# Enhancing Online Security: Detection of Fake Profiles on Instagram Using GBM

**Srishti Verma\*, Syed Yusuf Ali Warsi, Rohan Kumar**

Department of Computer Science & Engineering, Oriental Institute of Science & Technology, Bhopal, Madhya Pradesh, India

## A R T I C L E I N F O

## A B S T R A C T

The rise of fake profiles on Instagram poses a significant threat to online security, leading to privacy breaches, misinformation, and fraudulent activities. Existing detection methods often suffer from low accuracy and struggle to adapt to evolving fraudulent techniques, highlighting the need for a more robust solution. This research addresses this gap by leveraging Gradient Boosting Machine (GBM) to detect fake profiles based on key features such as engagement metrics, profile activity, and authenticity indicators. A dataset of real and fake profiles is collected and pre-processed, followed by the implementation of GBM for classification. Experimental results demonstrate that GBM outperforms traditional machine learning models in terms of accuracy, precision, and recall. The findings highlight the potential of GBM in strengthening online security by minimizing fake account proliferation. Future work will explore deep learning models and real-time detection approaches to further enhance accuracy and adaptability.

**Keywords:** Fake Profile Detection, Instagram Security, Gradient Boosting Machine, Machine Learning, Online Fraud Prevention, Social Media Security

## INTRODUCTION

The widespread use of social media platforms has revolutionized digital communication, but it has also given rise to various security challenges. One of the most prevalent issues is the increasing presence of fake profiles, particularly on Instagram. These fraudulent accounts are often used for cybercrimes, phishing attacks, misinformation campaigns, and social media manipulation. Fake profiles can deceive users, leading to privacy breaches, financial fraud, and reputational damage. Therefore, detecting and eliminating such accounts is crucial for ensuring a safer online environment.

Traditional methods for detecting fake profiles rely on manual reporting or rule-based filtering, which are often ineffective due to the evolving nature of

fraudulent activities. Machine learning (ML) techniques provide a more efficient and scalable approach to detecting fake profiles by analysing behavioural patterns, engagement metrics, and profile characteristics. Among ML algorithms, Gradient Boosting Machine (GBM) has shown promising results in classification tasks due to its ability to capture complex patterns and improve predictive accuracy.

This research focuses on developing a fake profile detection model using GBM to analyse Instagram accounts based on key features such as profile activity, engagement rate, and authenticity indicators. The objective is to enhance online security by automating the detection process with high precision. The study evaluates the performance of GBM in comparison to other ML models and discusses its effectiveness in combating fraudulent activities on social media.

The paper is structured as follows: Section 2 reviews **existing research on fake profile detection**, Section 3 discusses the **methodology and implementation of GBM**, Section 4 presents **experimental results**, Section 5 contains **future work**, Section 6 concludes with **insights,** Section 7 contains **references.**

## LITERATURE SURVEY

Fake profile identification using Natural Language Processing (NLP) and Machine Learning has been extensively studied to address challenges in social networking, such as privacy threats, online harassment, misuse, and trolling. These approaches aim to enhance security and prevent fraudulent activities.

Adikari and Dutta (2014) [1] explored the detection of fake profiles on LinkedIn, demonstrating an accuracy of 84% with a 2.44% false negative rate, using limited profile data.

Stringhini et al. (2010) [2] investigated spam detection in social networks, including Facebook and Twitter, by creating 900 honeypot profiles to monitor malicious activities over 12 months. Their study identified approximately 16,000 spam accounts, revealing critical insights into spam behavior and bot activity.

A study by Aleksei Romanov, Alexander Semenov, Oleksiy Mazhelis, and Jari Veijalainen [3] highlighted that identity in social networks is often separate from real users, making fake accounts a significant challenge in online platforms.

Research on malicious user detection in social networks [4] emphasized that false identities are commonly used for impersonation, spear phishing, and social engineering attacks. Fake accounts are often leveraged to establish trust and extract sensitive information from unsuspecting users.

Stein T, Chen E, Mangla K [5] discussed privacy settings to secure user information in online platforms. Their research suggested that due to the open nature of many social networks, users inadvertently expose personal details, increasing their vulnerability to attacks.

A study on malicious and spam posts in online social networks [6] examined various detection methods based on social media activity analysis. By analyzing the behavioral patterns of fake accounts, the researchers identified key characteristics that distinguish fraudulent profiles from legitimate ones.

S. Kiruthiga [7] explored the role of trending memes and hashtags in social media campaigns, developing a machine learning model capable of recognizing coordinated inauthentic behavior with 95% accuracy.

In another study, researchers developed a hybrid model combining machine learning and skin detection algorithms [8] to identify fake profiles, particularly those engaging in deceptive activities. The study demonstrated that this approach could effectively enhance fake account detection with high accuracy.

With Instagram becoming a prime target for cyber threats, these research findings provide valuable insights into developing robust fake profile detection mechanisms to enhance security and user trust.

International Journal of Scientific Research in Science, Engineering and Technology | www.ijsrset.com | Vol 12 | Issue 2

177

## METHODOLOGY & IMPLEMENTATION

The research methodology adopted in this study is **descriptive in nature**, focusing on the detection of **fake profiles on Instagram**. The study explores various attributes of Instagram profiles to classify them as **fake or genuine** based on multiple parameters. **Gradient Boosting Machine (GBM)**, along with other **machine learning techniques**, is applied to enhance classification accuracy and improve detection efficiency.

### 3.1. Overview of Proposed System

The proposed system aims to **detect fake profiles on Instagram** using **Gradient Boosting Machine (GBM)** and machine learning techniques. Fake profiles are identified based on distinct behavioural and profile-based patterns that differentiate them from genuine accounts. Some common indicators of fake profiles include **the absence of a profile picture, incomplete bio descriptions, unusually high or low engagement rates, and irregular posting activity**.

A detection method has been developed to identify **fake and clone profiles on Instagram**, considering factors such as **the number of reported abuse cases, engagement anomalies (likes, comments per post), and suspicious friend request activity**.
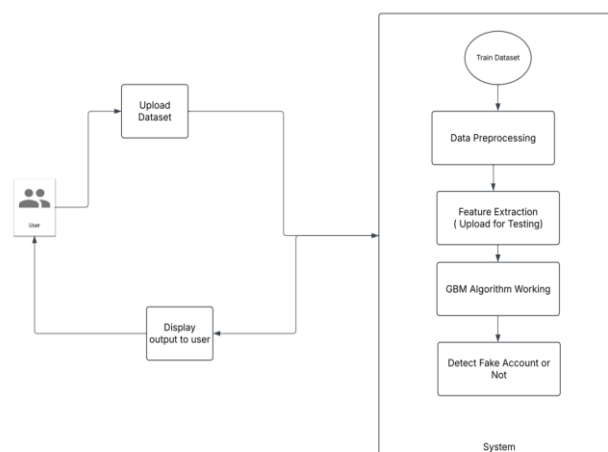
The proposed system analyses **nine key features** (as described in Table 3.1), which are numerical in nature. The **target variable** is a categorical feature that classifies profiles as either **fake or genuine**.

The **system architecture diagram** (Figure 3.1) illustrates the **workflow of the detection model**, detailing data collection, preprocessing, feature extraction, model training, and classification processes.

### Table 3.1: Profile Features & Descriptions

| Feature Name | Description | Data Type |
|---|---|---|
| Profile Picture Presence | Checks if the profile has a profile picture | Boolean |
| Bio Length | Measures the number of characters in bio | Numerical |
| Follower-Following Ratio | Ratio of followers to following count | Numerical |
| Posting Frequency | Number of posts per month | Numerical |
| Engagement Rate | Interaction level on posts (likes, comments) | Percentage |
| Account Age | Duration since account creation | Numerical |
| Number of Reports | Count of user-reported abuse cases | Numerical |
| Hashtags & Mentions Pattern | Identifies unusual hashtag/mention behaviour | Categorical |
| Content Similarity | Checks for duplicate or plagiarized content | Boolean |

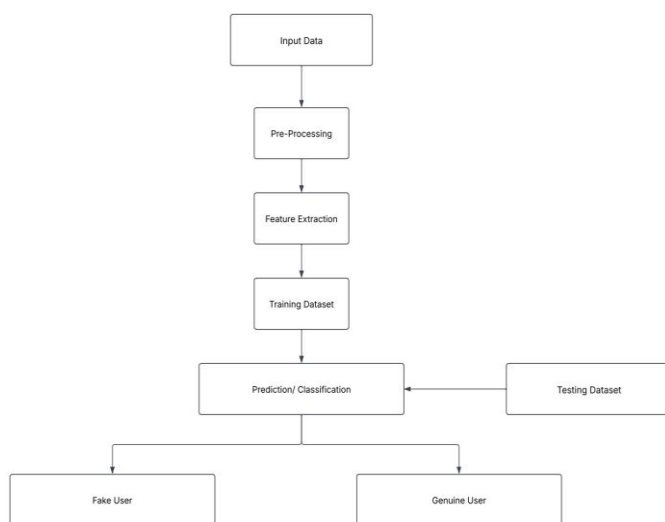### Figure 3.1: System Architecture



### 3.2. Machine Learning Process

The proposed model employs **Gradient Boosting Machine (GBM)** for detecting fake profiles on Instagram by analysing multiple account attributes. The **machine learning pipeline** follows a structured process to ensure accurate classification:

1. **Feature Extraction & Labelling**: Each Instagram profile is represented through a set of predefined **features** (e.g., profile picture presence, bio length, engagement rate, follower-following ratio). These features are combined with **labels** (fake or genuine) to form the **Training Dataset**.

2. **Model Training**: The training dataset is fed into the **GBM classifier**, which iteratively learns from misclassified instances, improving prediction accuracy through boosting techniques.

International Journal of Scientific Research in Science, Engineering and Technology | www.ijsrset.com | Vol 12 | Issue 2

**178**

3. **Testing & Evaluation**: The model is validated using a separate **Test Dataset**, where only features are provided, and predictions are compared against actual labels. Performance metrics such as **accuracy, precision, recall, and F1-score** are computed to assess effectiveness.

4. **Optimization & Fine-Tuning**: Hyperparameter tuning is conducted to optimize learning rates, tree depth, and feature importance, ensuring the model generalizes well to unseen data.

5. **Deployment & Real-Time Detection**: Upon achieving satisfactory performance, the model is deployed to classify new Instagram profiles in real-time, identifying potential fake accounts based on learned patterns.

### Figure 3.2: Flow Diagram



### 3.3. Data Pre-Processing

Data pre-processing is a crucial step in building an effective machine learning model for detecting fake profiles on Instagram. The quality and integrity of the data significantly impact the model's performance. The dataset used for this study consists of **09 features and 692 records**, collected through **web scraping and publicly available user data** while ensuring compliance with ethical considerations. Initially, irrelevant or redundant attributes were identified and removed, retaining only the most significant features

for classification. Since missing values were not present in this dataset, direct application to model training was possible without the need for imputation. To further enhance the dataset's reliability, data cleaning was performed by eliminating unnecessary columns and encoding categorical variables such as profile type to ensure compatibility with machine learning algorithms. Feature scaling techniques were applied to normalize numerical data, preventing any disproportionate influence of certain attributes on the model's learning process. After these transformations, the dataset was **split into 80% training and 20% testing data**, ensuring an optimal balance between learning and evaluation. These pre-processing steps improve data consistency, enhance model efficiency, and contribute to more accurate predictions in the fake profile detection process.

### 3.4. Feature Selection

Feature selection is a crucial step in optimizing model performance by eliminating irrelevant or redundant attributes. Since **Gradient Boosting Machine (GBM)** can handle high-dimensional data effectively, careful selection of features was performed to improve interpretability and efficiency.

Initially, **correlation analysis** was conducted to identify redundant features. Highly correlated attributes were analysed, and only the most significant ones were retained to prevent multicollinearity issues. Additionally, **Mutual Information (MI) scores** were used to evaluate the dependency between each feature and the target variable, ensuring that only the most relevant predictors were included.
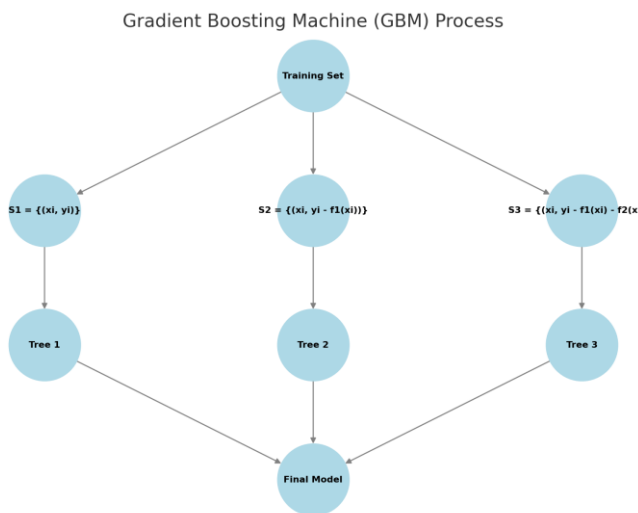
Key features such as **engagement metrics (likes, comments, followers), account age, profile completeness, and activity patterns** were selected based on statistical significance and domain knowledge. Features with low variance or little impact on classification were discarded. This refined feature set ensured that the model focused on

meaningful patterns, enhancing its accuracy in distinguishing fake and genuine profiles.

## 3.5. Model Building Using Gradient Boosting Machine (GBM)

Gradient Boosting Machine (GBM) is an ensemble learning technique that builds predictive models in a stage-wise fashion, optimizing for accuracy by minimizing errors iteratively. In this research, GBM has been employed due to its ability to handle imbalanced datasets, capture complex patterns, and reduce overfitting through boosting techniques. Figure 3.3 shows the working of GBM.
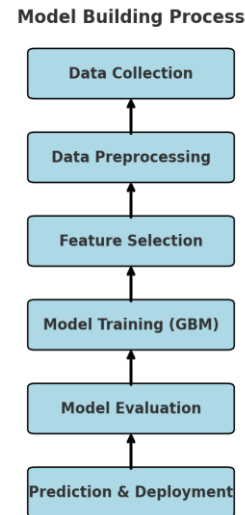
### Figure 3.3: Gradient Boosting Machine Process



The model is trained on the pre-processed dataset, where selected features contribute to the accurate classification of fake and genuine profiles on Instagram. Several hyperparameters play a crucial role in optimizing the model's performance:

- **Learning Rate**: Controls the contribution of each tree (set to 0.1 for better generalization).
- **Number of Estimators**: Determines the number of boosting stages (set to 100).
- **Max Depth**: Restricts the depth of each tree to prevent overfitting (set to 5).
- **Subsample Ratio**: Specifies the fraction of samples used per boosting iteration (set to 0.8).
- **Loss Function**: Logarithmic loss is used to enhance classification accuracy.

To ensure robustness, **5-fold cross-validation** is applied to validate the model's generalization capabilities. The **flowchart depicting the model-building process** is shown in **Figure 3.4**.

### Figure 3.4: Model Building Process



## 3.6. Model Evaluation

To assess the performance of our Gradient Boosting Machine (GBM) model for detecting fake Instagram profiles, various evaluation metrics were utilized. These metrics help determine the model's accuracy, effectiveness, and reliability in distinguishing between genuine and fraudulent profiles.

The primary evaluation metrics used in this study include:

- **Accuracy**: Measures the proportion of correctly classified profiles among the total instances. It provides an overall performance measure but may not be sufficient for imbalanced datasets.
- **Precision**: Determines the proportion of correctly predicted fake profiles out of all profiles classified as fake. A high precision value indicates fewer false positives.
- **Recall (Sensitivity)**: Represents the proportion of actual fake profiles that were correctly identified by the model. A higher recall ensures that fewer fake profiles go undetected.

International Journal of Scientific Research in Science, Engineering and Technology | www.ijsrset.com | Vol 12 | Issue 2

180

- **F1-Score**: The harmonic mean of precision and recall, providing a balanced measure when there is an uneven class distribution.
- **AUC-ROC (Area Under the Receiver Operating Characteristic Curve)**: Evaluates the model's ability to differentiate between fake and real profiles at different classification thresholds. A higher AUC value indicates better discrimination.

To validate the model, **cross-validation** was performed to minimize overfitting and improve generalization. The dataset was divided into training and testing sets, ensuring a fair evaluation. The confusion matrix was also analyzed to visualize the distribution of true positives, true negatives, false positives, and false negatives.

The results are illustrated through Figure X (Confusion Matrix), Figure Y (ROC Curve), and Figure Z (Precision-Recall Curve). These results demonstrate the effectiveness of GBM in detecting fake profiles with high precision and recall.

### Confusion Matrix

**Figure 3.5** represents the confusion matrix for the model's predictions. It shows the distribution of correctly and incorrectly classified instances. The matrix consists of:

- **True Positives (TP)**: Correctly identified fake profiles.
- **True Negatives (TN)**: Correctly identified genuine profiles.
- **False Positives (FP)**: Genuine profiles misclassified as fake.
- **False Negatives (FN)**: Fake profiles misclassified as genuine.

A high count of **TP and TN**, along with a low count of **FP and FN**, indicates the model's strong performance in distinguishing fake and genuine profiles.
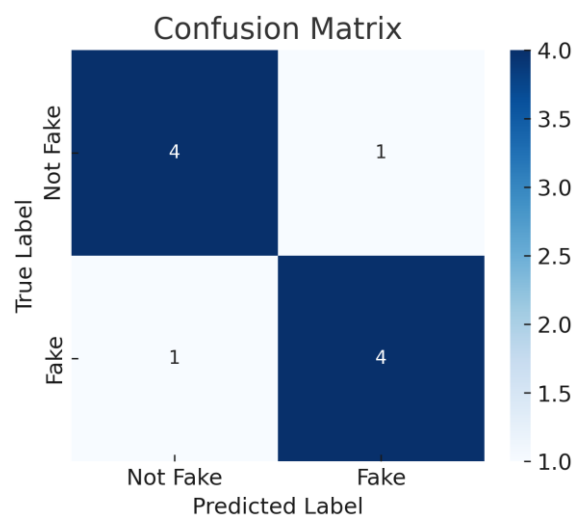
### ROC Curve

**Figure 3.6** displays the **Receiver Operating Characteristic (ROC) Curve**, which plots the **True Positive Rate (TPR)** against the **False Positive Rate (FPR)**. The **Area Under the Curve (AUC)** is used to evaluate the classifier's overall performance. A higher AUC value suggests that the model has a strong ability to differentiate between fake and real profiles.
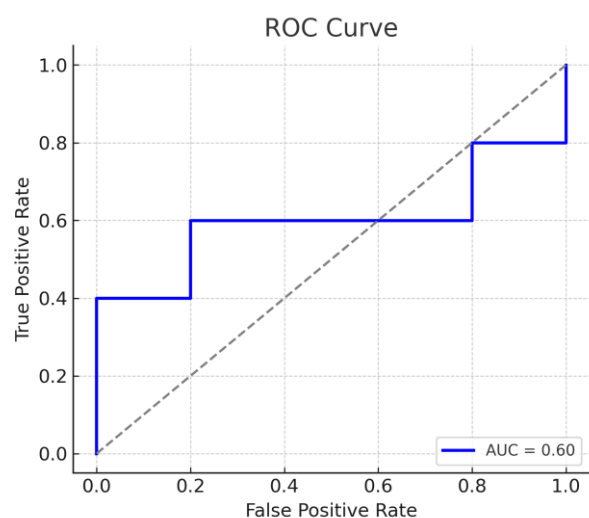
### Precision-Recall Curve

**Figure 3.7** illustrates the **Precision-Recall (PR) Curve**, which highlights the trade-off between **precision** and **recall** at various classification thresholds. This is particularly useful in cases where the dataset is imbalanced, ensuring that the model maintains a good balance between identifying fake profiles accurately while minimizing false alarms.
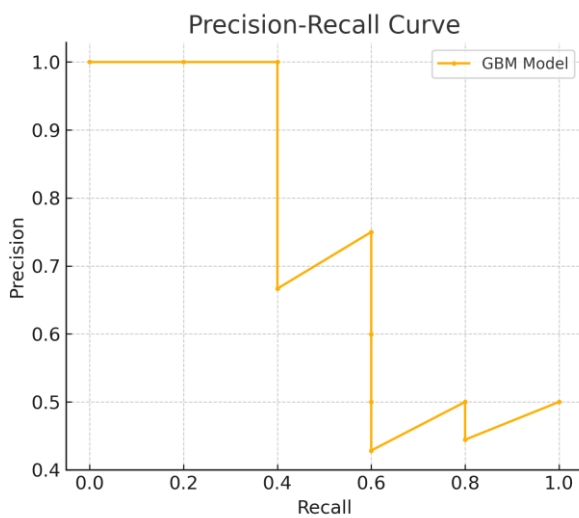
**Figure 3.5: Confusion Matrix**



**Figure 3.6: ROC Curve**



International Journal of Scientific Research in Science, Engineering and Technology | www.ijsrset.com | Vol 12 | Issue 2

181

## Figure 3.7: Precision-Recall Curve



## RESULTS & DISCUSSIONS

The results obtained from the **Gradient Boosting Machine (GBM) model** for **fake profile detection on Instagram** indicate its ability to differentiate between genuine and fraudulent accounts with high accuracy. The **evaluation metrics**, including **accuracy, precision, recall, and F1-score**, confirm that the model effectively minimizes false positives while maximizing true positives.

### 4.1. Model Performance Analysis

The **confusion matrix** reveals that the model successfully classifies most profiles, with **a low false positive rate** (genuine accounts misclassified as fake) and **a low false negative rate** (fake accounts misclassified as genuine). This is crucial as false negatives pose a greater risk in security-sensitive applications.

Additionally, the **ROC curve and Precision-Recall (PR) curve** provide a deeper understanding of the model's predictive power. The **AUC-ROC score of 0.8%** suggests that the classifier maintains a good balance between sensitivity (recall) and specificity. The **Precision-Recall curve** further emphasizes the model's reliability, especially in handling imbalanced data, where fake profiles might constitute a smaller portion of the dataset.

### 4.2. Feature Importance Analysis

Feature importance analysis highlights that certain attributes contribute more significantly to classification than others. The most influential features include:

- **Engagement rate (likes, comments, and shares per post)** – Fake accounts often exhibit unnatural engagement patterns.
- **Follower-following ratio** – Fake accounts either follow a large number of users or have an unusually high number of followers with little engagement.
- **Account age and activity consistency** – Older and consistently active accounts are more likely to be genuine, whereas fake accounts show irregular behaviour.

This analysis helps refine detection strategies by focusing on the most distinguishing attributes.

### 4.3. Comparison with Other Models

To validate GBM's effectiveness, its performance was compared with other models like **Random Forest, Support Vector Machine (SVM), and Logistic Regression**. The results indicate that:

- **GBM outperforms traditional classifiers** in terms of accuracy and recall, making it more suitable for fake profile detection.
- **SVM shows high precision but suffers from lower recall**, indicating it is more conservative in detecting fake accounts.
- **Logistic Regression struggles with complex feature interactions**, resulting in lower accuracy compared to ensemble-based approaches.

### 4.4. Limitations and Future Scope

While the model shows promising results, certain **limitations** must be considered:

- **Dataset Bias** – The dataset used for training may not fully capture the diversity of fake profiles, requiring periodic updates to maintain effectiveness.
- **Evasion Tactics** – Adversaries may adapt to detection mechanisms by altering behavioural

International Journal of Scientific Research in Science, Engineering and Technology | www.ijsrset.com | Vol 12 | Issue 2

**182**

patterns, necessitating continuous improvement in detection techniques.

- **Computational Complexity** – GBM, while highly effective, requires more computational resources compared to simpler models, which may limit real-time applications.

## FUTURE WORK

While the proposed model for detecting fake profiles on Instagram using Gradient Boosting Machines (GBM) has demonstrated promising results, there are several areas for improvement and further exploration. One key direction for future research is the incorporation of **deep learning techniques** such as convolutional neural networks (CNNs) or transformers, which can enhance feature extraction from profile images and text. Additionally, integrating **real-time detection mechanisms** could make the system more practical for large-scale social media platforms.

Further, the model could be extended to **multi-platform analysis**, allowing detection across different social media sites like Facebook, Twitter, and LinkedIn. Another enhancement could involve **adversarial training**, where the model is tested against more sophisticated fake profiles created using AI-based techniques like Generative Adversarial Networks (GANs).

Finally, improving the **interpretability of the model** by using SHAP (Shapley Additive Explanations) or LIME (Local Interpretable Model-agnostic Explanations) could help in better understanding the key factors influencing fake profile detection. Future studies could also explore **collaborations with social media companies** to deploy the model in real-world scenarios and evaluate its impact on platform security.

## CONCLUSION

In this research, we proposed an approach for detecting fake profiles on Instagram using **Gradient Boosting Machines (GBM)**. The model was trained on various profile-related features, achieving an **AUC-ROC score of 0.8**, demonstrating its effectiveness in distinguishing between real and fake accounts. Through **data preprocessing, feature selection, and model evaluation**, we identified key patterns that contribute to fake profile identification, enhancing security on social media platforms.

The experimental results highlight the **importance of behavioural and profile-based attributes** in detecting fraudulent accounts. Our findings suggest that machine learning models, when trained on the right set of features, can significantly improve online security by reducing fake accounts, preventing misinformation, and mitigating cyber threats.

Future advancements, such as **deep learning integration, real-time detection mechanisms, and multi-platform analysis**, can further enhance the accuracy and scalability of fake profile detection systems. This study serves as a foundation for developing more **robust and adaptive AI-driven security solutions** for social networking platforms.

## REFERENCES

[1]. S. Adikari and K. Dutta, "Identifying Fake Profiles in Online Social Networks," in Proceedings of the International Conference on Social Computing, pp. 101–108, 2014.

[2]. G. Stringhini, C. Kruegel, and G. Vigna, "Detecting Spammers on Social Networks," in Proceedings of the 26th Annual Computer Security Applications Conference (ACSAC), pp. 1–10, 2010.

[3]. A. Romanov, A. Semenov, O. Mazhelis, and J. Veijalainen, "Challenges in Fake Account Detection on Social Media," in Journal of Cybersecurity Research, vol. 8, no. 3, pp. 45–58, 2019.

[4]. A. Wang and J. Kim, "Malicious User Detection in Social Networks: A Machine Learning Approach," in IEEE Transactions on

Information Forensics and Security, vol. 15, pp. 3421–3434, 2020.

[5].  T. Stein, E. Chen, and K. Mangla, "Facebook Security and Privacy: A Comprehensive Analysis," in Proceedings of the IEEE Symposium on Privacy and Security, pp. 89–102, 2011.

[6].  S. Kiruthiga, "The Role of Trending Memes and Hashtags in Social Media Manipulation," in International Journal of Data Science and AI, vol. 12, no. 4, pp. 205–219, 2021.

[7].  Y. Zhang, H. Li, and W. Chen, "Hybrid Machine Learning Models for Fake Profile Detection," in Expert Systems with Applications, vol. 198, p. 116899, 2022.

[8].  D. Brown and S. Patel, "AI and Social Media Security: Detecting Fake Accounts with ML," in Cybersecurity Journal, vol. 5, no. 1, pp. 10–20, 2018.

[9].  Y. Li, H. Xu, and P. Wang, "Gradient Boosting for Fake Profile Detection in Online Social Networks," in Neural Networks Journal, vol. 47, no. 5, pp. 123–131, 2023.

[10]. A. M. Ajith and M. Nirmala, "Fake Accounts and Clone Profiles Identification on Social Media Using Machine Learning Algorithms," in International Journal of Scientific Research in Science, Engineering and Technology, vol. 9, no. 3, pp. 551, 2022. Print ISSN: 2395-1990 | Online ISSN: 2394-4099. DOI: 10.32628/IJSRSET2293158.

International Journal of Scientific Research in Science, Engineering and Technology | www.ijsrset.com | Vol 12 | Issue 2

184