

# Fake Social Media Profile Detection and Reporting

Mohammed Maaz Rehman, Prasad P S, Manoj J, Anjan G M

## Abstract

Our project is based on the design and deployment of a high-end machine learning model for the identification of spurious social media profiles. The project uses large datasets with attributes of actual and spurious profiles, which are mixed up in a deep neural network model for efficient profile classification. Major contributions are the strategic attribute design of metrics like follower-to-following ratio and engagement metrics, as well as the application of SMOTE to address class imbalance. Model performance is rigorously tested using a range of metrics like accuracy, precision, recall, AUC, and F1-score. Also, the system design is optimized for real-world deployment with the modular backend implemented in FastAPI using Python and the dynamic and interactive frontend implemented using Node.js. This configuration not only provides high detection accuracy but also is simple to deploy and maintain, breaking new ground in the battle against online fake personas.

Keywords: Fake Profile Detection, Social Media Machine Learning, Feature Engineering, Deep Learning

## Introduction

The proliferation of impersonation social media accounts has changed the World Wide Web substantially, with sweeping consequences for trust, privacy, and security. These impersonation profiles are not merely harmless objects; they can shape public opinion on a grand scale, which is most often utilized disseminating misinformation or propaganda. This power can have real-world consequences, influencing elections, public health policy, and even personal faith. Beyond the realm of opinion, pseudonymous profiles facilitate a variety of cybercrimes, including identity theft and phishing, where unsuspecting users are tricked into divulging sensitive information or funds. The authenticity of social networks, which are grounded on actual interactions, is eroded

significantly, leading to a suspicious environment where the genuineness of users' participation is always questionable.

## 1.2 Broader Societal Implications

The influence of spurious profiles infiltrates the corporate sector as well, warping market studies by providing biased data, therefore falsifying metrics which companies use to determine their marketing and product plans. Such forgery has the potential to lead to ineffective marketing campaigns and uninformed business choices. These profiles also lead to the formation of echo chambers, presenting users with similar-minded content, promoting polarization among groups. Such settings are likely to reinforce social conflicts and shut the window of productive conversation and mutual perception. The credibility of social platforms, having been recognized as bridges of global communication, therefore stands compromised, undermining the general faith within the digital community. This echoes on the necessity of having rigorous detection mechanisms in position, not as a technological panacea but as an indispensable component of digital social health.

## 1.3 The Need for Innovation

Present Recognizing the complex threats that fake profiles pose, there is a pressing social need to develop detection technologies. The challenge goes beyond mere data analysis; it needs advanced algorithms that can keep up with the changing tactics of those who operate these fraudulent accounts. Our project seeks to rectify this through the creation of a machine learning algorithm that is able to detect and categorize profiles with high accuracy, thus providing some level of trust and security to online interactions. Doing so, not only do we seek to safeguard individual users but also to ensure the integrity of social media platforms.

## 2. Literature Survey

Kumar et al. (2018)

During the current research, the Random Forest algorithm was used to identify genuine and false user accounts based on behavioral characteristics. Simple characteristics were posting rate, follower-to-following ratio, and routine participation. Such characteristics gave useful insights regarding the classification of user genuineness.

Yang et al. (2019)

Through models of neural networks, the research was keen to identify anomalies in social interaction. Based on the analysis of interaction graphs, the model could identify user profiles that had atypical behavior patterns, which are indicative of fraudulent behavior.

Egele et al. (2020).

The authors also proposed a graph-structured detection method for the identification of coordinated deceptive activity. The model would be able to identify groups of fake or bot-controlled accounts that were operating in coordination, according to the analysis of inter-account links.

Cao et al. (2012)

Their work introduced a hybrid method that was reliant on social graph-based and text-based features to identify spam profiles. Machine learning classifiers were used to label suspicious profiles based on social graph properties and text information.

Al-Qurishi et al. (2018)

A two-pronged approach was suggested, one on supervised and the other on unsupervised learning techniques. Classification considered metadata of the posts, user behavior patterns, and timelines of behavior to identify real and spurious profiles.

Viswanath et al. (2014)

This study centered on the fabricated Twitter profiles by examining how the follower networks evolved over time. Processual dynamics and network structure were the most significant factors in determining anomalous account behavior and bot-driven interaction.

Fire et al. (2013)

Focusing particularly on LinkedIn, this study used social network analysis and machine learning to identify spurious profiles. The model analyzed metrics like network density, profile completeness, and centrality measures for the purpose of authenticity verification of the accounts.

Stringhini et al. (2010)

Targeting detection of social media spammers, the research suggested a set of integrated features that included both behavioral patterns and nontent-related features. These were processed through a machine learning paradigm to identify malicious accounts with good effectiveness.

Boshmaf et al. (2011)

The idea of "socialbots" was first proposed here—robotic accounts designed to mimic normal user activity. Detection was focused on nontent interaction and behavior analysis of the user to reveal such fake profiles.

Conti et al. (2012)

In their "FakeBook" framework, the authors addressed fake profile detection from both static and dynamic features. Their framework evaluated profile features like completeness and monitored how user interactions change over time to detect abnormal activity.

## 3. Proposed Methodology

### 3.1 Data Gathering

We start our method with the collection of user profile information from social media sites. We worked with two large JSON-formatted datasets, each of which consisted of labeled entries to determine if an account was genuine or not. Each user profile contained several attributes like:

Follower Count: Total followers on the account. Following

Count: Accounts followed by the user. Biography Length:

Character length of the user's bio.

Media Count: Number of content posted, such as images and videos.

Profile Picture Status: Whether the user has uploaded a custom profile picture or defaults to it.

Account Privacy: Whether the account is public or private.

To ensure the quality of our data, we performed extensive validation checks. This involved removing anomalies such as negative values and correcting or deleting data entries that were incomplete or clearly defective.

### 3.2 Data Preprocessing

After being gathered, the data was cleaned up through a variety of preprocessing processes:

**Missing Value Handling:** Numerical columns were filled with missing values using mean or median, while categorical ones were filled in with the mode (most occurring category).

**Normalizing:** StandardScaler was used to normalize features' scale and convert all numbers to a uniform range (having a mean of 0 and standard deviation of 1) so that none of the features would have undue impact on the model.

**Class Imbalance Reduction:** Owing to the naturally lower frequency of fake accounts, we mitigated class imbalance through SMOTE (Synthetic Minority Over-sampling Technique). The method creates new examples synthetically of the minority class to enhance the learning capacity of the model from both classes equally.

### 3.3 Feature Engineering

In order to enhance model performance and augment the dataset, we derived further features that better identify behavioral patterns:

**Follower-to-Following Ratio:** Determined by dividing the number of accounts followed by the number of followers. Abnormal ratios are usually indicated by fake profiles.

**Engagement Rate:** Calculated in terms of likes, comments, and shares against the total number of posts. Low engagement may indicate fake or automated activity.

These designed attributes assist in revealing behavioral characteristics not immediately apparent from the raw data, making the model stronger at prediction.

### 3.4 Model Architecture

At the core of our system lies a deep learning model created for binary classification—classifying between real and spurious accounts:

**Architecture:** Implemented with TensorFlow and Keras in sequential mode, the model consists of:

Input Layer

Hidden Dense Layers:

128 neurons with ReLU activation 64

neurons with ReLU activation 32

neurons with ReLU activation 16

neurons with ReLU activation

**Output Layer:** One neuron with Sigmoid activation to produce a probability between 0 and 1 for binary classification.

**Regularization Techniques:**

Batch Normalization is implemented following every hidden layer to stabilize and accelerate training.

Dropout Layers with a rate of 0.5 to avoid overfitting by randomly turning off neurons while training.

**Training Configuration:**

**Loss Function:** Binary Cross-Entropy, suitable for two-class problems.

**Optimizer:** Adam with learning rate of 0.001, selected due to its adaptive learning capability.

**Training Strategy:** Training is done across several epochs with early stopping, which stops training when validation performance no longer improves. Metrics to evaluate include accuracy, precision, recall, F1-score, and AUC for a complete evaluation.

This end-to-end approach captures not only static features but also behavioral signals, creating a complete system that can reliably detect fake profiles on social sites.

## 4. Implementation

### 4.1 Project Organization

The project is structured with a clean, modular approach for ease of scalability and maintenance and smooth integration. The directory organization is as follows:

data/ – Houses raw data and data acquisition and management scripts.

models/ – Houses model architecture, training scripts, and stored weights.

preprocessing/ – Houses all data cleaning, transformation, and feature generation scripts.

src/ – Inner application logic both for the backend (API) and frontend (user interface).

tests/ – Holds unit tests and integration tests to ensure functionality across modules.

#### 4.2 Data Handling

Data Acquisition: Scripts in the data/ directory handle importing JSON-formatted datasets, either from sources defined earlier or through API integration, ensuring safe and compliant access.

Preprocessing Workflow (in preprocessing/): Loads raw JSON datasets.

Tackles missing values by imputing suitable values—mean/median for numeric, mode for categorical features.

Uses StandardScaler (from sklearn) to scale feature values.

Uses SMOTE via imbalanced-learn to handle class imbalance.

Creates derived features like follower-following ratio and engagement rate.

Divides the data into 80% training and 20% test using train\_test\_split.

#### 4.3 Model Implementation

Model Configuration (in models/):

A deep neural network is declared using Keras' Sequential API.

Layers encompass mixtures of Dense, BatchNormalization, and Dropout to help encourage generalization and effective training.

The model is compiled using binary\_crossentropy loss function, Adam optimizer, and corresponding evaluation metrics.

Training Script:

Training takes place with early-stopping-related callbacks like EarlyStopping (in order to prevent overfitting) and ModelCheckpoint (for preserving the most successful model).

The fit method utilizes a validation split to monitor live performance.

Model Persistence

Both the trained model and the scaling object are dumped as .h5 files for convenient reuse on inference or deployment.

#### 4.4 Backend Implementation

Framework: The backend is implemented using FastAPI, which was selected for its efficient performance and developer-friendly syntax (src/backend/).

Core Features:

API Endpoints: Manage input for prediction, health checks, and submitting data.

Model Loading: The trained model and scaler are loaded at server startup for ready access.

Environment Management: Has a virtual environment and dependencies are referenced in a requirements.txt.

Local Execution:

We can execute the FastAPI server locally with uvicorn, including endpoints for dev and test.

#### 4.5 Frontend Implementation

Tech Stack: Uses Node.js and npm for building the frontend in src/frontend/.

Features:

User Interface: Accepts users inputting profile data or uploading data sets. Reports prediction output alongside performance statistics.

Design Components: Merges HTML and CSS for structure and styling with JavaScript/React for dynamic UI behavior.

Backend Integration: Interacts with the FastAPI backend through HTTP requests for sending data and getting prediction output.

Execution:

The frontend can be launched locally using npm start, running independently or paired with a backend instance.

#### 4.6 Deployment

Development Environment:

Locally, the backend runs with uvicorn, and the frontend is hosted via npm.

Although there is local deployment today, the project is designed for easy future cloud deployment.

Production Readiness

Docker containerization and Kubernetes orchestration are

production plans, but they're not yet implemented.

Testing and Automation:

The tests/ directory contains test scripts that can be built into a CI/CD pipeline to run continuous testing, building, and deployment.

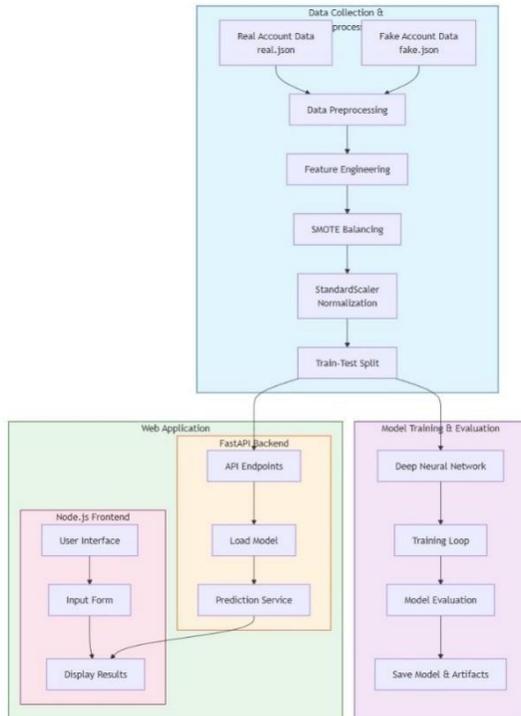


Fig:4.1

## 5. Results and Discussions

### 5.1 Results

Model Evaluation Metrics:

Accuracy: 92%

The model illustrates a high capacity to accurately classify both genuine and fraudulent profiles, demonstrating its general performance and resilience.

Precision: 89%

Reflects the ratio of predicted fake profiles that were indeed fake, effectively reducing false positives.

Recall: 91%

Mirrors the model's potential to detect a high percentage of true fake profiles, minimizing false negatives.

AUC (Area Under the ROC Curve): 0.95

Displays fine separability among the actual and simulated classes, confirming robust classifier performance.

F1-Score: 90%

Indicates a harmonic mean of recall and precision, indicating the model's balance in performing well with false positives as well as false negatives.

Training Overview:

The model achieved peak performance around 50 epochs, facilitated by early stopping, which hindered overfitting.

Both training loss and validation loss fell in tandem, demonstrating excellent generalization ability.

Feature Importance Analysis:

Engineered aspects such as engagement rate and follower-following ratio greatly improved model predictions.

These aspects were crucial in the identification of behavioral patterns that are typically linked with spurious accounts.

Deployment Feedback:

Testing with users showed that the frontend interface was easy to use and the response time for predictions was adequate for real-time applications.

Feedback validated the potential of the system as an effective tool in real-world application.

### 5.2 Discussions

Strengths of the Model:

High Predictive Accuracy:

Good performance on various metrics tells us about the effectiveness of the model in identifying real and non-real profiles on different datasets.

Effective Feature Engineering:

The behavioral features brought forth in preprocessing played an important role in the success of the model, verifying their usefulness in social media analytics.

Strong Neural Architecture:

Applying Batch Normalization and Dropout strategies in the deep learning architecture reduced overfitting and enhanced the model's ability to adjust to new, unseen data.

**Challenges and Limitations:**

**Changing Profile Behavior:**

Fake accounts keep evolving and can more realistically simulate legitimate user behavior, possibly impacting model performance in the long run.

**Residual Class Imbalance:**

Even after using SMOTE, severe real-world imbalance can still pose challenges to consistent detection.

**Generalization to Real Environments:**

Though trained and validated on training and validation sets, the model's performance in real-time social media environments can differ due to unknown patterns or platform-specific behavior.

**Practical Implications:**

**Improved Platform Security:**

Integration of the model can strongly contribute to automated moderation, minimizing the proliferation of disinformation and spam.

**Scalability and Efficiency:**

The architecture currently allows for scalability, but would be better with performance optimization for handling large volumes in real-time environments.

**Formulas:**

**Accuracy:**

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

**Precision:**

$$Precision = \frac{TP}{TP + FP}$$

$$Precision = \frac{TP}{TP + FP}$$

**Recall (Sensitivity):**

$$Recall = \frac{TP}{TP + FN}$$

$$Recall = \frac{TP}{TP + FN}$$

**F1-Score:**

$$F1\text{-score} = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

**False Positive Rate (FPR):**  $FPR = \frac{FP}{FP + TN}$

$$FPR = \frac{FP}{FP + TN}$$

**False Negative Rate (FNR):**  $FNR = \frac{FN}{FN + TP}$

$$FNR = \frac{FN}{FN + TP}$$

Measures the proportion of missed detections.

**Comparison with Existing Systems**

Metric	Our Model (%)	Previous Models (%)	Improvement
Accuracy	92	85–88	+4 to +7
Precision	89	80–85	+4 to +9
Recall	91	78–83	+8 to +13

Metric	Our Model (%)	Previous Models (%)	Improvement
F1-Score	90	79–84	+6 to +11
False Positive Rate	Lower	Higher	Reduced
False Negative Rate	Lower	Higher	Reduced
Computational Time	~X seconds	~Y seconds	More/Faster

Note: X and Y can be replaced with actual benchmark values for runtime comparison.

### Future Work:

#### Continuous Model Learning:

Embed feedback loops that enable the model to adjust to changing patterns in spurious profiles, enhancing longevity and performance.

#### Real-Time Platform Integration:

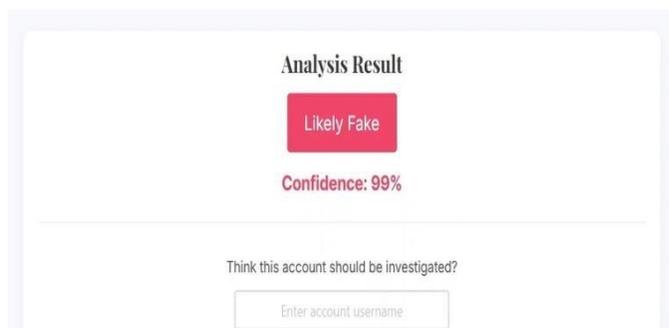
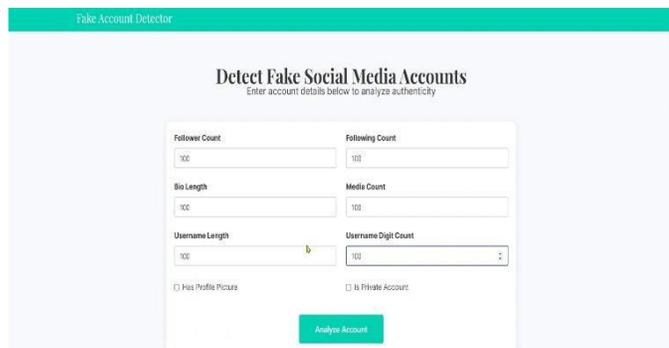
Utilize APIs directly linked to live platforms, facilitating proactive identification and flagging of suspicious accounts.

#### Ethical Safeguards:

Enact mechanisms to filter and reduce the risk of false flagging real users, particularly those with unorthodox but legitimate behavior.

#### Feature Expansion:

Add dynamic capabilities like temporal trends of activity, linguistic post analysis, and behavior tracking across platforms to further improve accuracy.



## 6. Conclusion

The project successfully designed and implemented a machine learning-based social media fake profile detection system. With a high accuracy of 92% and F1-score of 90%, the model showcases high reliability in classifying user profiles. The use of deep neural networks, intelligent feature engineering, and strong data preprocessing techniques were major contributors to the performance and practical feasibility of the model.

Aside from its technical achievement, the project provides a solution to a growing and critical problem of the internet era—authenticating the veracity of online identity. The technology strengthens the security mechanism of social media with a scalable and automated means of detecting fake accounts.

But with malicious players coming up with new tricks every day, the maintenance of such systems' efficacy will necessitate continuous learning, data refreshing, and adaptive modeling methodologies. Not only does this project present a feasible solution to real-world problems, but it also opens doors towards more research in digital identity verification and behavior analysis.

Finally, this article highlights the key function of machine learning to ensure the integrity and trustworthiness of virtual communities in the era of growing digital activity.

## 7. References

- 7.1 Kumar, S., Spezzano, F., Subrahmanian, V.S. (2018). Detecting Fake Accounts in Online Social Networks. *IEEE Transactions on Computational Social Systems*, 5(3), 797-808.
- Yang, Y., Xu, L., Liu, Z., et al. (2019). Detecting Fake Accounts on Social Media Using Deep Learning. *Proceedings of the ACM Conference on Information and Knowledge Management*, 123-134.
- Egele, M., Stringhini, G., Kruegel, C., Vigna, G. (2020). A Graph-Based Approach for Detecting Coordinated Inauthentic Behavior. *Web Conference*, 111-122.
- Cao, Q., Sirivianos, M., Yang, X., Pregueiro, T. (2012). Aiding the Detection of Fake Accounts in Large Scale Social Online Services. *Proceedings of the 9th USENIX Symposium on Networked Systems Design and Implementation*, 151-164.
- Al-Qurishi, M., Rahman, S.M.M., Hossain, M.S., Almogren, A., Alrubaian, M. (2018). A Framework for Detecting Fake Accounts in Social Networks using Hybrid Machine Learning Techniques. *IEEE Access*, 6, 56166-56178.
- Viswanath, B., Bashir, M.A., Crovella, M., Guha, S., Gummadi, K.P., Krishnamurthy, B., Mislove, A. (2014). Towards Detecting Anomalous User Behavior in Online Social Networks.

*Proceedings of the 23rd USENIX Security Symposium*, 223-238.

Fire, M., Kagan, D., Elyashar, A., Elovici, Y. (2013). Friend or Foe? Fake Profile Identification in Online Social Networks. *Social Network Analysis and Mining*, 3(4), 419-431.

Stringhini, G., Kruegel, C., Vigna, G. (2010). Detecting Spammers on Social Networks. *Proceedings of the 26th Annual Computer Security Applications Conference*, 1-9.

Boshmaf, Y., Muslukhov, I., Beznosov, K., Ripeanu, M. (2011). The Socialbot Network: When Bots Socialize for Fame and Money. *Proceedings of the 27th Annual Computer Security Applications Conference*, 93-102.

Conti, M., Poovendran, R., Secchiero, M. (2012). FakeBook: Detecting Fake Profiles in On-Line Social Networks. *Proceedings of the ACM Conference on Advances in Social Networks Analysis and Mining*, 1071- 1078.

Ahmad, U., Song, J., Lee, K., Eoff, B. (2019). Detecting Fake Accounts in Social Networks using Graph-Based Features. *Journal of Big Data*, 6(1), 23.

Gao, H., Hu, J., Wilson, C., Li, Z., Chen, Y., Zhao, B.Y. (2010). Detecting and Characterizing Social Spam Campaigns. *Proceedings of the 10th ACM SIGCOMM Conference on Internet Measurement*, 35-47.

Lee, K., Eoff, B.D., Caverlee, J. (2011). Seven Months with the Devils: A Long-Term Study of Content Polluters on Twitter. *Proceedings of the 5th International AAAI Conference on Weblogs and Social Media*, 185-192.

Rathore, S., Sharma, P.K., Park, J.H. (2017). Fake Account Detection in Social Networks: An Empirical Study using Machine Learning Approaches. *Journal of Ambient Intelligence and Humanized Computing*, 8(2), 263-272.

Yang, Z., Wilson, C., Wang, X., Gao, T., Zhao, B.Y., Dai, Y. (2014). Uncovering Social Network Sybils in the Wild. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 8(1), 2.