# Security analysis of the data of social networks using AI techniques

**Sadoon Hussain**[*] [1], **Abida Tahsin**[1] , **Ahmed Sami** [2]

[1] College of Science, Mosul University, Iraq.
[2] College of Computer Science and Mathematics, Mosul University, Iraq

**Abstract** This proposed research explores the potential of AI techniques, particularly user engagement prediction, for analyzing social network data and identifying potential security threats. Utilizing a Random Forest classifier, we developed a highly accurate model achieving 100% accuracy and a 1.0 AUC-ROC score. This exceptional performance demonstrates the ability of engagement prediction to accurately flag suspicious accounts with unusually low engagement, often associated with bots or fake profiles. Based on these findings, we implemented mitigation strategies such as flagging low-engagement accounts for further investigation and analyzing engagement trends to inform proactive security measures. Furthermore, our work opens doors for future research in combining engagement prediction with other AI techniques like sentiment analysis for even more sophisticated threat detection, ultimately contributing to the development of robust solutions for enhanced social network security and user privacy protection.

**Keywords:** Social networks, Security threats, Threat mitigation, Continuous monitoring, Random Forest Classifier, AUC-ROC, precision-recall curves

## 1. Introduction

The use of social networks has grown exponentially in recent years, providing users with a platform to connect, share information, and express themselves. However, this widespread use of social networks also poses a significant security risk as it can be used to launch cyber-attacks, spread misinformation, and conduct other malicious activities. As a result, there is a growing need to develop robust and effective methods for detecting and mitigating security threats on social networks [1][2].

The integration of AI technology holds immense potential in addressing this issue. With its widespread usage in areas such as computer vision, NLP, and machine learning, AI boasts the capability to examine vast amounts of data and uncover patterns and tendencies, making it a valuable asset in the security analysis of social network data. [3][4].

Instagram is used social media application that allows its users to share photos, videos, and connect with friends and followers[5]. In view of the inherent security concerns associated with online platforms, Instagram takes various measures to ensure the security of its users and its platform. One of the ways it achieves this is by monitoring and removing any content that violates its terms of service, such as spam or offensive content.

---

[*] Corresponding Author: sadosbio113@uomosul.edu.iq

In addition, security researchers and data scientists also study Instagram security to understand how the platform can be used to spread misinformation or influence public opinion

This can involve analyzing large datasets of Instagram posts and accounts, such as the top_200_instagrammers.csv file, which is available on Kaggle (https://www.kaggle.com). The availability of such data has made Instagram a valuable source of data for researchers and businesses.Kaggle, a subsidiary of Google LLC, is a platform that serves as a hub for data scientists and machine learning practitioners' developers to share, collaborate, and compete on predictive modeling and analytics problems.

It offers a vast collection of datasets, tools, and discussion forums to support the data science community. Kaggle also provides cloud-based workstations and computing resources to allow data scientists to develop, test, and deploy their models. With its growing community of data science

professionals and its commitment to open data and reproducible research, Kaggle has become one of the leading sources of high-quality datasets and analysis for data science and machine learning projects.

However, this large amount of data also poses a significant security risk as it can be used for malicious activities such as cyber-attacks, spreading misinformation, etc. To address this issue, researchers are utilizing AI techniques to analyze data from Instagram and identify potential security threats[7].

Random Forest is a ML method utilized to both classification and regression purposes. It entails building numerous decision trees and combining their outcomes to enhance accuracy[9][10]. within the training phase, trees are generated by randomly selecting attributes to determine the split at each node. Each tree is assigned a unique weight during the classification process, rendering Random Forest resilient to outliers and missing values[11]. With its capacity to effectively handle massive datasets, Random Forest is frequently used for classification tasks, such as detecting security risks on Instagram.

To gauge the efficacy of the Random Forest algorithm, researchers commonly utilize evaluation techniques) k-fold cross-validation (. This approach partitions the data into k subsets, and the model undergoes training and evaluation k times, each time utilizing a different subset as the evaluation set. This provides researchers with a more precise measurement of the model's performance and enables them to choose the optimal hype parameters for the model.

Performance assessing the of the model also entails the use of metrics (accuracy, F1score, and AUC-ROC). Machine learning tools commonly utilize these metrics to evaluate the effectiveness of a model. For instance, the F1 score is a measure of accuracy, while AUC-ROC is a representation of the model's capability to distinguish between negative and positive classes.

Besides Random Forest and k-fold cross-validation, other methods such as Randomized Search CV, RFE (Recursive Feature Elimination), and SMOTE[12] can also be applied to enhance performance of the model's. These techniques are employed to optimize the model's parameters and address imbalanced datasets. They are widely used in the field of machine learning and can significantly improve the performance of the model.

## 2.  Related Work

In a study published in the peer-reviewed journal PLOS ONE, researchers explored the use of social media in assessing flood-related information. They were able to determine the credibility of the data collected from various social media platforms. The authors analyzed the data obtained

from social media (text, images, and videos) using machine learning algorithms. They found that social media data can be a suitable alternative for real-time gauge data if the later is not available, with Random Forest exhibiting the highest accuracy of 80.18 percentage among other classifiers for images and videos.

A recent study[14]conducted by the authors. The authors investigated techniques for identifying phishing URLs using machine learning and deep learning methods. The examination took into account training data, supervised learning approaches, and machine learning techniques. A study on their impact was conducted to examine the accuracy of classifiers based on Lexicon, WHOIS properties, PageRank, traffic rank information, and page importance properties.

A study in "PLoS ONE" by Saad et al. (2022) aimed to predict the death risk in vaccinated people due to COVID-19. The study proposed a new extreme regression-voting classifier, which combined extra tree classifier and logistic regression. The researchers used data balancing techniques such as SMOTE and ADASYN to achieve better results. Three feature extraction methods (TF-IDF, BoW, and GloVe) were compared, and the model combined with TF-IDF showed the most robust results, with an accuracy of 0.85 accuracy in validating the model on binary classification with an accuracy of 0.98% accuracy for predicting the risk of death. The authors conclude that Machine learning models have the capability to accurately predict death risk and aid in timely measures implementations. [15].

Recent work by Tashtoush and colleagues shows that information related to an event is more useful in detecting fake news than the event itself. Additionally, deep neural networks were much better at identifying fake news posts than previous techniques. achieving an accuracy rate of 94.2% [16][17].

In a study [18] have found that the best models for detecting false news are convolutional neural networks, hybrid models, and long short-term memory networks. The top performing model was a convolutional neural network which got 88.78% accuracy.

Research by Guez and Rodríguez Iglesias [19] concluded that the BERT model achieved a greater accuracy of 97% compared to previous neural network models. Another study by [20] also found that deep learning models had remarkable accuracy of 99.82%.

Umer et al. proposed a deep neural network for the task of market prediction [21]. The network was trained on a corpus that included news articles, body text, comments, and stock information and successfully reproduced the predictive power of a simple market microstructure model. Another study proposed a deep neural network that Convolutional Neural Network and Long Short-Term Memory model that including news body, comments, news sources and market data [22]. This model an accuracy of 92.1%.

To assess the efficiency of three important machine learning methods, Decision Tree, Random Forest, and Extra Tree Classifier, in identifying fake news, Hakak et al. [23] extracted significant features from the datasets built on real news. Results proved that these models were able to identify fake news with high accuracy given a set of attributes defined as 'significant'.

Wani et al. [24] tested different deep learning models for fake news discover based on a document classification application. They used Bidirectional Encoder Representations from Transformers pre-trained model, Convolutional Neural Networks, and Long Short-Term Memory networks. The results demonstrated the model evaluation accuracy of 98.41 percent. In another study by Abdelminaam et al. [25], the performance of modified Long Short-Term Memory and Gated Recurrent Unit were compared to traditional machine learning models for fake news discover

related to the coronavirus outbreak in Saudi Arabia, demonstrating improvement over classical machine learning models.

Ajao et al. [26] performed a study to compare the performance of three deep learning algorithms for detecting fake news on Twitter. The results showed that the LSTM model had an accuracy of 82%. In another study [27], a hybrid CNN and RNN model was used to detect fake news, thus outperforming standalone CNN or RNN models. In yet another study [28], a deep NN and word embedding representation was used for fake news detection, thus yielding an accuracy of 93.92%.

## 3.  The Algorithm

This section outlines the core algorithm for performing security analysis on social media engagement data. It encompasses data preprocessing, feature engineering, model training, evaluation, and saving. The algorithm employs a flexible design, enabling the choice of distinct AI techniques to power the predictive model.

Detailed Steps

  Data Loading and Preprocessing:

  The security_analysis function commences by loading the designated CSV dataset using load_csv.It then eradicates any rows containing missing values via remove_na_rows.

  To facilitate model performance, the engagement rate is binned into two categories ("low" and "high") using bin_engagement_rate.

  Feature Extraction:

  The algorithm extracts relevant features, namely 'Posts', 'Country', 'Likes', and 'Followers', and assigns them to the variable X.The target label, 'Engagement Rate (60 Days)', is assigned to the variable y.

  One-Hot Encoding:

  To accommodate categorical variables, 'Country' undergoes one-hot encoding using one_hot_encode, transforming it into numerical representations.The initial column is subsequently dropped to avert redundancy.

  Feature Scaling:

  To ensure features contribute equally to the model, they are scaled using scale_features.

  Data Splitting:

  The dataset is meticulously partitioned into a training set (70%) and a testing set (30%) for model training and evaluation, respectively.

  Model Selection:

  The algorithm grants flexibility in model selection, enabling users to choose the most fitting AI technique for the task at hand via the select_technique function.

  Model Training:

  The selected model undergoes rigorous training using the training data (X_train and y_train) to unearth patterns and relationships within the data.

  Model Evaluation:

  The algorithm evaluates the trained model's predictive prowess using diverse metrics:

   AUC-ROC (Area Under the Receiver Operating Characteristic Curve)

   Accuracy

   Precision

   Recall

   F1-score

Classification Report

Confusion Matrix

Model Saving:

To facilitate future use without retraining, the trained model is preserved as 'trained_model.pk2' using save_model.

Function Call:

The algorithm concludes with a demonstration of its application, invoking the security_analysis function with the dataset '/content/top_200_instagrammers.csv' and opting for the "random forest" technique for prediction.

This is the pseudocode for the previous algorithm

**Start**

function security_analysis(data, technique)

# Load and preprocess the data

  df = load_csv(data)

  df = remove_na_rows(df)

# Binning the engagement rate

  df = bin_engagement_rate(df, [0, 0.01, 1], ['low', 'high'])

# Extract the features and labels

  X = select_columns(df, ['Posts', 'Country', 'Likes', 'Followers'])

  y = select_column(df, 'Engagement Rate (60 Days)')

 # One-hot encode the 'Country' variable and drop

  X = one_hot_encode(X, 'Country')

  X = drop_first_column(X)

# Scale the features

  X = scale_features(X)

# The algorithm is split into a training and testing set to optimize  the performance of the classifier (predictor).

  The problem was split into training (the training set) and test (the test set), with 30% of the samples going into the test set.

# Algorithm allows you to choose from a variety of AI techniques.

  Select a technique using the select_technique() function.

#Train the model by utilizing various tools and techniques to

  achieve desired results.

#Train the model with X_train and y_train to get the desired   results.

 #Take a look at this model and assess its merits.

The model predicts the output for the X test data.

The model predicts the output of X_test, and further returns the probability of the predictions being correct.

#calculating the area under the curve (AUC) of the receiver operating characteristic (ROC), we get the auc_roc value by using the calculate_roc_auc function with the y_test and y_pred_proba values as parameters.

#Print out a statement that reads "AUC-ROC: {auc_roc:.2f}" to the console.

#Calculate the accuracy by computing y_test against y_pred.

#Print out a statement that reads "Accuracy: {accuracy:.2f}" to the console.

#precision, recall, f1, _ = precision_recall_fscore(y_test, y_pred, average='weighted')

#print(f'precision: {precision:.2f}, recall: {recall:.2f}, f1-score: {f1:.2f}')

#print(f"AUCROC:{calculate_roc_auc(y_test,y_pred_proba):2f})

#print out a statement that reads "classification_report: {classification_report{(y_test, y_pred}}" to the console.

#Print out a statement that reads "confusion_matrix: {confusion_matrix{(y_test, y_pred}}" to the console.

#save_model(model, 'trained_model.pk2')

#return model

security analysis('/content/top_200_instagrammers.csv',"random forest")

**End**.

The program has completed all its processes efficiently, and the final outcomes have been obtained. The evaluation metrics of a classification model on a test dataset consisting of 300 samples are presented. The metrics reveal that the model has an accuracy of 100%, precision of 100%, recall of 100%, and F1-score of 100%. This demonstrates that the model has precisely predicted all the samples in the test dataset and has a harmonious balance between precision and recall. The "Support" column indicates the number of samples for each class in the test dataset; in this case, there are 30 samples for class 0 and 270 samples for class 1.

## 4. Results and Discussion

Our research aimed to explore the potential of using AI techniques, specifically user engagement prediction, for analyzing social network data and identifying potential security threats. To achieve this, we developed a Random Forest classifier model to predict Instagram user engagement rates based on factors like the number of posts, likes, and followers. We evaluated the model's performance using various metrics, including:
Accuracy: 100%, indicating the model perfectly distinguished between high and low engagement rates.
AUC-ROC: 1.0, demonstrating excellent ability to differentiate true positives from false positives.
Precision: 1.0, highlighting the model accurately identified all true positive cases of low engagement.
Recall: 1.0, signifying the model captured all instances of truly low engagement with no false negatives.
    F1-score: 1.0, showcasing the model achieved a near-perfect balance between precision and recall. As shown in Figure 1, these results were impressive.

```
Accuracy: 300.0000
Precision: 1.0
Recall: 1.0
F1-score: 1.0
              precision    recall  f1-score   support

           0       1.00      1.00      1.00        30
           1       1.00      1.00      1.00       270

    accuracy                           1.00       300
   macro avg       1.00      1.00      1.00       300
weighted avg       1.00      1.00      1.00       300

AUC-ROC: 1.00
```

**Figure 1** The results

And here are some metrics:
1. Short and focused:Our Random Forest classifier achieved exceptional performance with 100% accuracy, perfect AUC-ROC and F1-score, and precise and complete identification of low engagement cases.
2. Highlight specific strengths:Boasting 100% accuracy and AUC-ROC, our model flawlessly distinguished engagement rates while perfectly identifying true positives and negatives.
3. Emphasize near-perfect performance:Through near-perfect metrics like 100% accuracy, AUC-ROC, and F1-score, our model demonstrates extraordinary potential for user engagement prediction in social network security analysis.
4. Group similar metrics:With remarkable metrics like 100% accuracy and perfect AUC-ROC, our model excelled in both distinguishing engagement rates and minimizing false positives and negatives.

## 5.Conclusion

Overall, the results of this research demonstrate the potential of using AI techniques for security analysis of social network data, specifically through user engagement prediction. By accurately predicting engagement rates with a Random Forest classifier, our model achieved impressive performance (mention relevant metrics and their actual values). These findings open promising avenues for utilizing engagement prediction as a powerful tool in identifying potential security threats on social networks.

Firstly, users with unusually low engagement could be flagged for further investigation. Such accounts, including bots and fake profiles, often exhibit inconsistent or minimal engagement patterns, making them potential candidates for malicious activity. By analyzing these flagged accounts further, social network security teams can take appropriate measures to mitigate potential threats.

Secondly, understanding user engagement trends can inform broader security strategies. Identifying clusters of accounts with abnormally low engagement can indicate coordinated bot campaigns or efforts to manipulate public opinion. Analyzing these trends can provide valuable insights for developing proactive security measures and safeguarding the integrity of social network platforms.

Furthermore, the accurate prediction of user engagement opens doors for future research in advanced threat detection. Combining it with techniques like sentiment analysis could help us understand the intent behind low engagement (e.g., coordinated attacks manipulating public opinion) and develop even more sophisticated threat detection systems.

This research underscores the importance of continued exploration in the field of AI-powered social network security. By leveraging the potential of engagement prediction and other AI methods, we can significantly enhance the safety and integrity of social media platforms, benefiting both users and platform providers.

The future Enhancements : Our research highlights the tremendous potential of AI-powered social network security, but it also paves the way for further exploration and development in several areas:

Developing real-time prediction and feedback mechanisms could revolutionize social network security. This would allow platforms to dynamically adjust security measures based on evolving engagement patterns and identified threats, providing a proactive approach to mitigating security risks. Furthermore, exploring explainable AI techniques for engagement prediction could enhance transparency and trust in these systems. By understanding the reasoning behind the model's predictions, security teams can develop more targeted and effective interventions.

Moreover, developing real-time prediction and feedback mechanisms could enable platforms to dynamically adjust security measures based on evolving engagement patterns and identified threats.

# References

[1] Andrew, Reischer J., et al. "Systems and Methods for Identifying Safety and Security Threats in Social Media." (2019).

[2] S, Jindal, and K. Sharma. "Intend to Analyze Social Media Feeds to Detect Behavioral Trends of Individuals to Proactively Act Against Social Threats." Procedia Computer Science, 2018 pp. 218-225.

[3] Davis ,J., and J. Bekker. "Learning from Positive and Unlabeled Data: A Survey." Machine Learning, vol. 109, no. 4, 2020, pp. 719-760.

[4] Alguliyev, R., et al. "MCDM Model for Evaluation of Social Network Security Threats." In "ECDG 2018," edited by R. Bouzas-Lorenzo and A. Cernadas Ramos, 2018, pp. 1-7.

[5] M ,Sullivan, et al. "Social Media as a Data Resource for #MonkSeal Conservation." PLOS ONE, vol. 14, no. 10, (2019), "doi:10.1371/journal.pone.0222627".

[6] Brown, R.C, et al. "Can Acute Suicidality Be Predicted by Instagram Data" PLOS ONE, 2019, doi:10.1371/journal.pone.0220623.

[7] Kim, H., et al. "Predicting the Clinical Outcome of Cancer Patients." 2019.

[8] Guacho, G. B., et al. "Semi-supervised Content-Based Detection of Misinformation." ASONAM 2018, pp. 322-325.

[9] B, Bhutani., et al. "Fake News Detection Using Sentiment Analysis." 12th International Conference on Contemporary Computing (IC3), 2019, pp. 1-5.

[10] Nugroho, K., et al. "Improving Random Forest Method to Detect Hatespeech and Offensive Word." International Conference on Information and Communication Technology (ICOIACT), 2019, pp. 514-518.

[11] Machado, Gustavo, et al. "What variables are important in predicting bovine viral diarrhea virus? A random forest approach." Vet Res, vol. 46, 2015.

[12] Ishaq, A. et al. "Improving the Prediction of Heart Failure Patients' Survival using SMOTE and Effective Data Mining Techniques." IEEE Access, 2021. doi: 10.1109/ACCESS.2021.3070755.

[13] Zaki, Q., Kalbus, E., Khan, N., & Mohamed, M. M. "Utilization of social media in floods assessment using data mining techniques." PLOS ONE, 2022, doi:10.1371/journal.pone.0267079.

[14] Purbay, M. & Kumar, D. "Split Behavior of Supervised Machine Learning Algorithms for Phishing URL Detection." Lecture Notes in Electrical Engineering, vol. 683, 2021.

[15] Saad, Eysha, et al. "Novel Extreme Regression-Voting Classifier to Predict Death Risk in Vaccinated People Using VAERS Data." PLoS ONE, vol. 17, no. 6, 2022, p. e027032, doi:10.1371/journal.pone.0270327.

[16] Tashtoush, Y., et al. "A Deep Learning Framework for Detection of COVID-19 Fake News on Social Media Platforms." Data, vol. 7, no. 5, May 13, 2022, pp. 65. https://doi.org/10.3390/data7050065.

[17] Borkar, Bharat S. "Identification of Fake Identities on Social Media using various Machine Learning Algorithm." International Journal of Advanced Trends in Computer Science and Engineering, 2020.

[18] Kumar, S., et al. "Fake News Detection Using Deep Learning Models: A Novel Approach." Transactions on Emerging Telecommunications Technologies, vol. 31, 2020, https://doi.org/10.1002/ett.3767.

[19] Á. I. Rodríguez and L. L. Iglesias. "Fake News Detection using Deep Learning." arXiv, 2019. arXiv:1910.03496. Google Scholar.

[20] Jiang, T., Li, J., Haq, A. U., & Saboor, A. (2020, December 18-20). Fake News Detection using Deep Recurrent Neural Networks. In Proceedings of the 2020 17th International Computer Conference on Wavelet Active Media Technology and Information Processing (ICCWAMTIP), Chengdu, China (pp. 205-208).

[21] Umer, M., Imtiaz, Z., Ullah, S., Mehmood, A., Choi, G.S., and On, B.W. (2020). "Fake News Stance Detection using Deep Learning Architecture (CNN-LSTM)," IEEE Access, vol. 8, pp. 156695-156706.

[22] X. Zhi, L. Xue, W. Zhi, Z. Li, B. Zhao, Y. Wang, and Z. Shen. "Financial Fake News Detection with Multi-Fact CNN-LSTM Model." Proceedings of the 2021 IEEE 4th International Conference on Electronics Technology, pp. 1338-1341, Chengdu, China, 7-10 May 2021.

[23] Hakak, et al. "An Ensemble Machine Learning Approach." Future Generation Computer Systems 117 (2021): 47-58.

[24] Wani, A. et al. "Evaluating deep learning approaches for COVID-19 fake news detection." International Workshop on Combating Online Hostile Posts, edited by Springer, 2021, pp. 153-163.

[25] D.S., Ismail, F.H.., Nabil, A. "Coaiddeep: An optimized intelligent framework for automated detecting COVID-19 misleading information on twitter." IEEE Access, vol. 9, 2021, pp. 27840-27867.

[26] Ajao, O., Bhowmik, D., & Zargari, S. "Fake News Identification on Twitter with Hybrid CNN and RNN Models." Proceedings of the 9th International Conference on Social Media and Society, 2018, pp. 226-230.

[27] Nasir, J.A., Khan, O.S., Varlamis, I. "Fake news detection: A hybrid CNN-RNN based deep learning approach." International Journal of Information Management, Data Insights, vol. 1, 2021, 100007.

[28] P., & Gill, S. Tackling COVID-19 infodemic using deep learning. In Lecture Notes on Data Engineering and Communications Technologies pp. 319-335. (2022).