

Cryptography and Cyber Security

Course Code: R1UC505C

Final Project Report : Spam Mail Detection

For

BACHELOR OF

ENGINEERING & TECHNOLOGY



SCHOOL OF COMPUTER SCIENCE AND ENGINEERING

GALGOTIAS UNIVERSITY, GREATER NOIDA

UTTAR PRADESH

Student Name:

ABHISHEK KUMAR

Admission No:

22SCSE1012568

Semester : V

Student Name : AVISHANT KUMAR

Admission NO. : 22SCSE1012569

Semester : V

Final Report: Spam Mail Detection

1. Introduction

Spam emails are a significant problem, causing security risks and cluttering email inboxes. This project focuses on developing a **Spam Mail Detection System** using **machine learning** techniques to classify emails as **spam** or **ham** (normal). The goal is to reduce the impact of phishing and junk emails on users by automating the filtering process.

2. Objective

The main objective of the project is to build an effective and efficient system that can detect spam emails with high accuracy, improving **email security** for users. The system will use **Naive Bayes**, a popular algorithm for text classification, and various techniques like **TF-IDF** (Term Frequency-Inverse Document Frequency) for feature extraction.

3. Methodology

- **Data Collection:** A publicly available dataset of labeled emails (spam and ham) is used. The dataset is preprocessed, including cleaning and tokenization of email text.
- **Feature Extraction:** **TF-IDF** is used to convert the text into numerical features, capturing important words and phrases in the emails.
- **Model Selection:** The **Naive Bayes** classifier is chosen for its simplicity and good performance in text classification tasks. The model is trained using the processed text data.
- **Training and Testing:** The dataset is split into training and testing sets. The model is trained on the training data and evaluated on the test set using metrics such as **accuracy, precision, recall, and F1 score**.

4. Results

After training the model, it was evaluated on the test set. The results showed that the model accurately detected **spam** emails with a high degree of accuracy. The confusion matrix and classification report indicated that the system was effective at identifying both **spam** and **ham** emails with minimal errors.

5. Conclusion

The project successfully demonstrated how **machine learning** can be used to detect **spam** and **ham** emails. By applying **Naive Bayes** and **TF-IDF**, the system achieved high accuracy in classifying emails. This solution enhances **email security** by reducing **phishing** and **junk emails**. Future work could focus on improving the model with larger datasets and more advanced techniques like **deep learning**.

6. Future Improvements

- Training on larger, more diverse datasets to improve model robustness.
 - Exploring **deep learning** methods for potentially better spam detection accuracy.
 - Implementing a **real-time email filtering system** for production environments.
-

This report summarizes the approach taken to build an **email spam detection system** using **machine learning** and lays the foundation for future improvements in **email security**.