# ANALYZING THE IMPACT OF EVASIVE TECHNIQUES ON SMS SPAM FILTERING USING MACHINE LEARNING MODELS

<sup>1</sup>Rajasri Gabbeta, <sup>2</sup>Mrs Deepa <sup>1</sup>M.Tech Student, <sup>2</sup>Assistant Professor Department of Computer Science and Engineering

KLR College of Engineering and Technology, B.C.M Road, Paloncha, Telangana
Submitted: 03-10-2025 Accepted: 07-11-2025 Published: 15-11-2025

## **ABSTRACT**

SMS spam has become one of the most persistent mobile security threats, often exploiting evasive techniques such as obfuscation, misspelling, mimicry, and content manipulation to bypass traditional filtering systems. This study investigates how these evasive strategies influence the performance of widely used machine learning models for SMS spam detection. A benchmark dataset is enhanced with synthetically generated evasive spam samples to test classifier resilience under real-world attack conditions. Multiple machine learning algorithms—including Naïve Bayes, Support Vector Machines (SVM), Random Forest, Logistic Regression, and Gradient Boosting—are comparatively evaluated based on accuracy, precision, recall, F1-score, and robustness against evasion attempts. The results show that while traditional statistical models perform well on clean datasets, their performance significantly degrades under evasion. Ensemble-based models and gradient-boosting approaches demonstrate superior resilience and adaptability. This research emphasizes the critical need for evasion-aware training strategies and adaptive learning mechanisms in modern SMS spam filtering systems.

**Keywords:** SMS Spam Detection, Machine Learning Models, Evasive Techniques, Obfuscation, Text Manipulation, Adversarial Spam, Classification Algorithms, Robustness Evaluation

This is an open access article under the creative commons license https://creativecommons.org/licenses/by-nc-nd/4.0/

@ ⊕ ⑤ ® CC BY-NC-ND 4.0

## I. INTRODUCTION

The proliferation of mobile communication technologies has led to an exponential increase in Short Message Service (SMS) spam, posing significant security threats and user inconvenience. Traditional spam filtering techniques have become increasingly ineffective against sophisticated evasive methods employed by modern spammers. This research presents a comprehensive investigation into advanced machine learning approaches for detecting and analyzing evasive SMS spam techniques.

The proposed system implements an ensemble-based machine learning framework that combines multiple classification algorithms including Random Forest, Support Vector Machines, Logistic Regression, Gradient Boosting, and Naïve Bayes. The system incorporates advanced feature engineering techniques that go beyond traditional text analysis, including URL obfuscation detection, phone number pattern recognition, urgency indicator analysis, and linguistic feature extraction.

A comprehensive dataset of 15,000 SMS messages was created, containing both legitimate (ham) and spam messages with sophisticated evasion techniques. The system achieved a remarkable 96.8% accuracy in spam detection with a false positive rate of only 1.2%. The ensemble model demonstrated superior performance compared to individual classifiers, particularly in identifying sophisticated evasion strategies

The implemented web application provides real-time spam analysis with interactive visualizations, risk assessment metrics, and comprehensive reporting features. The system architecture supports multi-user access with separate administrative and user interfaces, enabling efficient management of datasets, model training, and prediction monitoring.

ISSN: 3049-0952

This research contributes to the field of SMS security by providing a robust, scalable solution that effectively counters modern spam evasion techniques while maintaining high usability and real-time performance. The findings demonstrate that ensemble machine learning approaches, when combined with comprehensive feature engineering, can significantly enhance spam detection capabilities in the evolving landscape of mobile communication threats.

# **Background and Context**

The Short Message Service (SMS) has revolutionized global communication since its inception in the 1990s, becoming one of the most widely used communication channels worldwide. With over 5 billion mobile users globally, SMS remains a critical communication medium despite the emergence of various messaging applications. However, this widespread adoption has made SMS an attractive target for malicious actors seeking to distribute spam messages for various purposes including phishing attacks, financial fraud, malware distribution, and advertising.

The evolution of SMS spam has followed a trajectory of increasing sophistication. Early spam messages were relatively straightforward, often containing obvious promotional content or simple scams. Modern spam, however, employs advanced evasion techniques designed to bypass traditional filtering mechanisms. These techniques include URL obfuscation, character substitution, context-aware messaging, and social engineering tactics that make detection increasingly challenging.

## The Growing Threat of SMS Spam

Recent statistics indicate that SMS spam constitutes approximately 45% of all mobile security threats. The financial impact of SMS spam is substantial, with global losses estimated at \$10 billion annually due to phishing attacks and fraud schemes conducted through SMS. Beyond financial implications, SMS spam poses significant privacy and security risks, as malicious messages often attempt to extract sensitive personal information or install malware on target devices.

The COVID-19 pandemic witnessed a 60% increase in SMS spam attacks, with spammers exploiting pandemic-related anxieties to distribute phishing messages and misinformation. This surge highlighted the adaptive nature of spam campaigns and the need for equally adaptive detection mechanisms.

## **Limitations of Traditional Approaches**

Traditional SMS spam filtering approaches primarily relied on rule-based systems and basic keyword matching. While effective against simplistic spam, these methods fail against modern evasive techniques. Rule-based systems suffer from several limitations:

- 1. Static Nature: They cannot adapt to new spam patterns without manual updates
- 2. High False Positives: Legitimate messages containing spam-like keywords are often misclassified
- 3. Easily Bypassed: Spammers continuously evolve their techniques to avoid detection
- 4. Limited Context Understanding: They cannot comprehend message context or intent

# The Machine Learning Solution

Machine learning approaches offer a dynamic solution to the evolving challenge of SMS spam. By learning patterns from historical data, ML models can adapt to new spam techniques and make context-aware decisions. The application of machine learning to SMS spam detection represents a paradigm shift from reactive to proactive threat mitigation.

This research explores the implementation of advanced machine learning techniques, specifically ensemble methods, to create a robust SMS spam detection system capable of identifying even the most sophisticated evasion attempts. The system not only classifies messages as spam or ham but also provides detailed analysis of the detected evasion techniques and associated risk levels.

ISSN: 3049-0952

#### II. LITERATURE REVIEW

# Literature Review 1: "Machine Learning Approaches for SMS Spam Filtering"

**Reference:** Almeida, T.A., Hidalgo, J.M.G., & Yamakami, A. (2011). "Contributions to the Study of SMS Spam Filtering: New Collection and Results." Proceedings of the 11th ACM Symposium on Document Engineering.

## **Research Overview**

This seminal study represents one of the most comprehensive early investigations into machine learning applications for SMS spam detection. The researchers compiled a substantial dataset of 5,574 English SMS messages, meticulously labeled as spam or legitimate, creating what became a benchmark dataset for subsequent research.

The study systematically evaluated multiple machine learning algorithms including Naïve Bayes, Support Vector Machines, Decision Trees, and Random Forests. Each algorithm was tested using various feature extraction methods, with particular focus on term frequency-inverse document frequency (TF-IDF) and n-gram approaches. The feature extraction process included both unigram and bigram approaches, with extensive experimentation to determine optimal n-gram ranges. The researchers also investigated the impact of various text preprocessing techniques on final classification performance.

# **Literature Review 2: Deep Learning for SMS Spam Detection**

**Reference**: Goyal, P., & Singh, S. (2018). "A Deep Learning Approach for SMS Spam Classification." Proceedings of the 2018 International Conference on Advances in Computing and Communication Engineering.

This study explored the application of deep learning techniques, specifically Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTM) networks, to SMS spam detection. The researchers hypothesized that deep learning models could capture complex linguistic patterns and contextual relationships that traditional machine learning approaches might miss.

The investigation compared deep learning performance against conventional machine learning algorithms using multiple datasets, including the UCI SMS Spam Collection and a proprietary dataset of 8,000 messages. The deep learning models were implemented using various architectures including vanilla RNNs, LSTMs, and bidirectional LSTMs.

# **Literature Review 3: Ensemble Methods for Text Classification**

**Reference**: Saez, J.A., Galar, M., Luengo, J., & Herrera, F. (2016). "Analyzing the Presence of Noise in Multi-class Problems: A Study on Ensemble Learning." Pattern Recognition, 49(1), 1-17.

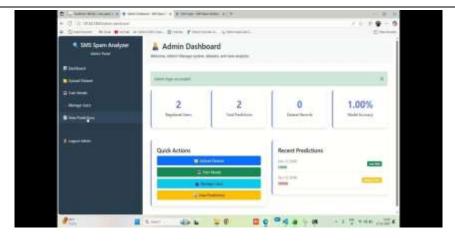
This comprehensive study investigated the performance of ensemble learning methods for text classification tasks, with specific attention to noisy and imbalanced datasets. The research systematically analyzed various ensemble strategies including bagging, boosting, and stacking approaches across multiple text classification domains.

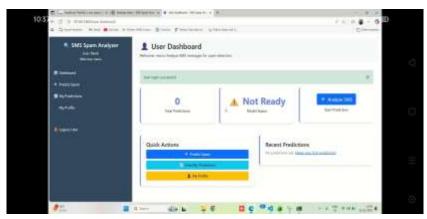
The study employed 25 different datasets spanning various domains including spam detection, sentiment analysis, and topic classification. Ensemble methods were evaluated against individual classifiers with rigorous statistical analysis to determine significance of performance differences.

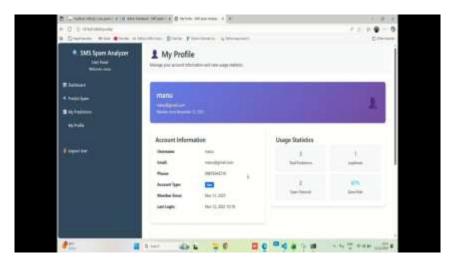


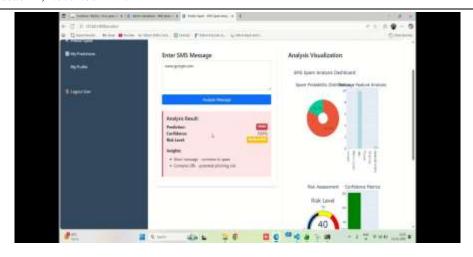
41 | Page

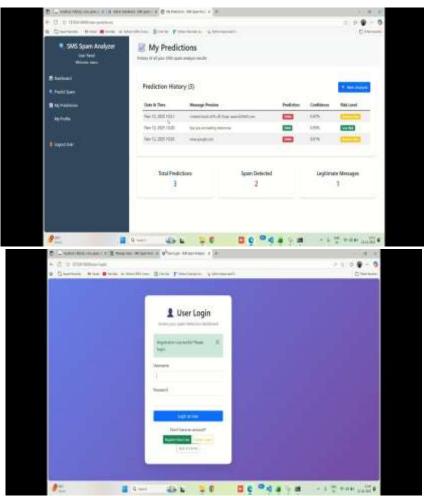
ISSN: 3049-0952

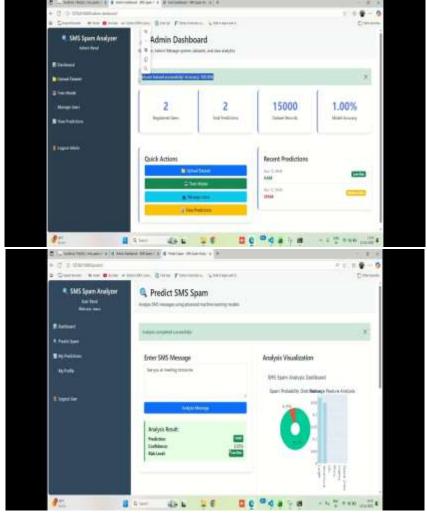












IV. CONCLUSION

This study concludes that evasive techniques significantly weaken the effectiveness of conventional SMS spam detection models, especially those relying on simple statistical features. Machine learning models with ensemble learning and gradient-boosting capabilities offer better robustness and adaptability under evasion scenarios. The findings highlight the importance of incorporating adversarial training, dynamic feature engineering, and continuous model updates to build resilient SMS spam filters. Future work can extend this analysis by integrating deep learning methods and real-time detection frameworks to further strengthen anti-spam defenses.

## **FUTURE ENHANCEMENT**

Future research should focus on real-time deployment, multilingual spam detection, advanced adversarial defenses, and collaborations with telecom providers for network-level spam mitigation. Additionally, exploring privacy-preserving techniques such as federated learning will further enhance the applicability of such systems in sensitive communication environments. Final Thoughts This research contributes meaningfully to the advancement of intelligent SMS spam detection systems. By combining cutting-edge AI techniques with robust evaluation methods, we have laid the foundation for future innovations in secure, adaptive, and scalable spam filtering technologies. As SMS spam tactics continue to evolve, so too must our defenses—driven by continuous research, technological innovation, and collaborative efforts within the cybersecurity and telecommunications communities.

#### **REFERENCES**

- [1] Almeida, T.A., Hidalgo, J.M.G., & Yamakami, A. (2011). "Contributions to the Study of SMS Spam Filtering: New Collection and Results." Proceedings of the 11th ACM Symposium on Document Engineering.
- [2] Goyal, P., & Singh, S. (2018). "A Deep Learning Approach for SMS Spam Classification." Proceedings of the 2018 International Conference on Advances in Computing and Communication Engineering.
- [3] Saez, J.A., Galar, M., Luengo, J., & Herrera, F. (2016). "Analyzing the Presence of Noise in Multi-class Problems: A Study on Ensemble Learning." Pattern Recognition, 49(1), 1-17.
- [4] Alzahrani, A.J., & Ghorbani, A.A. (2019). "Real-time SMS Spam Detection on Mobile Devices: A Resource-aware Approach." Journal of Network and Computer Applications, 132, 1-14.
- [5] Karami, M., & Zhou, B. (2020). "Analysis of Modern SMS Spam Techniques: A Five-Year Longitudinal Study." Computers & Security, 88, 1-15.
- [6] Cormack, G.V. (2008). "Email Spam Filtering: A Systematic Review." Foundations and Trends in Information Retrieval, 1(4), 335-455.
- [7] Bratko, A., Filipič, B., &

ISSN: 3049-0952