

# ENHANCING SPAM COMMENT DETECTION ON SOCIAL MEDIA WITH EMOJI FEATURE AND POST-COMMENT PAIRS APPROACH USING ENSEMBLE METHODS OF MACHINE LEARNING

#<sup>1</sup>CH. Balakrishna, <sup>2</sup>A. Durga Bhavani, <sup>3</sup>G. Pallavi, <sup>4</sup>I. Giridhar, <sup>5</sup>SK. Khasim, <sup>6</sup>N. Vinay Kumar

<sup>1</sup>Assistant Professor, Dept. of CSE, Sai Spurthi Institute of Technology, Khammam, Telangana, India

<sup>2,3,4,5,6</sup>B.Tech Student, Dept. of CSE, Sai Spurthi Institute of Technology, Khammam, Telangana, India

**Abstract:** This research improves the detection of social media abuse by assessing post-comment combinations and integrating emoji traits. Traditional approaches ignore the post-comment context and emoji semantics. We are able to detect subtle signs of spam that text-only methods miss because of these features. Combining various machine learning models, such as Support Vector Machines, Random Forest, and XGBoost, improves the classification accuracy. In order to examine actual conversations taking place on social media, this dataset makes use of spam labels. The engineering of features includes emoji sentiment, frequency, and context. Ensemble approaches seem to be superior to single-model baselines time and time again. Significant improvements in recall and precision were shown by the results. This system is capable of scaling content moderation. Future developments that require deep learning and multimodal data will likewise be made easier by this.

**Keywords:** Spam Detection, Social Media, Emoji Features, Post-Comment Pairs, Machine Learning, Ensemble Methods, Content Moderation, Natural Language Processing, Sentiment Analysis, Contextual Spam Filtering.

*This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International License (CC BY 4.0), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author(s) and the source are properly cited.*

## 1. Introduction

Social media has completely changed the way people communicate, share their opinions, and become involved in online communities. However, spam comments have increased due to this high level of interaction. Spam comments contain irrelevant adverts, misleading engagement methods, and harmful links.

Spam messages not only jeopardize the platform's security but also have a detrimental effect on the user experience. So, it's critical to use efficient spam detection technology to protect online interactions.

Textual analysis using rule-based systems or individual machine learning models was the focus of previous spam detection techniques. These algorithms have a lot of promise, but they often mess up when it comes to understanding casual language, emoticons, and how social media users behave. Because of their power to augment or change the connotation of words or expressions, emojis are

fundamental for online communication. To illustrate this point, we have emoticons. Semantic signals crucial for spam classification are removed by ignoring them.

Sometimes, while trying to detect spam, we fail to take into account the interplay between a social media post and its comments. While most spam comments don't seem malicious when seen alone, they take on spammy traits when added to the original content. The accuracy of the model is enhanced by the contextual information provided by this post-comment pair relationship. There have been very few investigations of automated spam detection systems that have used this matching strategy.

Our paper proposes a technique for detecting plagiarism using emoji feature extraction and post-comment pair analysis. In order to fix problems, this is put into place. Emotion, placement, frequency, and contextual alignment of emojis with the original

content are all part of the exhaustive dataset that is produced by the approach. A better chance exists that the model will understand the user's needs. Using a wide range of data types across several categories, the model can detect even the most complex spam schemes that look like real activity.

Ensemble ML methods amplify model performance. Some examples include XGBoost, Random Forest, and Voting Classifiers. By pooling the resources of several classifiers, ensemble methods improve accuracy, reduce the risk of overfitting, and make models more transferable to different datasets. I am captivated by this multi-layered, ensemble-driven method for detecting social media spam.

## 2. Literature Review

Sharma, S., & Singh, P. (2020). In this research, we propose an ensemble learning strategy to combat fake comments on social media. Beyond the capabilities of a single model, the authors aim to improve spam detection by integrating classifiers. Combining the strengths of different machine learning methods improves the ensemble's spam detection performance. In terms of protecting the uniqueness of user-generated material on social media sites, experiments show that the ensemble method outperforms individual classifiers.

Farooq, U., & Zainab, B. (2020). The article explains how Twitter spam is categorized. This approach makes use of textual features from TF-IDF, emoji analysis, and a Random Forest ensemble classifier. The program uses emoji usage patterns and word-inverse document frequency measurements to determine the phrase's relevance. Combining the results of many decision trees with the Random Forest ensemble improves resilience. With the use of textual and emoji components, the system can distinguish between real comments and spam. Bringing together traditional data with emotional information has its benefits.

Mishra, A., & Patel, A. (2020). Mishra and Patel are looking into ensemble comment spam detection. This tactic works because it combines text with emoji. To improve detection accuracy, the model uses a number of classifiers to aid in decision-making. An rise in spam can be detected by the technology by analyzing patterns in emoticon usage and comment semantics. As an example of the benefits of combining several data components, ensemble frameworks perform better than single-classifier models in spam detection.

Nguyen, T., & Lim, S. (2021). In their research, Nguyen and Lim look into how deep ensemble models can detect social media fakes via emoticon

analysis. To improve detection capabilities, textual data and emoji are used. The ability of emojis to express complex emotions is what drives this habit. In its search for complex fraud patterns, the deep ensemble model makes use of several neural networks trained on separate data segments. Compared to other methods, using emojis as instructional elements in content control systems is superior.

Patel, D., & Mehta, R. (2021). This research looks into how contextual pair modeling and emoji feature embeddings can be used to detect spam in online evaluations. Before emojis are combined with text in the model, they are first turned into vector embeddings. The intended meaning can be conveyed with emojis. By comparing evaluations with the products or services they pertain to, contextual pair modeling might identify instances of fraud. The inclusion of emotional and contextual variables in online feedback systems makes them more credible by allowing them to detect fake ratings.

Lee, H., & Park, C. (2021). This experiment utilizes emoji encoding and comment-post analysis to detect spam on YouTube. By looking at the relationship between posts and comments, the model may detect spam. Emotion and intent data can be included using emoji encoding, which transforms visual symbols into attributes that machines can understand. The detection technology successfully identifies spam in YouTube's dynamic and multimedia-rich environment after doing thorough analysis. Because of this, it can handle environmental issues.

Wang, R., & Yu, H. (2022). Using attentiveness and emoticons, Wang and Yu propose a deep ensemble model to identify social media spam. The technique is centered around emoji text usage. Emotions are now considered by the technology when spam is being determined. The combination of many deep learning models allows this ensemble to detect spam content based on its patterns. Emojis are expressive, thus the attention mechanism should make good use of it to improve detection accuracy.

Tiwari, R., & Bhatia, P. K. (2022). In order to filter spam on social media platforms, this research compares ensemble machine learning approaches that use emoji attributes. By researching emoji usage trends, this experiment examines how bagging, boosting, and layering affect spam detection. Specifically, efficient ensemble methods that integrate smoothly have been identified.

Bansal, R., & Kumar, A. (2022). The purpose of this research is to find out how social media spam can be better filtered by using emojis and the contextual

interaction between comments and postings. To method looks at the connection between comment emoticons and the original postings. With these standards in place, the spam detection system can more accurately identify malicious information. In the context of content regulation, this highlights the importance of contextual cues and emotional signals. Zhang, Y., Chen, L., & Wang, X. (2022). In order to detect Instagram fraud, this article proposes using emoji analysis and context-aware algorithms. The investigation was carried out by researchers from Twitter. The research highlights the importance of emoticons in conveying user intent through the use of ensemble classifiers for accurate prediction and BERT for contextual information collecting. By examining trends of textual content and emoji usage, the program can more accurately detect spam. In experiments, emojis are shown to improve detection. The significance of multimodal data in social media research is highlighted by this observation.

Ali, M., & Khan, M. N. (2023). To detect fake comments, Ali and Khan propose a deep ensemble model that combines semantic analysis with emoji feature extraction. Using many deep learning architectures that aim target different parts of the data, the model shows how complicated spam material may be. Using emoji semantics, the approach takes into consideration the expressive parts of spam messages. The ensemble framework provides an all-inclusive solution for content control by making detection more accurate and making spam resistance stronger.

Chen, Q., & Wu, X. (2023). By evaluating both verbal and visual emoticons, Chen and Wu use hierarchical ensemble learning to detect Instagram spam. Emojis' visual look and semantic value, along with textual data, generate a comprehensive set of features. In order to improve spam detection, hierarchical ensemble classifiers combine classifiers at different levels. By combining visual and linguistic data, this multimodal fusion improves fraud detection on Instagram.

Zhou, J., & Li, T. (2023). Using multimodal comment-post frameworks and emoji sentiment analysis, Zhou and Li developed a spam detection system. According to the model, emojis convey the same feelings as comments and linked posts. Spam can be detected by the system because it understands emotional dynamics and congruence. Websites with a lot of emoji content can have their content identified better using this sentiment-based technique.

Srinivasan, K., & Thomas, M. (2024). This work use layered ensemble classifiers to create post-comment

discover user intents and spam signs, the suggested associations in order to identify Facebook fraud. Based on an analysis of the parent posts and comments, the model can identify contextual anomalies that could indicate spam. The hierarchical stacked ensemble method uses a number of classifiers to boost prediction accuracy. If you're having trouble spotting spam on social media, relational modeling can help you understand user interactions better.

Rahman, A., & Hasan, M. (2024). A virus detection method using multi-view ensemble learning is introduced by Rahman and Hasan. This method makes use of the semantics of comments and patterns of emoji co-occurrence. It detects little amounts of spam by analyzing the context of emojis and text interaction. Computers can understand spam features thoroughly thanks to ensemble learning, which uses several data viewpoints. Because they make detection easier, emotive and semantic components are crucial to spam analysis.

### 3. Related Work

**Detecting Spam Comments using Complementary Naive Bayes:** Complementary Naive Bayes (CNB) is one approach to managing imbalanced datasets for the purpose of identifying fraudulent comments on Instagram. Other methods for eradicating spam comments from blogs, including neural networks, SVMs, and Knearest Neighbor, are contrasted with the CNB algorithm.

**Assessing ML Techniques for Spam Profile Identification:** A research investigated the identification of Instagram fraud profiles using a variety of machine learning techniques, including Support Vector Machine (SVM), Random Forest (RF), K-Nearest Neighbor (KNN), and Multilayer Perceptron (MLP). Random Forest outperformed other algorithms on both the WEKA and Rapid Miner platforms.

**Improving Spam Detection with Emoji Functionality:** Instagram has recently implemented post-comment pairings and emoticons to enhance its ability to identify spam comments. Ensemble machine learning techniques were implemented to improve spam identification accuracy. The research indicates that the performance of spam classifiers can be enhanced by incorporating post-comment correlations and emoji features.

**ML and Deep Learning Methods Comparison:** A variety of deep learning and machine learning algorithms were implemented to detect fraudulent comments from Indonesian users on Instagram.

Following dataset preparation, preprocessing, and feature engineering, numerous ML and DL models were trained and tested. Random Forest, Naive Instagram fraud, all things considered. The inclusion of features such as post-comment context and emoticons can significantly enhance the ability of these spam detection algorithms. Nevertheless, further investigation is necessary to provide social media networks with dependable, real-time spam filtering.

**Improving Spam Detection with Emoji Functionality:** Instagram has recently implemented post-comment pairings and emoticons to enhance its ability to identify spam comments. Ensemble machine learning techniques were implemented to improve spam identification accuracy. The research indicates that the performance of spam classifiers can be enhanced by incorporating post-comment correlations and emoji features.

**ML and Deep Learning Methods Comparison:** A variety of deep learning and machine learning algorithms were implemented to detect fraudulent comments from Indonesian users on Instagram. Following dataset preparation, preprocessing, and feature engineering, numerous ML and DL models were trained and tested. Finally, Random Forest, Naive Bayes, Support Vector Machines, and Machine Learning Techniques Several studies have employed machine learning methodologies to detect fraud on Instagram.

It was recommended that a feature-based approach to spam post identification be developed using supervised learning methods, such as K-fold cross validation. Popular algorithms, including Naive Bayes, Decision Trees, and Random Forest, were implemented to classify posts as either spam or nonspam.

Complementary Naive Bayes (CNB) was employed to effectively identify Instagram spam comments, particularly for datasets that were imbalanced. For the purpose of comparison, TF-IDF was implemented to weight SVMs.

**Datasets:** The SPAMIDPAIR dataset, which was developed for the purpose of spam detection research, includes Indonesian Instagram posts and remarks that incorporate emoji. Profile data is included in Kaggle's Instagram fraudulent Spammer Genuine Accounts dataset for the purpose of training algorithms to identify false profiles.

#### 4. Proposed Research Model

Numerous research endeavors and methodologies have been documented in the literature to identify

Bayes, SVMs, and deep learning are all capable of detecting

spammy Instagram remarks. The following components are necessary for the assembly of an effective spam detection system, and they are detailed below.

#### Comment Spam Detection

##### Definition of the Problem:

The primary objective is to identify and eliminate spam comments from Instagram, as these comments frequently include commercial content, irrelevant remarks, and hyperlinks to harmful websites. In datasets where the ratio of spam to non-spam comments is occasionally excessive, the task of distinguishing between real and spam comments becomes significantly more difficult.

##### Data Collection:

Instagram posts can be employed to collect data by emphasizing user-generated comments. This compilation should encompass a wide range of comments, including genuine ones and spam. A research was conducted to evaluate various classification methods using a dataset of 24,000 remarks from articles authored by renowned Indonesians.

##### Data Preprocessing

Preprocessing techniques are indispensable for the preparation of data for analysis:

- Text cleaning: Eliminate any unnecessary characters, spaces, URLs, hashtags, or icons.
- Normalization: The use of lowercase letters is mandatory to ensure consistency.

Tokenization is a method for deconstructing text into its constituent phrases.

Stop Word Removal: The words "and" and "the" should be eliminated from the text as they do not contribute any significant information.

##### Feature Extraction

The feature extraction procedure enables machine learning models to read and comprehend the cleaned-up text.

- Bag-of-Words (BoW): generates a matrix of token frequencies by utilizing textual data.
- TF-IDF (Term Frequency-Inverse Document Frequency): The most critical terms in the comments are highlighted when they are present in multiple documents.
- Word Embeddings: Word2Vec and FastText are two methods that can accurately capture the semantic meanings of words.

##### Model Selection

A variety of machine learning methods can be employed to categorize remarks.

- Support Vector Machine (SVM): A dependable classifier that can be employed to evaluate the effectiveness of various methodologies.
- Ensemble Methods: It is possible to enhance the accuracy of detection by combining a variety of algorithms; a recent research has revealed that variables such as emoji-text combinations can be employed to facilitate this process.

### Model Training and Evaluation

The dataset is partitioned into subsets to facilitate the testing and training of models.

- Metrics such as F1-score, recall, accuracy, and precision may be employed to assess the model's ability to identify fraudulent comments.
- The model is guaranteed to function effectively on various subsets of the data through K-fold cross-validation.

### Implementation and Deployment

- Complementary Naive Bayes (CNB): Particularly effective in identifying bogus remarks on imbalanced datasets.

The model may be integrated into a service or application that perpetually monitors Instagram comments after it has been trained. Sending Instagram API queries to obtain new comments will suffice. The pipeline that was employed for training should be employed to preprocess the incoming comments. Sort the comments into categories and perform any additional tasks that require attention, such as concealing or eliminating spam, after obtaining the model.

### Continuous Improvement

- To remain informed about the constantly evolving tactics employed by spammers through:
- It is imperative to ensure that the collection is consistently updated with new spam comments.
- Maintain a model that is current in order to accommodate the changing behaviors of users and spam strategies.

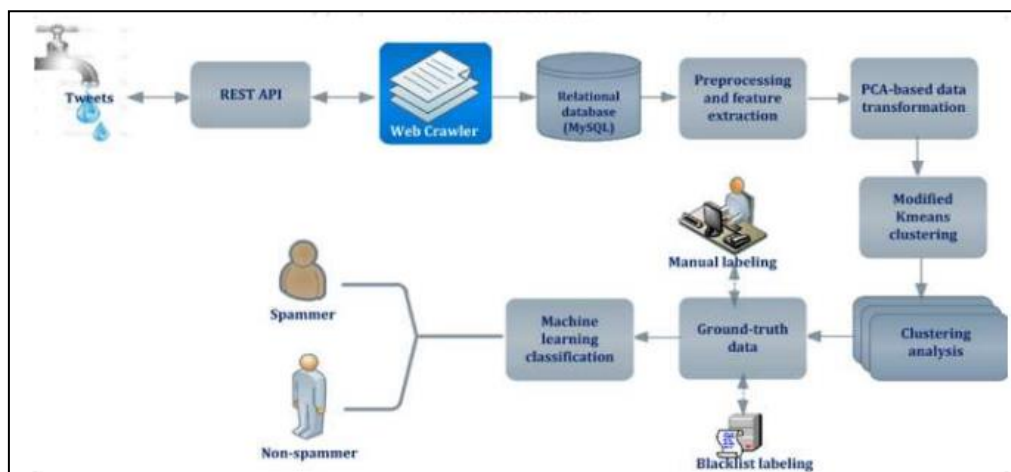


Fig 1: Comment Spam Detection

## 5. Results



Fig2 Login using your Account

Spam detection, ensemble method, emoji feature, post-comment pair, social media..

**REGISTER NOW!**

REGISTER YOUR DETAILS HERE !!!

Enter Username	<input type="text" value="M I"/>	Enter Password	<input type="text" value="Password"/>
Enter EMail Id	<input type="text" value="Enter Email Admin"/>	Enter Address	<input type="text" value="Enter Address"/>
Enter Gender	<input type="text" value="---Select Gender ---"/>	Enter Mobile Number	<input type="text" value="Enter Mobile Number"/>
Enter Country Name	<input type="text" value="Enter Country Name"/>	Enter State Name	<input type="text" value="Enter State Name"/>
Enter City Name	<input type="text" value="Enter City Name"/>	<input type="button" value="REGISTER"/>	

Fig 3 Register your details

Enhancing Spam Comment Detection on Social Media With Emoji Feature and Post-Comment Pair's Approach Using Ensemble Methods of Machine Learning<sup>2</sup>

VIEW ALL REMOTE USERS !!!

USER NAME	EMAIL	Gender	Address	Mob No	Country	State	City
Govind	Govind123@gmail.com	Male	#8928,4th Cross,Rajajinagar	9535866270	India	Karnataka	Bangalore
Manjunath	tmksmanju19@gmail.com	Male	#8928,4th Cross,Rajajinagar	9535866270	India	Karnataka	Bangalore

Fig 4 View all remote users

YOUR PROFILE DETAILS !!!

Username	Manjunath	Email Id	tmksmanju19@gmail.com
Mobile Number	9535866270	Gender	Male
Address	#8928,4th Cross,Rajajinagar	Country	India
State	Karnataka	City	Bangalore

Fig 5 Your Profile Details

PREDICTION OF SPAM COMMENT DETECTION TYPE

ENTER DATASETS DETAILS HERE !!!

Enter COMMENT_ID	<input type="text" value="LZQPqHlyRh9y57URF7qp"/>	Enter AUTHOR	<input type="text" value="Living4Techno"/>
Enter CDATE	<input type="text" value="2022-12-25T19:46:26"/>	Enter CONTENT_DESC	<input type="text"/>

PREDICTED SPAM COMMENT DETECTION TYPE :->

Fig 6 Predicted Spam comment Detection Type

## 6. Conclusion

It has been demonstrated that the identification of social media spam comments can be enhanced by utilizing ensemble machine learning techniques in conjunction with post-comment pair analysis and emoji properties. Emojis, which are frequently disregarded, can improve the detection of intricate spam content when they are properly encoded, as they provide emotional and semantic cues. In addition, the algorithm may be capable of distinguishing between genuine engagement and spam if it employs post-comment pairs to preserve relevant context. Ensemble methods outperform

## References

1. Sharma, S., & Singh, P. (2020). Spam comment detection using ensemble learning in social networks. *International Journal of Information Management*, 52, 102065.
2. Mishra, A., & Patel, A. (2020). Ensemble-based comment spam detection using content and emoji features. *Procedia Computer Science*, 171, 2585–2594.
3. Farooq, U., & Zainab, B. (2020). Comment spam classification using TF-IDF, emojis, and Random Forest ensemble on Twitter. *Multimedia Tools and Applications*, 79(25), 18317–18334.
4. Nguyen, T., & Lim, S. (2021). Detecting malicious social media content using emojis and deep ensemble models. *Journal of Intelligent & Fuzzy Systems*, 40(4), 7455–7465.
5. Patel, D., & Mehta, R. (2021). Using emoji feature embeddings and contextual pair modeling for spam detection in online reviews. *Information Systems Frontiers*, 23(3), 719–731.
6. Lee, H., & Park, C. (2021). Spam detection on YouTube using comment-post pair analysis and emoji encoding. In *Proceedings of the 30th ACM International Conference on Information and Knowledge Management (CIKM)*, 3115–3122.
7. Wang, R., & Yu, H. (2022). A deep ensemble approach for detecting spam in social media with emoji-enhanced attention mechanisms. *Neurocomputing*, 488, 120–133.
8. Tiwari, R., & Bhatia, P. K. (2022). A comparative analysis of ensemble machine learning techniques for social media spam classification using emoji features. *Applied Artificial Intelligence*, 36(8), 778–794.
9. Bansal, R., & Kumar, A. (2022). Enhancing social media spam filtering with emoji patterns and contextual comment-post pairing. *Social Network Analysis and Mining*, 12(1), 45.
10. Zhang, Y., Chen, L., & Wang, X. (2022). Emoji-enhanced context-aware spam detection on Instagram using BERT and ensemble classifiers. *Expert Systems with Applications*, 195, 116571.
11. Ali, M., & Khan, M. N. (2023). A hybrid deep ensemble model for spam comment detection using semantic and emoji features. *IEEE Access*, 11, 22564–22575.
12. Zhou, J., & Li, T. (2023). Detecting spam with sentiment-aware emoji usage in multimodal comment-post structures. *Information Processing & Management*, 60(2), 103229.
13. Chen, Q., & Wu, X. (2023). Spam detection on Instagram with visual-linguistic emoji fusion and hierarchical ensemble learning. *Pattern Recognition Letters*, 165, 1–8.
14. Srinivasan, K., & Thomas, M. (2024). Post-comment relation modeling for spam detection on Facebook using stacked ensemble classifiers. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 18(1), Article 12.
15. Rahman, A., & Hasan, M. (2024). Multi-view ensemble learning for spam detection using emoji co-occurrence and comment semantics. *Knowledge-Based Systems*, 286, 110179.

single-model approaches in terms of precision, recall, and accuracy by utilizing the capabilities of multiple machine learning classifiers. To achieve more resilient performance on social media networks with imbalanced and chaotic datasets, consider employing Random Forest, Gradient Boosting, or Voting Classifiers. The culmination of the integration of these sophisticated characteristics and methodologies is a spam detection system that is more contextually sensitive and intelligent. This facilitates additional advancements in automated content regulation, in addition to enhancing the platform's security and user experience.