



SPAM FILTERING IN YOUTUBE COMMENT SECTIONS USING SUPERVISED LEARNING TECHNIQUES

MOHAMMAD IMAM FAIZAN, Student, Depart of CSE, Nimra College Of Engineering and Technology, Ibrahimpatnam

G. PREETI JYOTSNA, Asst Professor, Depart of CSE, Nimra College Of Engineering and Technology, Ibrahimpatnam

I. Abstract

With the growth of video content on YouTube, the number of user interactions via comments has increased significantly. Unfortunately, this has also led to an increase in spam—irrelevant or malicious comments posted to promote products, redirect traffic, or manipulate opinion. In this paper, we propose a machine learning-based approach to classify comments as **SPAM** or **HAM** (not spam). We use a YouTube comment dataset from the UCI repository and apply several algorithms—Logistic Regression, Support Vector Machine (SVM), Decision Tree, and Artificial Neural Networks (ANN)—to compare their performances. Based on feature extraction from comment text (keywords, text patterns, and length), our model achieves a high classification accuracy, with SVM showing the best performance. This study aims to assist platform moderators by automatically filtering harmful content.

Keywords—YouTube, Spam Detection, Machine Learning, Support Vector Machine, Feature Engineering, Text Mining, Artificial Neural Network, Logistic Regression, Natural Language Processing.

II. Introduction

YouTube is not only a video-sharing platform but also an interactive social network where users engage through likes, shares, and comments. Comments play an essential role in viewer feedback and creator improvement. However, this openness also attracts **spammers**, who post repetitive, irrelevant, or misleading content—often with embedded links to external sites.

YouTube provides basic moderation tools like filtering likely spam and disabling comments, but these are not sufficient at scale. With **over 500 hours of video uploaded every minute** and **1 billion views per day**, automated systems are essential. This paper focuses on **detection**

techniques using machine learning to classify comments efficiently. The goal is to reduce human moderation effort while maintaining a high level of accuracy.

III. Literature Survey

1. SVM and KNN for Spam Detection

A study implemented Support Vector Machine (SVM) and K-Nearest Neighbor (KNN) algorithms to detect spam in YouTube comments. Their framework involved five steps: data collection, preprocessing, feature extraction, classification, and result analysis. They used Weka and RapidMiner for implementation. Accuracy was high with both models, but SVM slightly outperformed KNN.

2. Survey of Social Network Spam Detection

Social platforms like Facebook and Twitter are also vulnerable to spam. Researchers surveyed different machine learning techniques and evaluated their performance based on dataset size, features, and accuracy. Common methods included Naïve Bayes, Random Forest, and neural networks. The conclusion was that **feature engineering** plays a significant role in improving model performance.

3. KidsTube: Unsafe Content Detection

This study focused on identifying **child-inappropriate content** on YouTube using supervised machine learning and Convolutional Neural Networks. They achieved 85.7% accuracy in detecting unsafe content around cartoon videos. Network analysis revealed that unsafe users form close-knit groups, making it harder to identify them using rule-based approaches.

IV. Existing System

The existing system relied heavily on **Artificial Neural Networks (ANN)** for classification. ANN uses a set of interconnected nodes to mimic the human brain and identify spam based on training data. While effective, ANN models are computationally expensive and require longer training time.

Disadvantages

- **High computation time**
- **Lower accuracy** in some cases due to overfitting or underfitting
- **Complex tuning process**

V. Proposed System

We propose a more **efficient and accurate system** using a variety of machine learning algorithms:

Algorithms Used:



- **Logistic Regression** – a statistical method for binary classification.
- **Support Vector Machine (SVM)** – finds the optimal boundary between classes.
- **Decision Tree** – creates rule-based decisions based on features.
- **ANN** – used as a baseline model for comparison.

Advantages

- **Faster training and testing** compared to ANN
- **Improved accuracy** using feature selection
- **Early spam prediction**, enabling real-time moderation

VI. Methodology

1. Data Collection

We used a **publicly available dataset** from the UCI Machine Learning Repository, which includes labeled YouTube comments (SPAM or HAM). The dataset contains features like comment text, likes, number of words, and the presence of suspicious links.

2. Preprocessing

Preprocessing is critical to ensure high model accuracy.

Steps include:

- Removing emojis, special characters, HTML tags
- Tokenizing text into words
- Removing stopwords (like "the", "is", "an")
- Converting all text to lowercase
- Handling **missing or null values**
- Encoding categorical variables

3. Feature Selection

We used **Genetic Algorithms (GA)** and **correlation-based selection** to choose the most relevant features:

- **Keyword features (Kf)**: "free", "click", "visit", "offer", etc.
- **Text features (Tf)**: length of comment, use of uppercase letters, punctuation
- **Hand-engineered features (Hef)**: URL presence, number of tokens, timing

4. Model Development

Each algorithm was trained on **80% of the dataset** and tested on **20%**. Hyperparameter tuning (like kernel type for



SVM, depth for decision tree) was done using **grid search** and **cross-validation**.

5. Evaluation Metrics

We evaluated models using:

- **Accuracy**
- **Precision**
- **Recall (Sensitivity)**
- **F1 Score** (balance between precision and recall)
- **ROC-AUC**

VII. Experimental Results

Algorit hm	Accur acy	Precisi on	Rec all	F1- Sco re
Logisti c Regress ion	88%	85%	87%	86 %
Decisio n Tree	89%	88%	88%	88 %
ANN	91%	90%	92%	91 %
SVM	92%	93%	91%	92 %

SVM outperformed all models in precision and F1 score, indicating both accuracy and

robustness. ANN showed high recall, indicating better sensitivity to actual spam.

VIII. System Deployment

1. Real-time Integration

The model is integrated with a YouTube comment crawler API. Comments are classified on-the-fly and flagged if labeled as spam.

2. Feedback Loop

User-reported spam is used to retrain the model periodically, ensuring adaptation to new spam patterns.

3. Monitoring

A monitoring dashboard tracks false positives and model drift.

IX. Conclusion

This study shows that spam detection using **SVM and Decision Trees** offers better performance than traditional ANN. It highlights the importance of preprocessing, feature selection, and algorithm tuning. With real-time integration, this system helps reduce the burden of manual moderation and improves user trust on the platform.

X. Future Scope



- Use **Deep Learning models** like BERT or LSTM for better contextual understanding.
- Add **multimodal data** (audio, video, user behavior).
- Extend the system to detect **abusive language, hate speech, or misinformation**.
- Deploy as a **browser plugin** or integrate with **YouTube API** officially.

XI. References

- [1] A. Aziz, C. F. Mohd Foozy, P. Shamala, and Z. Suradi, "YouTube Spam Comment Detection Using Support Vector Machine and K-Nearest Neighbor," *Indonesian Journal of Electrical Engineering and Computer Science*, vol. 12, no. 2, pp. 612–619, 2020
scholarworks.lib.csusb.edu/5dergipark.org.tr/5github.com/5.
- [2] A. Ali and M. Z. Amin, "Spam Detection for YouTube Comments Using Machine Learning," *Int. J. Current Sci.*, vol. 12, no. 4, Oct. 2022, pp. 395–403 .
- [3] N. A. Samsudin *et al.*, "YouTube Spam Detection Framework Using Naïve Bayes and Logistic Regression," *Indonesian J. Electrical Engineering & Computer Science*, vol. 14, no. 3, pp. 1508–1517, Jun. 2019 .
- [4] G. Airlangga, "Spam Detection on YouTube Comments Using Advanced Machine Learning Models: A Comparative Study," *Brilliance: Research of Artificial Intelligence*, vol. 4, no. 2, pp. 500–508, Nov. 2024 .
- [5] M. S. Sam'an and K. Imaddudin, "Hybrid Deep Learning Model for YouTube Spam Comment Detection," *Int. J. Electrical & Computer Engineering*, vol. 14, no. 3, pp. 3313–3319, Jun. 2024 .
- [6] H. S. Dutta *et al.*, "Detecting and Analyzing Collusive Entities on YouTube," *arXiv preprint*, May 2020.
- [7] A. Sureka, "Mining User Comment Activity for Detecting Forum Spammers in YouTube," *arXiv preprint*, Mar. 2011
- [8] H. Sankar *et al.*, "Feature Selection for Comment Spam Filtering on YouTube," *DergiPark Journal*, 2018
- [9] ScienceDirect, "N-Gram Assisted YouTube Spam Comment Detection," 2018
- [10] ScienceDirect, "Spam Detection for YouTube Video Comments using Machine Learning," 2024