

Report On

Spam message Classifier using Ensemble Learning(XGBoost)

Submitted in partial fulfillment of the requirements of the Course project in
Semester VII of fourth year Computer Science Engineering (Data Science)

by

Dikshant Buwa (04)
Mayannk Jadhav(20)
Yash Sankhe (53)
Arpit Sutariya(59)

Supervisor

Dr. Tatwadarshi P.N.



University of Mumbai

Vidyavardhini's College of Engineering & Technology

Department of Computer Science Engineering (Data Science)



(2023-24)

**Vidyavardhini's College of Engineering & Technology Department of
Computer Science Engineering (Data Science)**

CERTIFICATE

This is to certify that the project entitled “Spam message Classifier using Ensemble Learning(XGBoost)” is a bonafide work of “Dikshant Buwa (Roll No. 04), Mayank Jadhav (Roll No. 20), Yash Sankhe (Roll No. 53) Arpit Sutariya (Roll No. 59)” submitted to the University of Mumbai in partial fulfillment of the requirement for the Course project in Semester VII of fourth year Computer Science Engineering (Data Science).

Supervisor

Dr. Tatwadarshi P. N.

Dr. Vikas Gupta
Head of Department

Table of Contents

Chapter No		Title	Page No.
1		Introduction	
	1.1	Introduction	1
	1.2	Problem Statement	1
	1.3	Objective	1
2		Proposed System	
	2.1	Introduction	2
	2.2	Architecture/Framework	2
	2.3	Algorithm and Process Design	2
	2.4	Details of Hardware and Software	2
	2.5	Experiments and Results	3
	2.6	Conclusion	6
		References	6

Chapter 1: Introduction

1.1 Introduction

Spam is still a major problem in today's digital world, interfering with communication and posing security risks. From rule-based to machine learning-based methods, spam message identification has progressed, with ensemble learning being essential to increasing resilience and accuracy. This paper investigates the use of XGBoost, a potent ensemble learning algorithm, to improve spam message detection efficiency. We explore the difficulties presented by spam messages, the development of anti-spam strategies, and the special benefits of using XGBoost as a fundamental part of our ensemble learning methodology. Our goal in doing this study is to offer a thorough grasp of XGBoost's potential for preventing spam communications.

1.2 Problem Statement & Objectives

Unwanted spam messages are a constant annoyance in the digital era. Not only can these communications overflow our text and email inboxes, but they may also include unsolicited adverts, malware, or frauds. To safeguard our time and security, we require an accurate and fast method for automatically identifying and removing these spam messages. In order to do this, we want to make use of XGBoost, a powerful algorithm that recognizes patterns in data and helps us differentiate between spam and real messages more accurately. Our objective is to reduce the intrusion of unsolicited messages in order to improve our digital communication experience.

Chapter 2: Proposed System

2.1 Introduction

The system introduces a novel approach for spam message detection, capitalizing on the robust XGBoost algorithm to enhance accuracy and reliability. It encompasses data preprocessing, model training, and cross-validation for improved performance. Evaluation metrics, including accuracy, F1 Score, and AUC-ROC, gauge the model's effectiveness in distinguishing spam from legitimate messages. Future optimization steps include hyperparameter tuning, feature engineering, and advanced data preprocessing techniques. This system presents a comprehensive solution to combat spam in digital communication, elevating security and user experience. Dataset is of email.csv and collected from Kaggle of 32 mb.

2.3 Algorithm and Process Design

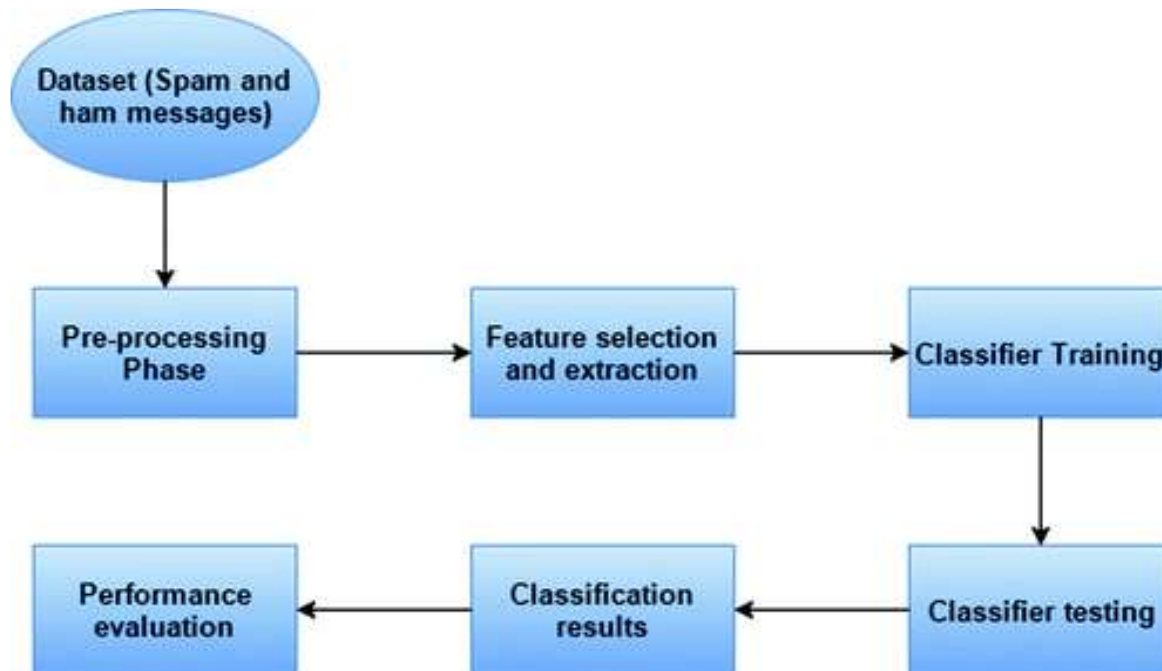


Fig 1: Process Diagram

2.4. Details of Hardware & Software

- Python
- Google Colab
- Sci-kit
- Pandas
- Numpy

- Matplotlib
- SeaBorn

2.5. Experiment and Results for Validation and Verification :

Algorithm achieves following evaluations :

Accuracy on Test Set: 0.9381642512077295

Mean AUC-ROC Score: 0.9635852704792015

Mean F1 Score: 0.9389262552387342

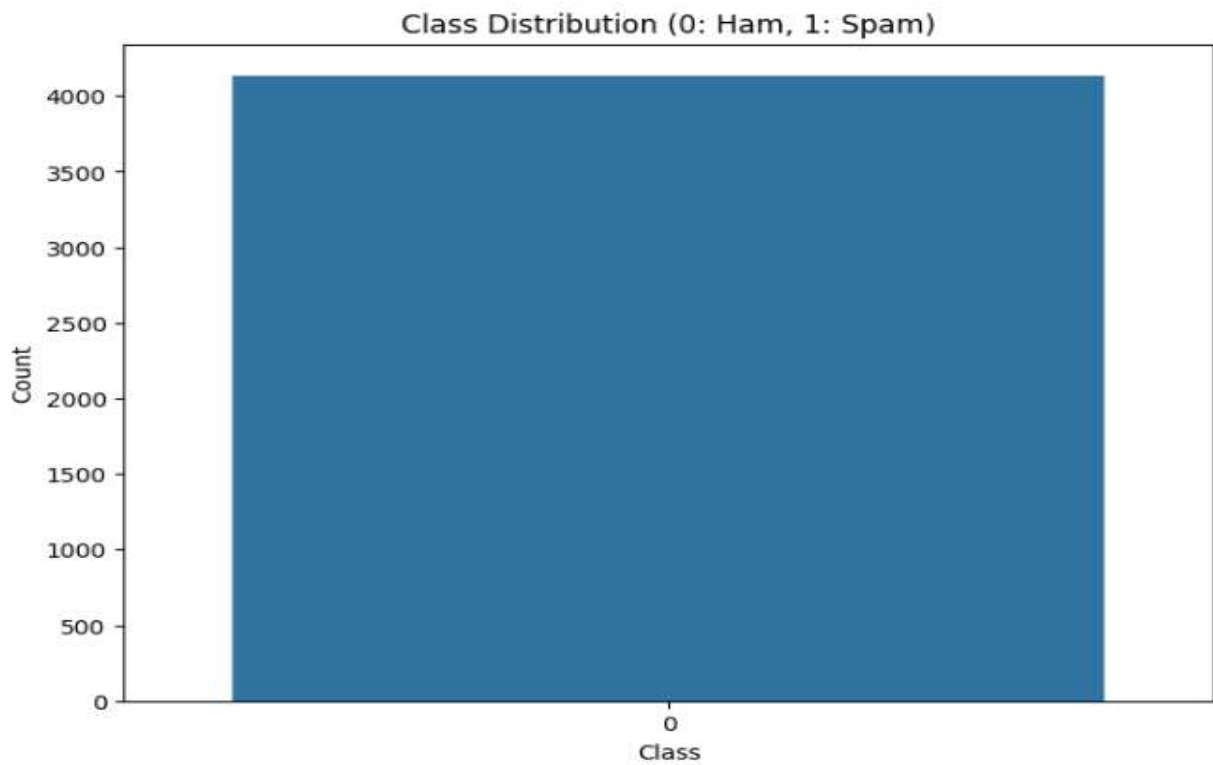


Fig 2: Visualize Class Distribution

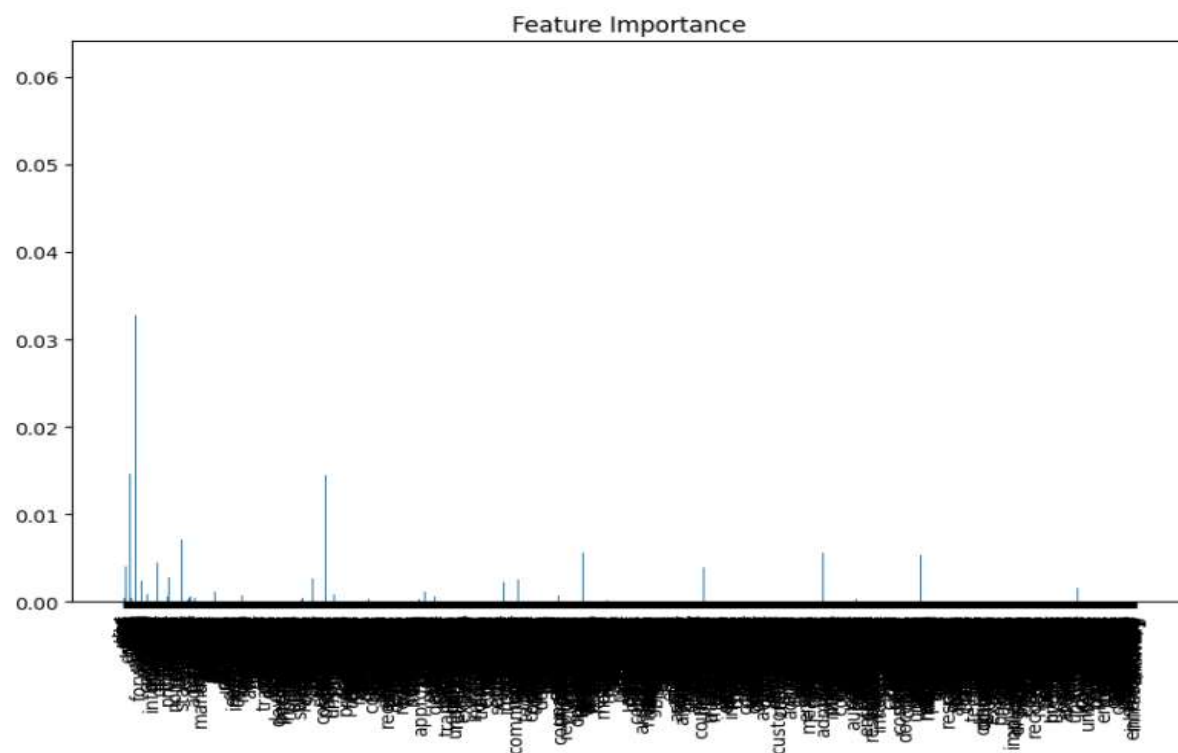


Fig 3 : Feature Importance

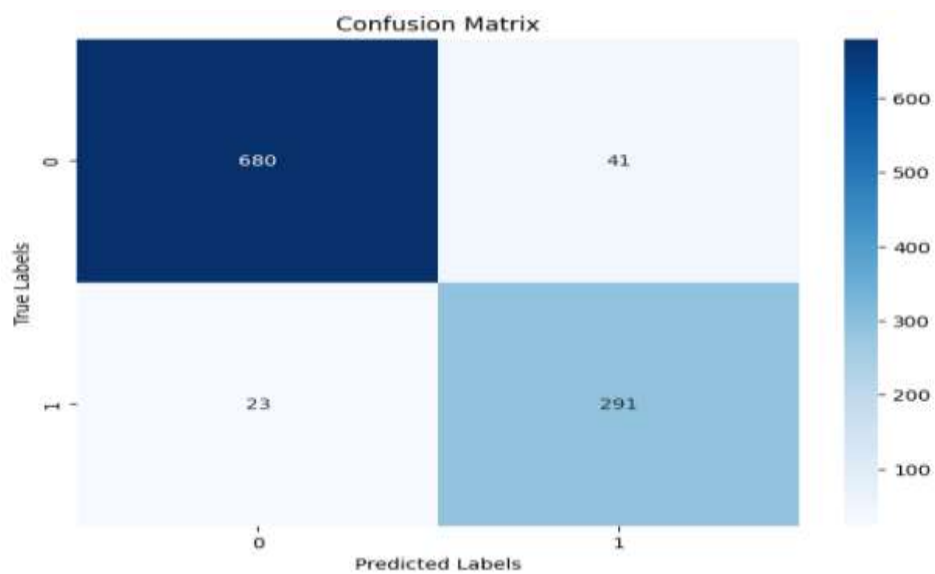


Fig 4: Confusion Matrix

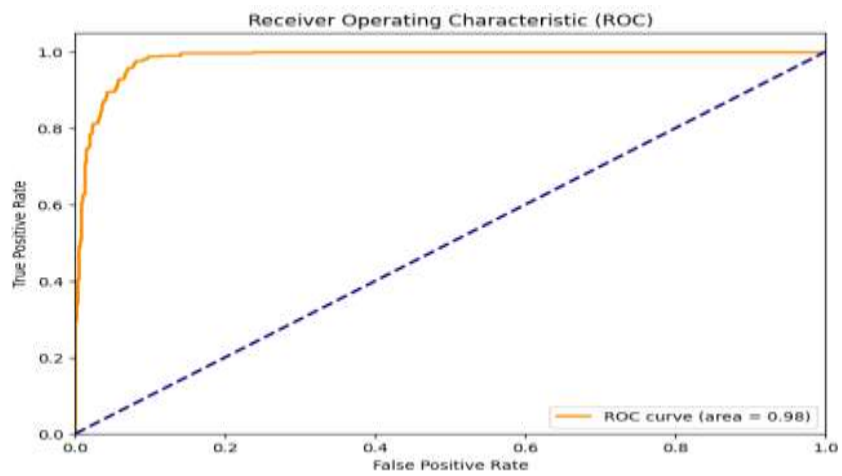


Fig 5 : ROC

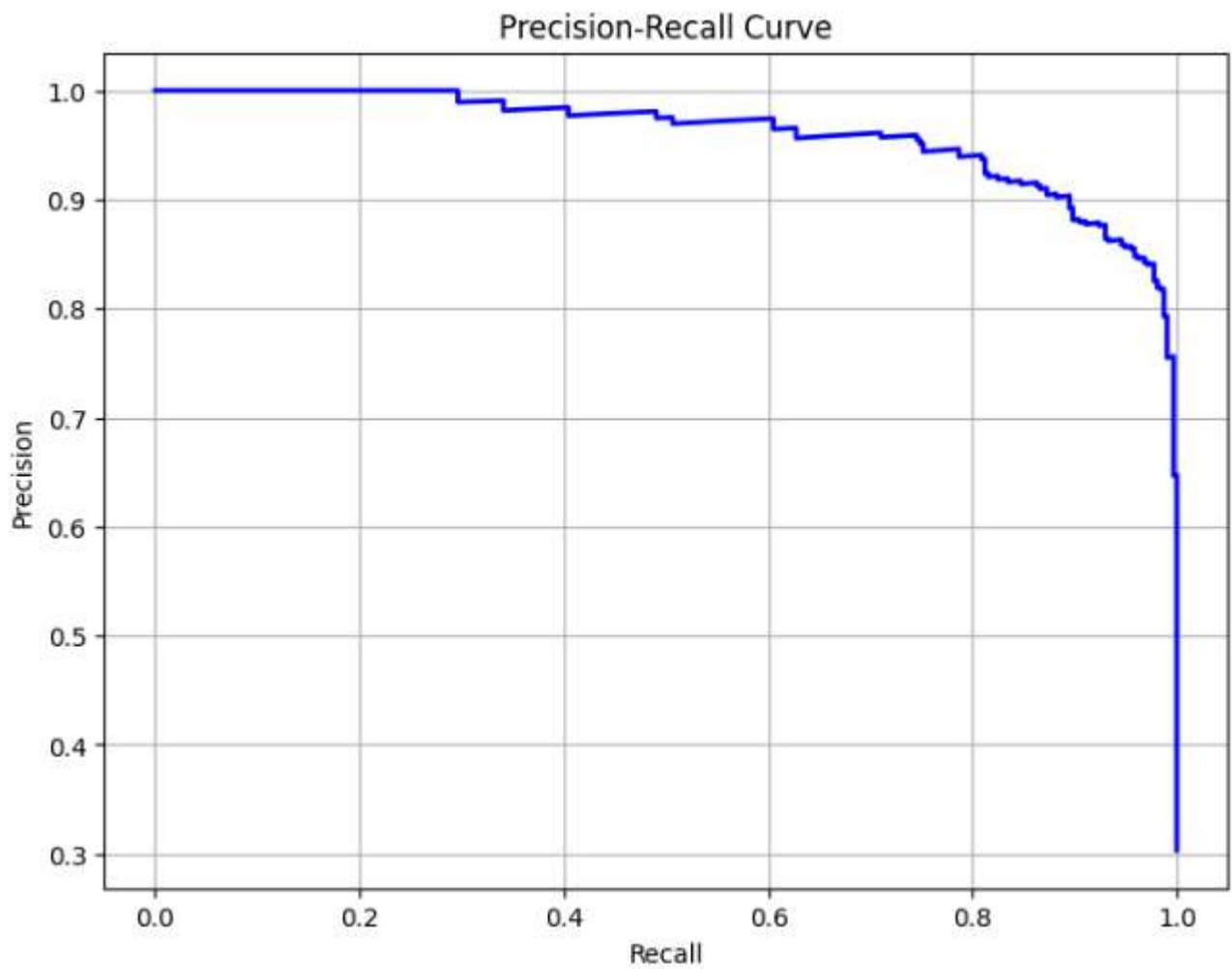


Fig 6 : F1- Score

2.6 Conclusion

When tested on a test set, the suggested XGBoost-based spam message detection system obtains a high degree of accuracy (Accuracy on Test Set: 0.938), demonstrating its capacity for accurate prediction. Additionally, the model's Mean AUC-ROC Score of 0.964 shows that it can efficiently differentiate between spam and non-spam communications. A balanced trade-off between recall and precision is shown by the Mean F1 Score of 0.939. Visualizations that offer information on feature relevance, class distribution, and the dynamic performance trends of AUC-ROC and F1 Score during cross-validation are also beneficial to the system. These visual aids improve our comprehension of the behaviour of the system and direct future improvements for more effective spam identification.

References:

1. Sridevi Gadde,” SMS Spam Detection using Machine Learning and Deep Learning Techniques”2021 7th International Conference on Advanced Computing & Communication Systems (ICACCS).
2. Suparna Das Gupta,” SMS Spam Detection Using Machine Learning” 2021 *J. Phys.: Conf. Ser.* 1797 012017.
3. M.Rubin Julis, S.Alagesan ,” Spam Detection In Sms Using Machine Learning Through Text Mining ” INTERNATIONAL JOURNAL OF SCIENTIFIC & TECHNOLOGY RESEARCH VOLUME 9, ISSUE 02, FEBRUARY 2020 ISSN 2277-8616
4. Sifat Ahmed, Faisal Muhammad,” Using Boosting Approaches to Detect Spam Reviews” 2019 1st International Conference on Advances in Science, Engineering and Robotics Technology (ICASERT)
5. Shaghayegh Hosseinpour, Hadi Shakibian, “An Ensemble Learning Approach for SMS Spam Detection” 2023 9th International Conference on Web Research (ICWR)