecessary information like email headers and convert the emails into a suitable format for analysis. Techniques like tokenization and stemming can be used to standardize the text and reduce its complexity.

4. **Feature Selection Techniques:** Feature selection techniques like mutual information, chi-square test, and information gain can be used to identify relevant features from the email data that can improve the accuracy of the Naive Bayes model.

5. **Naive Bayes Algorithm:** The Naive Bayes algorithm will be used to train a model on the dataset of emails. The algorithm can be implemented using different variants like Multinomial Naive Bayes, Bernoulli Naive Bayes, or Gaussian Naive Bayes depending on the type of features being used.

6. **Evaluation Metrics**: Appropriate evaluation metrics such as precision, recall, and F1 score will be used to evaluate the accuracy of the system in identifying spam and non-spam emails.

7**. User Interface:** The system can be integrated with an email client to provide users with a seamless experience. A user interface can be provided to allow users to customize the spam filter settings and manage their email preferences.

Overall, the system specifications for email spam detection using Naive Bayes algorithm will depend on the specific requirements of the project. The system can be designed to be flexible and scalable to handle large volumes of emails and to adapt to changing email patterns over time.

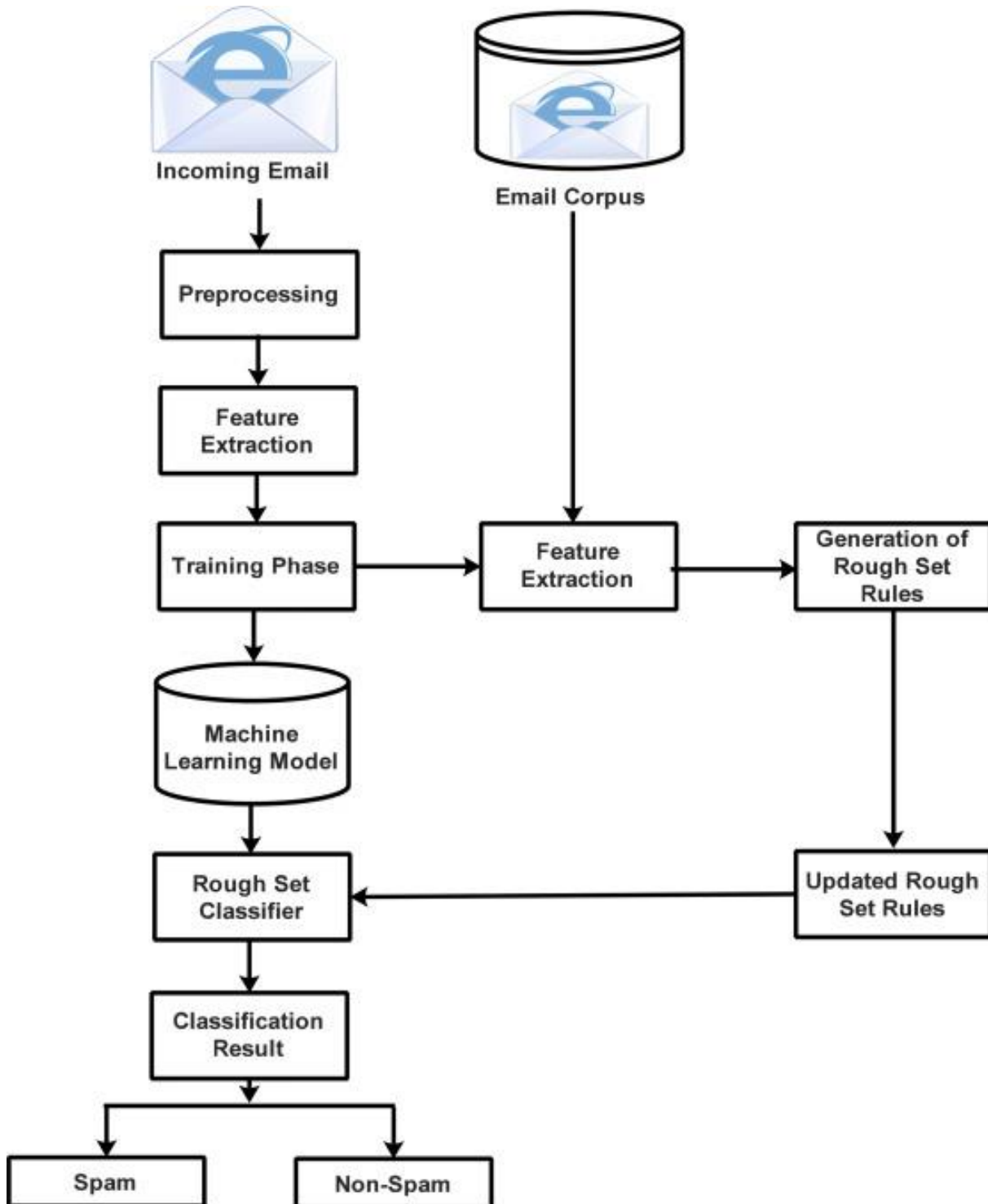# CHAPTER 4
# SYSTEM DESIGN

## 4.1 OVERALL ARCHITECTURE



**FIGURE 4.1**

The overall architecture of email spam detection using naive Bayes algorithm can be broken down into the following steps:

1.Data Preprocessing: The first step in building a spam detection model is to preprocess the data. This involves cleaning and preparing the email data by removing unwanted characters, converting all letters to lowercase, and tokenizing the text.

2. Feature Extraction: The next step is to extract features from the preprocessed email data. The most common approach is to use the Bag-of-Words (BoW) model, which represents each email as a vector of word counts. Other feature extraction techniques include TF-IDF (Term Frequency-Inverse Document Frequency) and n-grams.

3. Training: Once the features have been extracted, the next step is to train the model. In the case of naive Bayes, the model is trained by estimating the probability of each feature occurring in spam and non-spam emails. This is done by calculating the conditional probability of each feature given the class (spam or non-spam).

4. Classification: After the model has been trained, it can be used to classify new emails as either spam or non-spam. This is done by calculating the probability of each email belonging to each class (spam or non-spam) using Bayes' theorem. The class with the highest probability is then assigned to the email.

5. Evaluation: The final step is to evaluate the performance of the model. This is typically done by splitting the data into training and testing sets and calculating metrics such as accuracy, precision, recall, and F1 score.

Overall, the naive Bayes algorithm is a simple yet effective method for email spam detection, and its implementation can be easily adapted to different types of text classification problems.

## 4.2 DATASET DESCRIPTION

The choice of dataset for email spam detection using naive Bayes can have a significant impact on the performance of the model. Ideally, the dataset should be representative of the types of emails that the model will be used to classify and should contain a sufficient number of examples of both spam and non-spam emails. Here are some common characteristics of a dataset used for email spam detection using naive Bayes:

1. Source: The dataset can be obtained from various sources, including publicly available repositories such as the UCI Machine Learning Repository, Kaggle, or from email providers.

2. Size: The size of the dataset can vary depending on the application and the available resources. Typically, a dataset with at least several thousand examples is required for effective training and testing of the model.

3. Class distribution: The dataset should contain a balanced representation of both spam and non-spam emails. A dataset with imbalanced class distribution can result in a biased model, where the model may predict the majority class for all examples.

4. Data format: The dataset should be in a format that can be easily processed by the model. Common formats include CSV, TSV, and JSON.

5. Data attributes: The dataset should contain relevant attributes for spam detection, such as email subject, body, sender, recipient, and date. Additional attributes may include email attachments, URLs, and HTML content.

6. Preprocessing: The dataset may require preprocessing to remove noise and irrelevant data. Common preprocessing techniques include removing stop words, stemming, and lemmatization.

7. Annotated data: The dataset may require manual annotation of the email data as spam or non-spam. This is typically done by human annotators and can be time-consuming and expensive.

Overall, a good dataset for email spam detection using naive Bayes should be representative of the target application, contain a balanced distribution of both spam and non-spam emails, and be in a format that is easily processed by the model.

## 4.3 MODULE DESCRIPTION

## 4.3.1 DATA COLLECTION

The first step of Module description is collection of spam data. Get the spam data. Data is essential ingredient before we can develop any meaningful algorithm. .For this project We are collecting data set from the website called Kaggle. The data which is collected Consist of spam and ham data. Further  more, in the ham data, they are easy and hard. Which mean, there is some non-spam data that has a very high similarity with spam data.

## 4.3.2 IDENTIFY THE SPAM EMAIL

There are some specific features to identity the received mail is spam / not. some of the features of identity the spam mail are
1. If the received mail is asking personal information
2. If mail consist of urgent / Threatening messages in the subject
3. Spam mails mostly consist of typos/strange phrasing(sending mails from unauthorized website by making small change in the website name.)
4. If the received mail consist of large storage
5. Is the sender know to you and does the email address match the name?

## 4.3.3 IMPORTING LIBRARIES

The python libraries that are used in the code as follows

• NumPy

• Pandas

• SK-learn

## NUMPY

NumPy is used for working with arrays. It also has functions for working in domain of linear algebra , Fourier transform , and matrices. NumPy stands for Numerical Python.

## PANDAS

Pandas used for working with data sets. It has a functions for analyzing ,cleaning , exploring , and manipulating data. The name "Pandas" has a reference to both "Panel Data", and "Python Data Analysis"

## SKLEARN

Scikit-learn (SK-learn) is the most useful and robust library for machine learning in python. It provides a selection of efficient tools for machine learning and statistical modelling including classification, regression , clustering and dimensionality reduction via a consistence interface in Python

## 4.3.4 CLASSIFICATION OF SPAM MAILS.

In this all the received mails are classified into ham and spam mails. Ham consist of useful information whereas spam consist of unsolicited messages. The classification is can be done by supervised learning technique. In this technique using a trained model .receiving mail can be labelled as spam or ham. If the mail is spam it is moved to spam folder
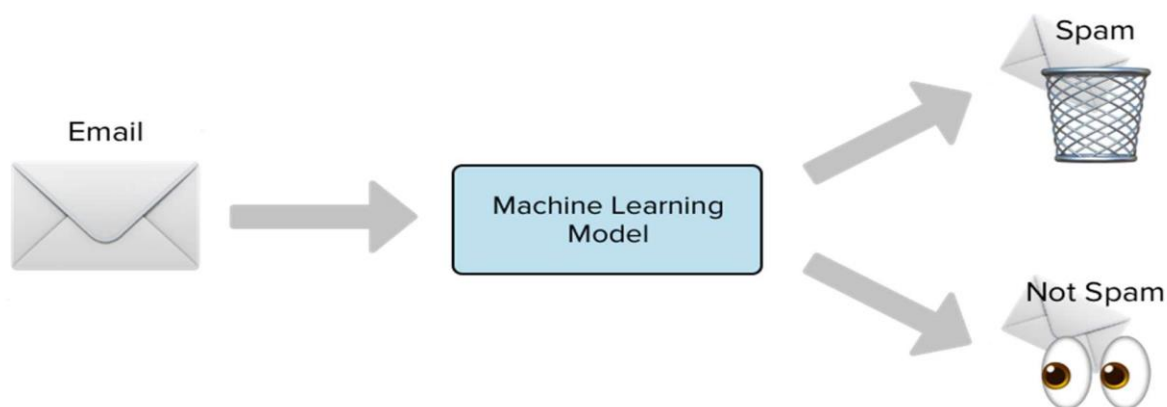


**FIGURE 4.3.4**

## 4.3.5 VERIFYING THE SPAM MAIL

Verifying the spam mail we use naive bayes algorithm work by bayes classifier algorithm. Naive bayes classifier works by correlating the use of tokens (typically words, or sometimes other things) with spam and not spam emails and then using bayes theorem to calculate a probability that an email is or not spam . Sometimes informative mails are stored in spam folder. For classifying those we use naive bayes algorithm. It works by comparing the incoming mail with a dataset of emails, which are categorized into spam and ham.

## 4.3.6 DENY THE SPAM EMAILS

 we are classifying into spam and ham mails. If the mail is found to be spam it enters spam folder. If the mail is found to be ham it deny from the spam folder. now if the mail present in spam folder is informative mail it moves to inbox deny from spam folder. if the mail present in spam folder is not informative mail allow spam folder.

# CHAPTER 5
# SYSTEM IMPLEMENTATION

## 5.1 REQIREMENT SPECIFICATION

**1. Problem Statement:** The purpose of this project is to develop a system that can accurately detect spam emails using the Naive Bayes algorithm.

**2. Functional Requirements:**

   a. The system shall be able to read email messages from a specified email folder.

   b. The system shall be able to classify the email messages into either spam or non-spam categories.

   c. The system shall use the Naive Bayes algorithm to determine the probability of each email being spam.

   d. The system shall use a pre-defined set of features (such as word frequency, sender information, etc.) to classify the emails.

   e. The system shall allow users to manually classify emails as spam or non-spam, to improve the accuracy of the system over time.

   f. The system shall provide a user interface for users to interact with the system, view classified emails, and manually classify emails.

**3. Non-functional Requirements:**

   a. The system shall be able to process emails in real-time with minimal delay.

   b. The system shall be able to handle large volumes of emails.

   c. The system shall be able to run on a variety of operating systems.

   d. The system shall be user-friendly, with a simple and intuitive user interface.

   e. The system shall be accurate, with a high precision and recall rate.

   f. The system shall be secure, with appropriate measures to protect user data and prevent unauthorized access.

g. The system shall be reliable, with minimal downtime and a low error rate.

**4. Performance Requirements:**

a. The system shall be able to classify emails with an accuracy of at least 95%.

b. The system shall be able to process at least 100 emails per minute.

c. The system shall be able to handle at least 10,000 emails per day.

d. The system shall respond to user input within 2 seconds.

**5. Constraints:**

a. The system shall only classify emails written in English.

b. The system shall not classify emails that contain highly sensitive or personal information, such as medical records or financial information.

c. The system shall not modify or delete any emails in the user's email account.

d. The system shall comply with all relevant laws and regulations regarding email privacy and data protection.

**6. Assumptions:**

a. The system assumes that the user has access to an email account with sufficient permission to read and classify emails.

b. The system assumes that the user has a basic understanding of the Naive Bayes algorithm and its implementation for spam detection.

## 5.2 HARDWARE REQUIREMENTS

- System: intel i5 Processor
- Monitor: 15'' LED
- Input Devices: Keyboard, Mouse
- Ram: 8 GB
- CPU: 2 GHz or faster

- For a server handling 20,000 emails/day, a 500MHz CPU and 512MB of RAM is the minimum recommended.

- For a server handling 200,000 emails/day, a dual core 2GHz Xeon CPU and 4GBMB of RAM is the minimum recommended.

- For a server handling 2 million emails/day, two dual-core 3GHz Xeon CPUs and 4GBMB of RAM is the minimum recommended.

## 5.3 SOFTWARE REQUIREMENTS

- Operating system: windows 7 and above or Linux based OS or MAC OS
- Coding Language: Python
- Python 3.5 in Google Collab is used for data pre-processing, model training and prediction.

## 5.3 TECHNOLOGIES USED

## 5.3.1 MACHINE LEARNING

Machine learning is a growing technology which enables computers to learn automatically from past data. Machine learning uses various algorithms for building mathematical models and making predictions using historical data or information. Currently, it is being used for various tasks such as image recognition, speech recognition, email filtering, Facebook autotagging, recommender system, and many more. Machine Learning is said as a subset of artificial intelligence that is mainly concerned with the development of algorithms which allow a computer to learn from the data and past experiences on their own. With the help of sample historical data, which is known as training data, machine learning algorithms build a mathematical model that helps in making predictions or decisions without being explicitly programmed. Machine learning brings computer science and statistics together for creating predictive models. Machine learning constructs or uses the algorithms that are learned from historical data. The more we provide the information, the higher will be the performance.

## 5.3.2 HOW DOES MACHINE LEARNING WORK

A Machine Learning system learns from historical data, builds prediction models, and whenever it receives new data, predicts the output for it. The accuracy of predicted output depends upon the amount of data, as the huge amount of data helps to build a better model which predicts the output more accurately.
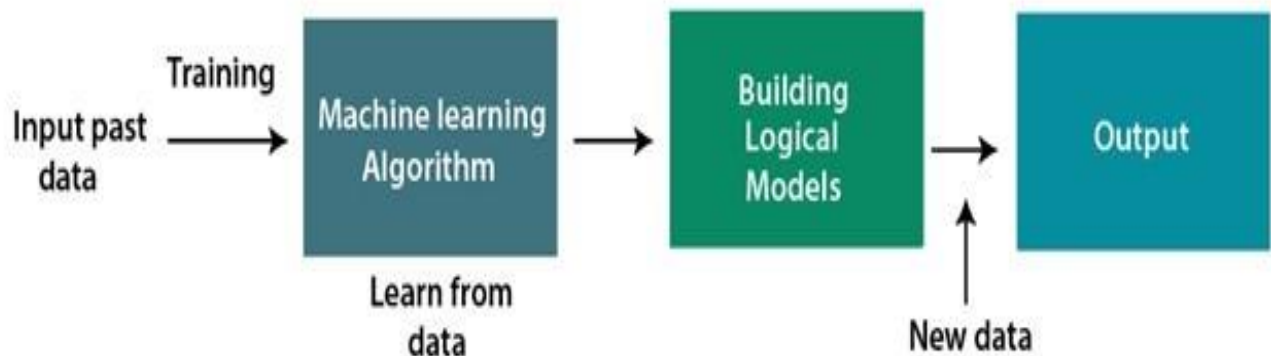


**FIGURE 5.3.2**

## 5.3.3 NAIVE BAYES CLASSIFIER ALGORITHM

Naive Bayes is a supervised algorithm, trained with the help of training data. It is a comparatively fast classification algorithm related to others. The Naive Bayes algorithm is a supervised learning algorithm, which is based on Bayes theorem and used for solving classification problems. It is mainly used in text classification that includes a high-dimensional training dataset. Naive Bayes Classifier is one of the simplest and most effective Classification algorithms which helps in building the fast machine learning models that can make quick predictions. It is a probabilistic classifier, which means it predicts based on the probability of an object. Some popular examples of Naive Bayes Algorithm are spam filtration, Sentimental analysis, and classifying articles.

### 5.3.4 WORKING OF NAÏVE BAYES ALGORITHM

A training dataset of known spam and non-spam email is used to train the algorithm

1. This algorithm uses probabilistic calculations to classify the email messages.

2. To calculate the probability of each word appearing in spam and non-spam mail.

3. It compares this probability to a predefined threshold. If the probability is above the threshold, the email is classified as spam. If the probability is below the threshold, the email is classified as non-spam.

### 5.4 PYTHON (Python 3.4)

Python is an interpreter, high-level, general-purpose programming language. Created by Guido van Rossum and first released in 1991, Python's design philosophy emphasizes code readability with its notable use of significant whitespace. Its language constructs and object-oriented approach aim to help programmers write clear, logical code for small and large-scale projects. Python is dynamically typed and garbage-collected. It supports multiple programming paradigms, including procedural, object-oriented, and functional programming. Python is often described as a "batteries included" language due to its comprehensive standard library. Python was conceived in the late 1980s as a successor to the ABC language. Python 2.0, released in 2000, introduced features like list comprehensions and a garbage collection system capable of collecting reference cycles. Python 3.0, released in 2008, was a major revision of the language that is not completely backward-compatible, and much Python 2 code does not run unmodified on Python 3. Python interpreters are available for many operating systems. A global community of programmers develops and maintains CPython, an open-source reference implementation. A non-profit

organization, the Python Software Foundation, manages and directs resources for Python and CPython development. The following figure illustrates how this works.
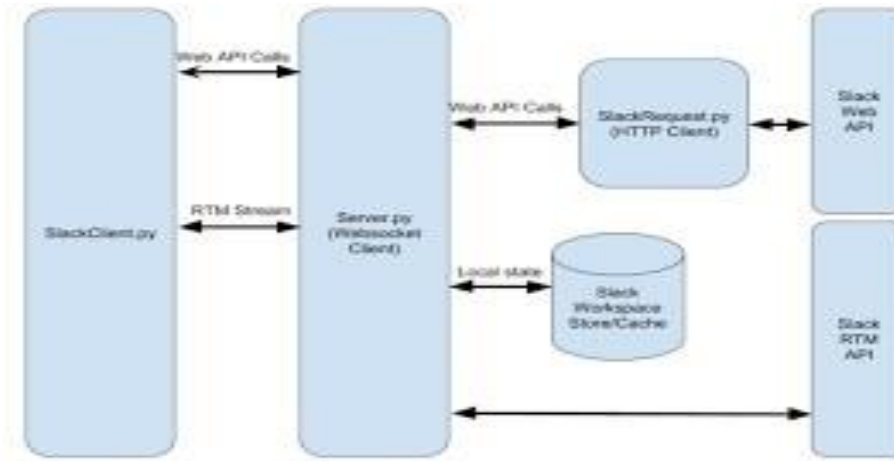
**WORKING MECHANISM OF PYTHON**



**FIGURE 5.4**

As a Machine Learning Engineer, I have been using Python for more than a year. Recently, I have also started learning C++, for fun. It made me realize how easy and intuitive Python is. I got more curious about how Python is different from other languages and its working. In this blog, I try to scratch Python's inner working.

Python started as a hobby project by Guido Van Rossum and was first released in 1991. A general purpose language, Python is powering large chunk of many companies like Netflix and Instagram. In Guido compares Python to languages like Java or Swift and says that while the latter two are a great choice for software developers — people whose day job is programming, but Python was made for people whose day job has nothing to do with software development but they code mainly to handle data
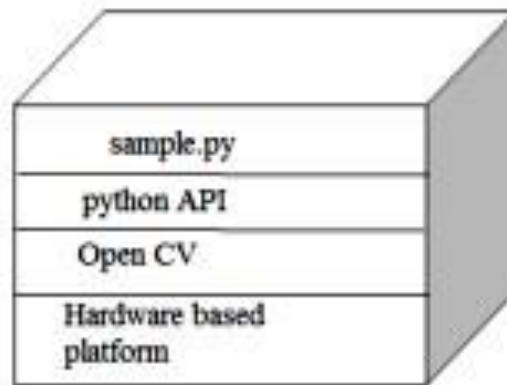
# PYTHON PLATFORM



## FIGURE 5.4

## 5.5 CODING

import numpy as np # linear algebra

import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)

# Input data files are available in the read-only "../input/" directory

# For example, running this (by clicking run or pressing Shift+Enter) will list all files under the input directory

import os

for dirname, _, filenames in os.walk('/kaggle/input'):

   for filename in filenames:

      print(os.path.join(dirname, filename))

```python
from google.colab import drive

drive.mount('/content/drive/')


data=pd.read_csv('/content/drive/MyDrive/Dataset/spam.csv')


data.columns

data.info()

data.isna().sum()


data['Spam']=data['Category'].apply(lambda x:1 if x=='spam' else 0)

data.head(5)


from sklearn.model_selection import train_test_split

X_train,X_test,y_train,y_test=train_test_split(data.Message,data.Spam,test_size=0.25)

from sklearn.feature_extraction.text import CountVectorizer

from sklearn.naive_bayes import MultinomialNB

from sklearn.pipeline import Pipeline
```

```python
clf=Pipeline([

    ('vectorizer',CountVectorizer()),

    ('nb',MultinomialNB())

])

clf.fit(X_train,y_train)


clf.score(X_test,y_test)


mails=['due to lot of request-,gold membership offer extended only today-final
call.','masterclass on database management system-mongoDB and SQL']


clf.predict(mails)


for x in mails:

  y= [x]

  if clf.predict(y)==[0]:

    print("ham")

  else:

    print ("spam")
```

# CHAPTER 6

## RESULT AND DISCUSSION

## SCREENSHOT-1

## DATA INFORMATION

```
data.columns
data.info()

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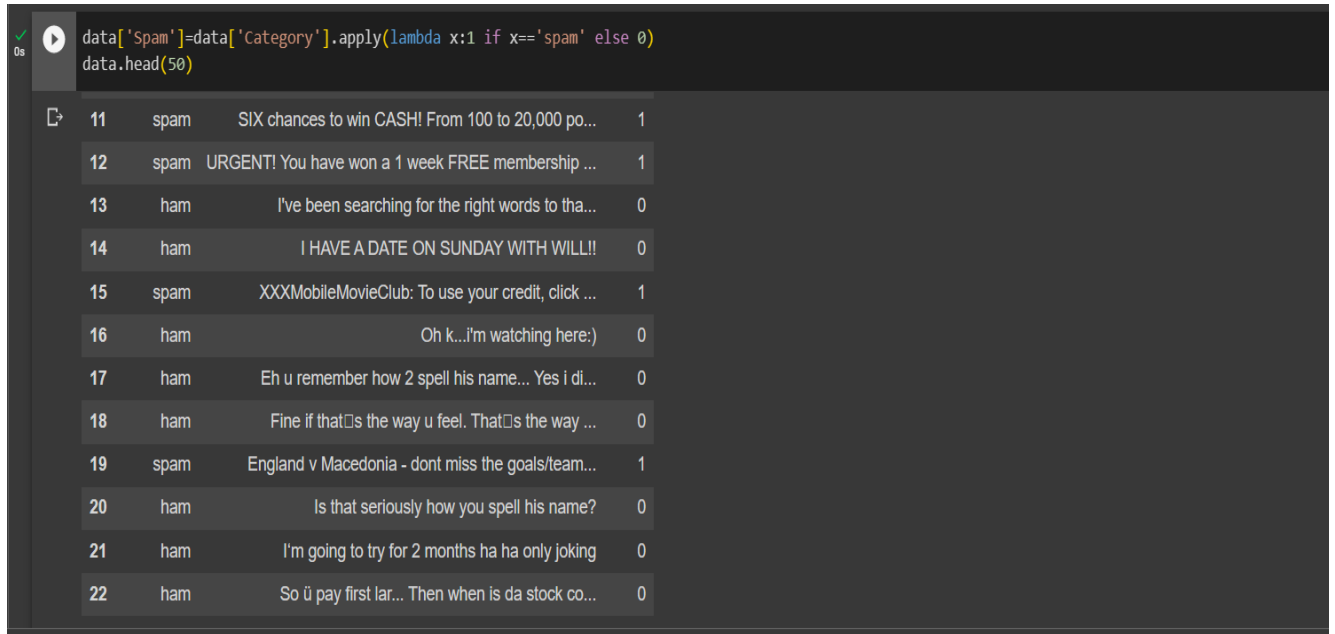
## SCREENSHOT-2

## DROPPED THE COLUMN UNNAMED:0

```
[6] data.isna().sum()

    Category    0
    Message     0
    dtype: int64
```

# SCREENSHOT-3

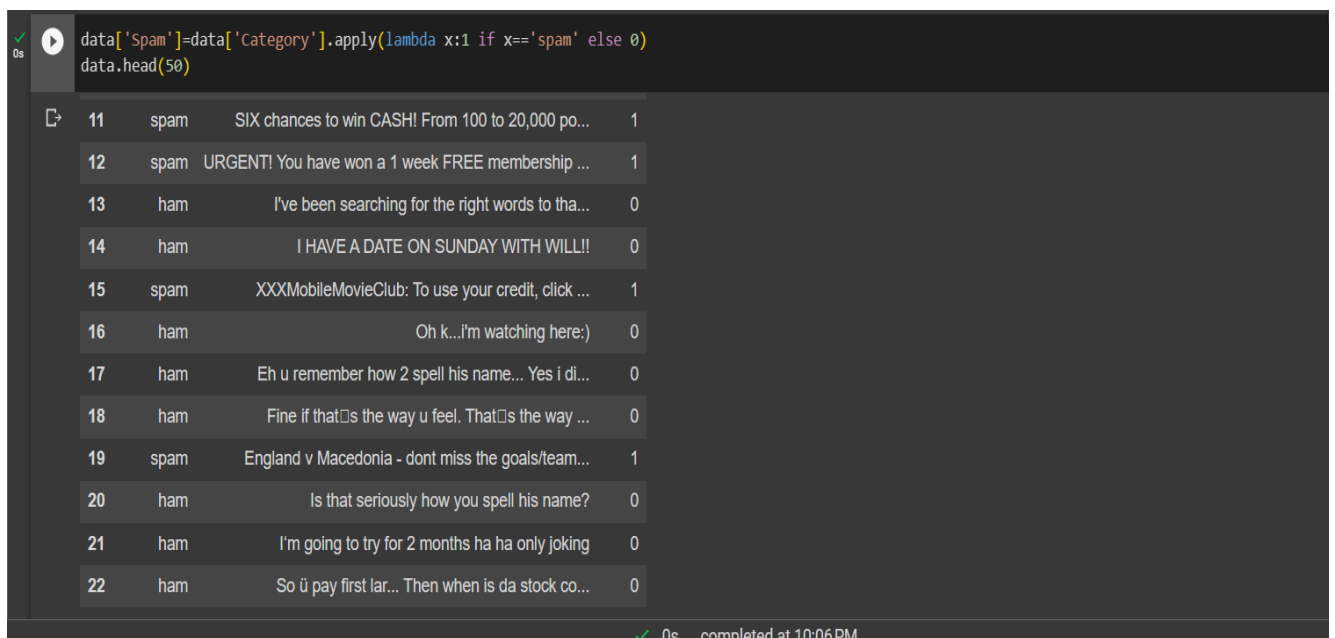## COLLECTION OF DATA

```
data['Spam']=data['Category'].apply(lambda x:1 if x=='spam' else 0)
data.head(50)
```

| 11 | spam | SIX chances to win CASH! From 100 to 20,000 po... | 1 |
| 12 | spam | URGENT! You have won a 1 week FREE membership ... | 1 |
| 13 | ham | I've been searching for the right words to tha... | 0 |
| 14 | ham | I HAVE A DATE ON SUNDAY WITH WILL!! | 0 |
| 15 | spam | XXXMobileMovieClub: To use your credit, click ... | 1 |
| 16 | ham | Oh k...i'm watching here:) | 0 |
| 17 | ham | Eh u remember how 2 spell his name... Yes i di... | 0 |
| 18 | ham | Fine if that□s the way u feel. That□s the way ... | 0 |
| 19 | spam | England v Macedonia - dont miss the goals/team... | 1 |
| 20 | ham | Is that seriously how you spell his name? | 0 |
| 21 | ham | I'm going to try for 2 months ha ha only joking | 0 |
| 22 | ham | So ü pay first lar... Then when is da stock co... | 0 |

# SCREENSHOT-4

```
data['Spam']=data['Category'].apply(lambda x:1 if x=='spam' else 0)
data.head(50)
```

| 11 | spam | SIX chances to win CASH! From 100 to 20,000 po... | 1 |
| 12 | spam | URGENT! You have won a 1 week FREE membership ... | 1 |
| 13 | ham | I've been searching for the right words to tha... | 0 |
| 14 | ham | I HAVE A DATE ON SUNDAY WITH WILL!! | 0 |
| 15 | spam | XXXMobileMovieClub: To use your credit, click ... | 1 |
| 16 | ham | Oh k...i'm watching here:) | 0 |
| 17 | ham | Eh u remember how 2 spell his name... Yes i di... | 0 |
| 18 | ham | Fine if that□s the way u feel. That□s the way ... | 0 |
| 19 | spam | England v Macedonia - dont miss the goals/team... | 1 |
| 20 | ham | Is that seriously how you spell his name? | 0 |
| 21 | ham | I'm going to try for 2 months ha ha only joking | 0 |
| 22 | ham | So ü pay first lar... Then when is da stock co... | 0 |

✓ 0s completed at 10:06 PM

# SCREENSHOT-5

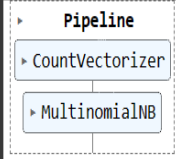| | 23 | ham | Aft i finish my lunch then i go str down lor. ... | 0 |
|---|---|---|---|---|
| | 24 | ham | Ffffffffff. Alright no way I can meet up with ... | 0 |
| | 25 | ham | Just forced myself to eat a slice. I'm really ... | 0 |
| | 26 | ham | Lol your always so convincing. | 0 |
| | 27 | ham | Did you catch the bus ? Are you frying an egg ... | 0 |
| | 28 | ham | I'm back &amp; we're packing the car now, I'll... | 0 |
| | 29 | ham | Ahhh. Work. I vaguely remember that! What does... | 0 |
| | 30 | ham | Wait that's still not all that clear, were you... | 0 |
| | 31 | ham | Yeah he got in at 2 and was v apologetic. n ha... | 0 |
| | 32 | ham | K tell me anything about you. | 0 |
| | 33 | ham | For fear of fainting with the of all that hous... | 0 |
| | 34 | spam | Thanks for your subscription to Ringtone UK yo... | 1 |
| | 35 | ham | Yup... Ok i go home look at the timings then i... | 0 |
| | 36 | ham | Oops, I'll let you know when my roommate's done | 0 |
| | 37 | ham | I see the letter B on my car | 0 |
| | 38 | ham | Anything lor... U decide... | 0 |
| | 39 | ham | Hello! How's you and how did saturday go? I wa... | 0 |
| | 40 | ham | Pls go ahead with watts. I just wanted to be s... | 0 |
| | 41 | ham | Did I forget to tell you ? I want you , I need... | 0 |
| | 42 | spam | 07732584351 - Rodger Burns - MSG = We tried to... | 1 |

✓ 0s    completed at 10:06 PM

# SCREENSHOT-6

| | 43 | ham | WHO ARE YOU SEEING? | 0 |
|---|---|---|---|---|
| | 44 | ham | Great! I hope you like your man well endowed. ... | 0 |
| | 45 | ham | No calls..messages..missed calls | 0 |
| | 46 | ham | Didn't you get hep b immunisation in nigeria. | 0 |
| | 47 | ham | Fair enough, anything going on? | 0 |
| | 48 | ham | Yeah hopefully, if tyler can't do it I could m... | 0 |
| | 49 | ham | U don't know how stubborn I am. I didn't even ... | 0 |

# SCREENSHOT-7

# TRAINING THE MODEL

```
[12] from sklearn.pipeline import Pipeline
     clf=Pipeline([
         ('vectorizer',CountVectorizer()),
         ('nb',MultinomialNB())
     ])
```

```
[13] clf.fit(X_train,y_train)
```

```
         Pipeline
    ▸ CountVectorizer

       ▸ MultinomialNB
```

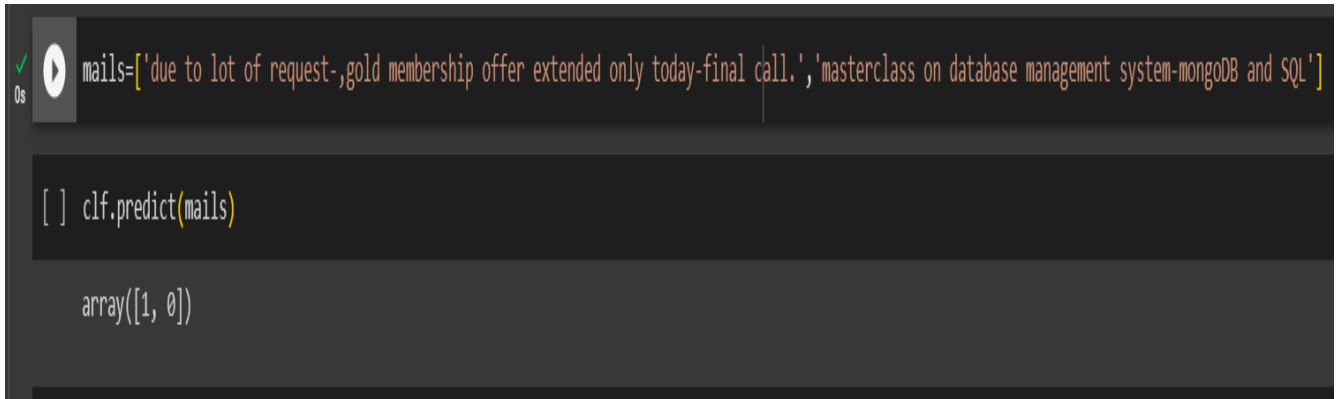# SCREENSHOT-8

# PREDICTION OF THE MODEL

```
clf.score(X_test,y_test)
```

```
0.9820531227566404
```
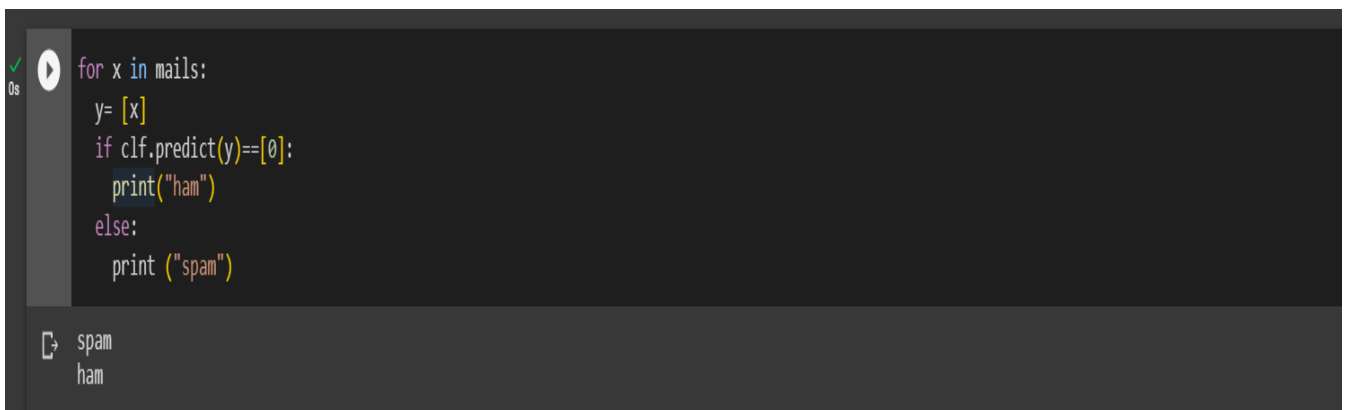
# SCREENSHOT-9

## PREDICT THE EMAILS

```
mails=['due to lot of request-,gold membership offer extended only today-final call.','masterclass on database management system-mongoDB and SQL']
```

```
[ ] clf.predict(mails)
```

```
array([1, 0])
```

# SCREENSHOT-10

## FINAL OUTPUT

```
for x in mails:
    y= [x]
    if clf.predict(y)==[0]:
      print("ham")
    else:
      print ("spam")
```

```
spam
ham
```

# CHAPTER 7

# CONCLUSION AND FUTURE WORK

## 7.1 CONCLUSION

In conclusion, Naive Bayes is a powerful and effective algorithm for email spam detection. Its accuracy, efficiency, and scalability make it an ideal choice for addressing the challenges posed by spam emails in today's digital landscape. By leveraging probabilistic calculations and assuming independence between features, Naive Bayes can accurately classify incoming emails as spam or non-spam. Its simplicity and speed enable it to process large volumes of data in real-time, making it suitable for handling the ever-increasing influx of emails. Implementing Naive Bayes for email spam detection not only helps protect users from the nuisances of spam emails but also safeguards against potential security threats such as phishing attempts and malware distribution. By filtering out unwanted and potentially harmful content, Naive Bayes plays a crucial role in maintaining the integrity and security of communication channels.

As the sophistication and prevalence of spam emails continue to grow, Naive Bayes remains a valuable tool in the fight against this issue. Its effectiveness, coupled with its ability to adapt and handle new types of spam, makes it a reliable solution for email providers, businesses, and individual users alike. In conclusion, Naive Bayes offers an efficient, accurate, and scalable approach to email spam detection, providing peace of mind and enhancing productivity in the digital realm.

## 7.2 FUTURE WORK

Future Work for Email Spam Detection using Naive Bayes

I hope this email finds you well. I am writing to discuss potential areas of future work for email spam detection using Naive Bayes. While Naive Bayes has proven to be an effective algorithm for this task, ongoing research and development efforts can further enhance its performance and address emerging challenges. Here are some avenues for future exploration:

**1. Feature Engineering**: The effectiveness of Naive Bayes heavily relies on the selection and quality of features used for classification. Future work can focus on refining the feature set by considering additional attributes such as email headers, metadata, and embedded links. Incorporating more advanced natural language processing techniques, such as semantic analysis or contextual understanding, can also improve the algorithm's ability to identify spam emails accurately.

**2. Handling Imbalanced Data**: Imbalanced datasets, where the number of spam emails is significantly lower than non-spam emails, can pose a challenge for Naive Bayes. Future research can explore techniques to address this issue, such as oversampling the minority class, under sampling the majority class, or using synthetic data generation methods. Additionally, employing ensemble methods that combine multiple Naive Bayes classifiers or other algorithms can help mitigate the impact of imbalanced data.

**3. Dynamic Adaptation**: Spam email tactics are continually evolving, requiring spam detection systems to adapt quickly. Future work can focus on developing adaptive Naive Bayes models that can automatically update their knowledge base to incorporate emerging spam patterns and tactics. This can be achieved through continuous learning and integration with real-time feedback mechanisms, allowing the system to adapt and improve its accuracy over time.

**4. Hybrid Approaches**: Combining Naive Bayes with other machine learning algorithms or techniques can potentially yield even better results. Hybrid models that leverage the strengths of Naive Bayes along with other classifiers, such as support vector machines (SVM) or deep learning models, can improve the overall performance of email spam detection systems. Exploring the benefits of such hybrid approaches and finding the optimal integration strategies is an area of future research.

**5. User Feedback Integration:** Incorporating user feedback in the spam detection process can enhance the accuracy and customization of the system. Future work can explore techniques to collect and utilize user feedback to improve the classification results. This can include allowing users to mark emails as spam or ham and incorporating their feedback into the Naive Bayes model to refine its predictions.

These are just a few directions for future work in email spam detection using Naive Bayes. As the spam landscape continues to evolve, ongoing research and innovation are essential to stay ahead of new spamming techniques and provide robust and reliable spam detection solutions.

# CHAPTER-8

# REFERENCES

[1]. Nikhil Govil, Kunal Agarwal, Ashi Bansal, Astha Varshney "A Machine Learning based Spam Detection Mechanism" Proceedings of the Fourth International Conference on Computing Methodologies and Communication (ICCMC 2020) IEEE

[2]. Nikhil kumar, sanket sonowal, nishant "email spam detection using machine learning algorithms" Proceedings of the Second International Conference on Inventive Research in Computing Applications (ICIRCA-2020) IEEE

[3]. Simran Gibson, Biju Issac, Li Zhang, Seibu Mary Jacob "Detecting Spam email with machine Learning Optimized with bio-inspired metaheuristic algorithms" IEEE VOLUME 8, 2020

[4]. Ganiev Salim Karimovich, Khamidov Sherzod Jaloddin ugli, Olimov Iskandar Salimbayevich "Analysis of machine leaning method for filtering spam messages in email services " 2020 International Conference on Information Science and Communications Technologies (ICISCT) | IEEE nov 2020

[5]. Mansoor RAZA and Nathali Dilshani Jayasinghe, Muhana Magboul Ali Muslam "A Comprehensive Review on email spam classification using Machine Learning Algorithms" 2021 International Conference on Information Networking (ICOIN) IEEE 02 February 2021

[6]. Nandhini.S, DR.Jeen Marseline.K.S "Performance Evaluation of machine algorithms for email spam detection" 2020 International Trends in Information Technology and Engineering (ic-ETITE) IEEE 2020

[7]. Mahammad Abdullahi, Abdulmalik D. Mohammed, Opeyemi O. Abisoye "A review on Machine Learning Techniques for images-based spam emails detection" Proceedings of the 2020 IEEE 2nd International Conference on Cyberspace (Cyber Nigeria) IEEE 2020

[8]. Priya.S, Annie Uthra.R "An Effective Concept Drift Detection Technique with Kernel Extreme Learning Machine for Email Spam Filtering" Proceedings of the Third  Sustainable Systems [ICISS 2020] IEEE 2020

[9]. Fahima Hossain, Mohammed Nasir Uddin, Rajib Kumar Halder "Analysis of Optimized Machine Learning and Deep Learning Techniques for Spam Detection" International IOT, Electronics and Mechatronics Conference (IEMTRONICS) IEEE 2021

[10]. Sunday Olusanya Olatunji "Extreme Learning Machines and Support Vector Machines Models for Email spam detection" Canadian Conference on Electrical and Computer Engineering (CCECE) IEEE 2017