

AI Based Detection of Fake Profiles on Social Media Using Machine Learning and Deep Learning

Kartik Sayal, Pravleen Kaur, Ramneek Kaur
Department of Artificial Intelligence and Machine Learning
Chandigarh Group of Colleges, Jhanjeri

Abstract

With the increasing influence of social media, the proliferation of fake profiles has become a pressing concern. These accounts contribute to cyber threats, misinformation, and identity theft. This research introduces an AI-based system that leverages machine learning (ML) and deep learning (DL) algorithms to detect fake profiles on platforms like Facebook and Twitter. The system analyzes profile metadata, activity patterns, linguistic content, and network structures using a hybrid model comprising Random Forest, Support Vector Machine (SVM), and Long Short-Term Memory (LSTM) networks. Among the evaluated models, LSTM achieved the highest accuracy of 94.5%, offering a scalable and effective solution for real-time detection. Future enhancements include the integration of transformer models and adversarial learning to improve robustness.

1. Introduction

Social media platforms are increasingly being exploited by malicious entities through the creation of fake profiles. These accounts are used for phishing, spreading propaganda, cyberbullying, and financial scams. Traditional rule-based systems are inadequate due to the evolving tactics of profile creators. Thus, AI-driven approaches that combine natural language processing, behavioural analysis, and network graph techniques are essential.

2. Related Work

Traditional detection relied on rule-based systems that analysed profile completeness, posting behavior, and friend connections. However, they

lack adaptability. Recent studies advocate supervised models (e.g., Decision Trees, SVMs), unsupervised techniques (e.g., clustering, graph analysis), and deep learning frameworks such as CNNs, RNNs, and Transformers. Hybrid approaches have shown promise in overcoming the limitations of standalone models.

3. Problem Statement

Fake profiles can be categorized as bots, spammers, or impersonators and are associated with misinformation, fraud, and privacy violations. The project aims to develop a robust AI-based detection system capable of accurately identifying these profiles through a multi-modal analysis of data.

4. Methodology

4.1 Dataset

The dataset includes 10,000+ user profiles sourced from Twitter, Facebook, and Kaggle. It consists of labeled real and fake profiles characterized by features such as account age, post frequency, image use, and engagement rates.

4.2 Feature Engineering

- Profile-based: Account age, completeness, bio length.
- Behavioural: Frequency and nature of posts, engagement metrics.
- Textual: NLP techniques including TF-IDF and sentiment analysis.
- Network-based: Follower patterns, friend connections, clustering coefficients.

4.3 Models Used

- ML Models: Logistic Regression, Random Forest, Support Vector Machine.
- DL Models: LSTM for sequential text analysis, CNN for image analysis.
- Hybrid Approach: Combines outputs from ML, DL, and network analysis modules.

5. Result and Discussion

Model	Accuracy	Precision	Recall	F1-Score
Logistic Regression	85.3%	83.7%	81.5%	82.6%
Random Forest	91.2%	90.8%	89.5%	90.1%
SVM	88.7%	87.9%	85.6%	86.7%
LSTM	94.5%	93.8%	92.2%	93.0%

- The **LSTM model** outperformed other methods, especially in text analysis.
- **Confusion matrix** analysis revealed high true positive rates, indicating robustness against false negatives.
- The hybrid model showed promising adaptability across various datasets.

6. Conclusion

The project demonstrates the feasibility of using AI to detect fake profiles on social media with high accuracy. LSTM models, enriched by feature engineering and NLP, offer superior performance. The system effectively classifies accounts by analyzing user behavior, content, and connectivity.

7. Future Work

- Integrate transformer-based models like BERT and GPT for enhanced contextual understanding.

- Employ Graph Neural Networks (GNNs) for sophisticated relationship modeling.
- Explore real-time deployment using APIs.
- Adopt adversarial learning to counter evolving tactics.
- Ensure ethical compliance and privacy preservation through data anonymization.

References

- Aggarwal, C. C. (2015). *Data Mining: The Textbook*. Springer.
- Al-Qurishi, M., Al-Rakhami, M., & Al-Rakhami, M. S. (2019). Detection of Fake Accounts in Social Media: A Survey. *IEEE Access*, 7, 21276–21293. <https://doi.org/10.1109/ACCESS.2019.2891504>
- Chowdhury, S. R., Alghamdi, M. A., & Khandoker, S. (2020). A Hybrid Machine Learning Approach for Fake Profile Detection. *Journal of Information Security and Applications*, 54, 102526. <https://doi.org/10.1016/j.jisa.2020.10252>
- Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press.
- Kaggle. (2023). *Fake and Real Profile Dataset*. <https://www.kaggle.com>
- Miller, Z. (2021). Adversarial Machine Learning for Social Media Security. *ACM Transactions on Cybersecurity*, 9(3), 45–61. <https://doi.org/10.1145/3447781>
- Rathore, S., et al. (2019). Social Media Security: Machine Learning Approaches to Fake Profile Detection. *Future Generation Computer Systems*, 96, 579–593. <https://doi.org/10.1016/j.future.2017.06.031>
- Sebastian, R., & Patil, P. (2021). AI-Powered Fake Account Detection on Twitter Using NLP. *International Journal of Computer Science & Information Technology*, 13(2), 112–128. <https://doi.org/10.5121/ijcsit.2021.13209>
- Scikit-learn Developers. (2023). *Scikit-learn: Machine Learning in Python*. <https://scikit-learn.org>

- Vaswani, A., Shazeer, N., Parmar, N., et al. (2017). Attention Is All You Need. *Advances in Neural Information Processing Systems (NeurIPS)*. <https://arxiv.org/abs/1706.03762>
- Zhang, C., et al. (2020). A Survey on Fake News Detection: Data, Methods, and Challenges. *ACM Computing Surveys (CSUR)*, 53(5), 1–40. <https://doi.org/10.1145/3395046>
- Ferrara, E. (2017). Disinformation and Social Bot Operations in the Run Up to the 2017 French Presidential Election. *First Monday*, 22(8). <https://doi.org/10.5210/fm.v22i8.8005>
- Cresci, S., et al. (2015). Fame for Sale: Efficient Detection of Fake Twitter Followers. *Decision Support Systems*, 80, 56–71. <https://doi.org/10.1016/j.dss.2015.09.003>
- Wu, L., et al. (2020). Graph-based Social Bot Detection. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 11(1), 1–31. <https://doi.org/10.1145/3375627>
- Monti, F., et al. (2019). Fake News Detection on Social Media Using Geometric Deep Learning. *arXiv preprint arXiv:1902.06673*. <https://arxiv.org/abs/1902.06673>
- Rudin, C. (2019). Stop Explaining Black Box Machine Learning Models for High Stakes Decisions and Use Interpretable Models Instead. *Nature Machine Intelligence*, 1(5), 206–215. <https://doi.org/10.1038/s42256-019-0048-x>
- Zhou, X., & Zafarani, R. (2019). A Survey of Fake News: Fundamental Theories, Detection Methods, and Opportunities. *ACM Computing Surveys*, 53(5), 1–40. <https://doi.org/10.1145/3395046>
- Tang, J., et al. (2015). LINE: Large-scale Information Network Embedding. *Proceedings of the 24th International Conference on World Wide Web (WWW)*, 1067–1077. <https://doi.org/10.1145/2736277.274109>
- Chen, E., et al. (2022). How Suspicious Are You? Towards Characterizing Suspicious Users in Social Media. *The Web Conference 2022*, 431–442. <https://doi.org/10.1145/3485447.351199>