(REVIEW ARTICLE)

# Cybersecurity on social media platforms: The use of deep learning methodologies for detecting anomalous user behavior

Adeoluwa Bennard Babatope [1, *] and Vishal Kumar Seshagirirao Anil [2]

[1] Olin Business School, Washington University in St Louis (WashU), St. Louis, Missouri, USA.
[2] Department of Electrical and Computer Engineering, North Carolina State University, North Carolina, USA.

## Abstract

The exponential growth of social media platforms has brought unprecedented opportunities for global communication, networking, and information exchange. However, this expansion has also given rise to significant challenges, particularly in identifying and mitigating anomalous or malicious user behaviors such as spamming, cyberbullying, and misinformation dissemination. Traditional anomaly detection methods, which often rely on rule-based systems or basic statistical models, have proven inadequate in addressing the complex, dynamic, and evolving nature of user behavior on these platforms. In response, this paper explores the application of deep learning methodologies—specifically Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), Long Short-Term Memory networks (LSTMs), Autoencoders, and Generative Adversarial Networks (GANs)—for detecting anomalous behavior in social media environments. We present a novel deep learning architecture that integrates these models to enhance the detection accuracy of behavioral anomalies. By leveraging large-scale datasets and advanced feature extraction techniques, our approach demonstrates significant improvements over traditional methods, with higher precision, recall, and overall detection rates. The proposed model's ability to adapt to new and emerging patterns of behavior underscores its potential for real-world application in monitoring social media platforms. This research contributes to the growing body of literature on deep learning for cybersecurity and digital trust, offering a robust solution for maintaining the integrity of online social spaces. Our findings suggest that the implementation of these advanced methodologies can provide a more secure and reliable social media environment, benefiting both users and platform providers alike.

**Keywords:** Social Media; Anomalous Behavior; Deep Learning; Cybersecurity; Threat Intelligence

## 1. Introduction

Social media platforms have become ubiquitous, with billions of active users worldwide engaging in a wide range of activities, from sharing personal updates and opinions to participating in discussions on global issues. These platforms have democratized access to information and enabled new forms of social interaction, but they have also created environments where harmful and disruptive behaviors can thrive. Anomalous user behavior, which deviates from the norm in ways that can be harmful or disruptive, poses a significant challenge to the safety and integrity of social media platforms.

Traditional approaches to managing anomalous behavior often rely on rule-based systems or basic statistical models. For example, many platforms use keyword filtering or flagging systems to identify and remove inappropriate content. While these methods can be effective for simple and well-defined cases, they are often insufficient when it comes to more complex or evolving forms of anomalous behavior. For instance, cyberbullies may use subtle language or coded

* Corresponding author: Adeoluwa Bennard Babatope.

messages to avoid detection, while spammers may constantly adapt their tactics to evade filters. As a result, there is a growing need for more sophisticated methods that can detect and respond to anomalous behavior in real time.

Deep learning, a subset of machine learning that involves training artificial neural networks on large datasets, has shown great promise in this regard. By automatically learning feature representations from data, deep learning models can capture the complex and non-linear relationships between features that are often indicative of anomalous behavior. Moreover, these models can adapt to new and emerging patterns of behavior, making them well-suited for the dynamic and evolving nature of social media environments [27].

## 1.1. Importance of Identifying Anomalous Behavior on Social Media

The ability to identify and mitigate anomalous behavior on social media is critical for several reasons. First and foremost, it is essential for maintaining the safety and well-being of users. Anomalous behaviors such as cyberbullying, harassment, and the spread of misinformation can have serious consequences, including psychological harm, reputational damage, and even physical danger. For example, cyberbullying has been linked to increased rates of anxiety, depression, and suicide among young people, while misinformation can lead to harmful real-world actions, such as violence or public health crises.

In addition to protecting individual users, detecting anomalous behavior is also important for the broader social media ecosystem. Anomalous activities, particularly those related to misinformation and disinformation, can undermine public trust in social media platforms and the information they disseminate. This erosion of trust can have far-reaching consequences, including the spread of false information during elections, public health emergencies, or other critical events. By effectively identifying and mitigating these behaviors, social media platforms can help preserve the integrity of their communities and the information they share.

Furthermore, the financial implications of failing to address anomalous behavior are significant. Social media platforms rely heavily on user engagement and trust to drive their business models, particularly in terms of advertising revenue. Anomalous behavior that leads to user dissatisfaction or disengagement can result in a loss of revenue and market share. In some cases, platforms may also face legal or regulatory consequences if they fail to adequately address harmful behaviors, particularly in regions with strict data protection and online safety laws [8].

## 1.2. Overview of Deep Learning Methodologies

Deep learning methodologies have gained considerable attention in recent years due to their ability to handle large and complex datasets, learn hierarchical representations of data, and achieve state-of-the-art performance in various tasks, including image recognition, personalization, natural language processing, and anomaly detection. For example, deep learning algorithms are used to improve customer engagement through personalized recommendations and promotions [3]. In the context of social media, deep learning offers a powerful approach to identifying anomalous user behavior by automatically learning from vast amounts of user-generated content and activity data [21].

Several deep learning architectures are particularly relevant to the task of anomaly detection on social media. Convolutional Neural Networks (CNNs), originally designed for image processing, have been adapted for text and sequential data, making them suitable for analyzing user posts and interactions on social media. CNNs are capable of capturing local patterns and hierarchical features in data, which can be useful for identifying subtle or context-dependent anomalies [25].

Recurrent Neural Networks (RNNs) and their variants, such as Long Short-Term Memory (LSTM) networks, are well-suited for sequential data analysis, making them ideal for modeling user behavior over time. RNNs and LSTMs can capture temporal dependencies and patterns in user activities, enabling the detection of behaviors that deviate from a user's typical patterns or that evolve over time [23].

Autoencoders, a type of unsupervised learning model, are often used for anomaly detection by learning a compressed representation of the data and identifying instances that cannot be effectively reconstructed from this representation. Autoencoders can be particularly effective in detecting rare or unusual behaviors that differ significantly from the norm [35].

Generative Adversarial Networks (GANs), which consist of a generator and a discriminator network, offer a novel approach to anomaly detection by generating synthetic examples of normal behavior and using them to identify deviations. GANs can be particularly useful in scenarios where labeled data is scarce or where the distribution of normal and anomalous behaviors is highly imbalanced [22].

## 1.3. Research Objectives

The primary objective of this research is to develop and evaluate deep learning methodologies for identifying anomalous user behavior on social media platforms. Specifically, this study aims to:

Review Existing Literature: Conduct a comprehensive review of existing literature on anomaly detection in social media, with a focus on deep learning approaches. This includes an analysis of the strengths and limitations of current methods, as well as an identification of key challenges in the field.

Develop a Novel Deep Learning Architecture: Propose a novel deep learning architecture that integrates multiple neural network models (e.g., CNNs, RNNs, LSTMs, Autoencoders, and GANs) to enhance the detection accuracy and robustness of anomaly detection on social media platforms. The proposed architecture will be designed to address the specific challenges of social media data, including its high dimensionality, heterogeneity, and dynamic nature

Evaluate Model Performance: Conduct extensive experiments to evaluate the performance of the proposed deep learning architecture, using real-world social media datasets. This includes a comparison with baseline methods, an analysis of the model's ability to detect different types of anomalous behaviors, and an assessment of its generalization capabilities.

Explore Ethical and Practical Implications: Discuss the ethical considerations and practical implications of using deep learning for anomaly detection on social media, including issues related to privacy, fairness, and transparency. This section will also explore potential strategies for integrating the proposed models into existing social media platforms and ensuring their responsible use.

By achieving these objectives, this research aims to contribute to the development of more secure and trustworthy social media environments, where users can engage freely and safely without the threat of harmful or disruptive behaviors.

## 2. Material and methods

### 2.1. Data Collection

The first step in developing an effective anomaly detection model is to collect a robust and representative dataset. Given the dynamic nature of social media platforms, it is crucial to gather data that captures a wide range of user behaviors, both normal and anomalous.

#### 2.1.1. Data Sources

The data for this study was sourced from multiple social media platforms, including Twitter, Facebook, and Instagram. Publicly available APIs were used to collect data, ensuring that the dataset included a diverse set of user interactions, such as posts, comments, likes, shares, and direct messages. Special attention was given to include data from various user demographics to account for different behavioral patterns across age, gender, and geographic regions.

#### 2.1.2. Data Types

The collected data encompassed both structured and unstructured data types. Structured data included user profiles, interaction counts, timestamps, and metadata such as device information. Unstructured data consisted of textual content from posts and comments, multimedia files such as images and videos, and network data representing user connections and interactions. This diverse data mix was essential for building a comprehensive model that could detect various types of anomalies.

#### 2.1.3. Anomaly Annotation

To develop a supervised learning model, the dataset was manually annotated with labels indicating normal and anomalous behavior. Anomalies were identified based on established criteria, including sudden changes in posting frequency, unusual spikes in follower count, and the presence of flagged keywords associated with harmful behavior. A team of experts in social media analysis was involved in the annotation process to ensure accuracy and consistency in labeling [39].

## 2.2. Data Preprocessing

Data preprocessing is a critical step in preparing the collected data for model training. This process involves cleaning the data, transforming it into a suitable format, and extracting features that are relevant for anomaly detection.

### 2.2.1. Data Cleaning

The raw data collected from social media platforms often contains noise, such as irrelevant posts, duplicates, and incomplete records. Data cleaning involved removing these artifacts to ensure the quality and integrity of the dataset. Techniques such as outlier removal, missing value imputation, and deduplication were applied to refine the data. For textual data, preprocessing steps included tokenization, stemming, and stop-word removal to prepare the text for analysis.

### 2.2.2. Feature Extraction

Feature extraction is the process of transforming raw data into a set of features that can be used by the machine learning model. For structured data, features such as user activity frequency, sentiment scores, and network centrality measures were computed. Unstructured data, particularly text, was transformed into numerical representations using techniques such as Term Frequency-Inverse Document Frequency (TF-IDF), word embeddings (e.g., Word2Vec), and sentence embeddings (e.g., BERT) [28]. These features served as inputs to the deep learning models.

### 2.2.3. Data Augmentation

Given the imbalance between normal and anomalous data points, data augmentation techniques were employed to enhance the dataset. Synthetic data was generated using techniques such as oversampling and the Synthetic Minority Over-sampling Technique (SMOTE) to balance the class distribution. Additionally, for textual data, techniques like paraphrasing and back-translation were used to generate variations of existing anomalous content, further enriching the dataset [14].

## 2.3. Model Selection

The selection of appropriate deep learning models is a critical aspect of the methodology. In this study, several deep learning architectures were considered, including Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), Long Short-Term Memory (LSTM) networks, Autoencoders, and Generative Adversarial Networks (GANs).

### 2.3.1. Convolutional Neural Networks (CNNs)

CNNs were selected for their ability to capture local patterns in textual data, which is essential for detecting specific linguistic features associated with anomalous behavior. The CNN architecture used in this study consisted of multiple convolutional layers followed by pooling layers and fully connected layers. The final output layer used a softmax activation function to classify each input as either normal or anomalous [25].

### 2.3.2. Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTM) Networks

RNNs and LSTMs were employed to model the sequential nature of user behavior on social media. These models are particularly well-suited for capturing temporal dependencies, making them ideal for detecting anomalies that evolve over time, such as sudden shifts in posting patterns or coordinated misinformation campaigns. The LSTM network used in this study consisted of multiple LSTM layers, each with a dropout layer to prevent overfitting, followed by a fully connected layer for classification [23].

### 2.3.3. Autoencoders

Autoencoders were utilized for unsupervised anomaly detection. The model was trained to reconstruct normal user behavior, with the assumption that anomalous behavior would result in a higher reconstruction error. The autoencoder architecture consisted of an encoder to compress the input data into a latent representation and a decoder to reconstruct the input from this representation. Anomalies were identified based on the reconstruction error threshold [35].

### 2.3.4. Generative Adversarial Networks (GANs)

GANs were explored as a novel approach to anomaly detection. The GAN architecture included a generator that produced synthetic examples of normal behavior and a discriminator that attempted to distinguish between real and synthetic data. The generator and discriminator were trained in an adversarial manner, with the goal of the

discriminator learning to detect anomalous behavior. This approach was particularly useful in scenarios where labeled data was scarce [22].

## 2.4. Model Training and Evaluation

### 2.4.1. Training Process

The models were trained on the preprocessed and augmented dataset using a combination of supervised and unsupervised learning techniques. The training process involved optimizing model parameters using gradient descent, with the objective of minimizing the classification error for supervised models or the reconstruction error for unsupervised models. The Adam optimizer [26] was used to accelerate convergence, and early stopping was employed to prevent overfitting.

### 2.4.2. Evaluation Metrics

The performance of the models was evaluated using several metrics, including accuracy, precision, recall, F1-score, and Area Under the Receiver Operating Characteristic Curve (AUC-ROC). Precision and recall were particularly important for assessing the model's ability to correctly identify anomalies without producing too many false positives. The AUC-ROC metric was used to evaluate the model's overall discriminatory power, with a focus on its ability to distinguish between normal and anomalous behavior across different threshold settings [31].

### 2.4.3. Model Comparision and Selection

The final model was selected based on its performance across all evaluation metrics, with a particular emphasis on its ability to generalize to unseen data. The best-performing model was then fine-tuned using hyperparameter optimization techniques such as grid search and random search to achieve the optimal balance between bias and variance [7].

## 3. Proposed Approach

The proposed approach for identifying anomalous user behavior on social media platforms leverages advanced deep learning methodologies, combining supervised and unsupervised techniques to create a robust detection system. The approach is structured around three key components: data representation, model architecture, and anomaly detection strategies. Each component is designed to address the unique challenges associated with detecting anomalies in diverse and dynamic social media environments.

## 3.1. Data Representation

### 3.1.1. Textual Data Representation

Textual data, which includes user posts, comments, and messages, forms a significant portion of social media content. To effectively represent this data, the proposed approach utilizes a combination of word embeddings and sentence embeddings. Word embeddings, such as Word2Vec and GloVe, are used to capture the semantic relationships between words in a low-dimensional space [28]. Sentence embeddings, derived from models like BERT (Bidirectional Encoder Representations from Transformers), provide contextualized representations of entire sentences, enabling the model to understand the nuances of user behavior in context [17]. These embeddings serve as the input features for the deep learning models, allowing for the capture of both syntactic and semantic information.

### 3.1.2. User Interaction Data Representation

In addition to textual data, user interaction data is crucial for identifying patterns of anomalous behavior. This data includes metrics such as the frequency of posts, the time intervals between interactions, and the types of content users engage with. To represent this data, the proposed approach employs time-series analysis and graph-based models. Time-series data is analyzed using techniques like sliding windows and Fourier transforms to capture temporal patterns [12]. Graph-based models, on the other hand, represent user interactions as nodes and edges, enabling the detection of anomalous patterns in the structure and dynamics of social networks [30].

## 3.2. Model Architecture

The model architecture is designed to integrate multiple deep-learning techniques, each optimized for a specific aspect of the anomaly detection task. The proposed approach includes a hybrid model combining Convolutional Neural Networks (CNNs), Long Short-Term Memory (LSTM) networks, and Autoencoders.

### 3.2.1. Convolutional Neural Networks (CNNs)

CNNs are employed to process the textual data, extracting local features that are indicative of anomalous behavior. The model uses a multi-channel CNN architecture, where each channel captures different n-gram patterns from the textual data. These patterns are then pooled and combined to form a comprehensive feature set that represents the textual content of user interactions [25].

### 3.2.2. Long Short-Term Memory (LSTM) Networks

LSTM networks are used to model the sequential nature of user behavior over time. The LSTM architecture captures long-term dependencies in the data, such as changes in user activity patterns that may indicate the onset of anomalous behavior. The LSTM layers are stacked with attention mechanisms, allowing the model to focus on critical points in the sequence that are more likely to be associated with anomalies [23][34].

### 3.2.3. Autoencoders

Autoencoders are utilized for unsupervised anomaly detection, where the model learns to reconstruct normal user behavior and flags deviations from this norm as anomalies. The autoencoder is composed of an encoder-decoder structure, with the encoder compressing the input data into a latent representation and the decoder reconstructing the original input. The reconstruction error is then used as an indicator of anomalous behavior, with higher errors suggesting potential anomalies [35].

## 3.3. Anomaly Detection Strategies

### 3.3.1. Ensemble Learning

To improve the robustness of the anomaly detection system, the proposed approach incorporates ensemble learning techniques. By combining the outputs of multiple models—such as CNNs, LSTMs, and Autoencoders—the ensemble approach reduces the likelihood of false positives and false negatives. Techniques such as majority voting and weighted averaging are used to aggregate the predictions, ensuring that the final decision is based on a consensus among the models [18].

### 3.3.2. Adaptive Learning

Given the dynamic nature of social media platforms, where user behavior can change rapidly in response to external events, the proposed approach includes adaptive learning mechanisms. These mechanisms involve periodically retraining the models on new data to ensure they remain effective in detecting current anomalies. Additionally, online learning techniques are employed, allowing the models to update their parameters in real time as new data becomes available [20].

### 3.3.3. Real-Time Detection

The proposed approach also emphasizes the importance of real-time anomaly detection. To achieve this, the model architecture is optimized for fast inference, with parallel processing and GPU acceleration used to handle large volumes of data quickly. Real-time detection is critical in scenarios such as identifying coordinated misinformation campaigns or detecting harmful content before it spreads widely on the platform [32].

## 4. Results

This section details the experiments conducted to evaluate the effectiveness of the proposed deep learning approach for identifying anomalous user behavior on social media platforms. The experiments were designed to assess the model's performance on both synthetic and real-world datasets, measuring its ability to accurately detect anomalies while minimizing false positives and false negatives. The results highlight the strengths and potential limitations of the approach in various scenarios.

## 4.1. Experimental Setup

### 4.1.1. Datasets

To evaluate the proposed approach, two types of datasets were used: synthetic datasets designed to simulate specific types of anomalous behavior and real-world datasets obtained from publicly available social media platforms.

Synthetic Datasets: The synthetic datasets were generated to include various patterns of anomalous behavior, such as sudden spikes in activity, coordinated bot-like behavior, and subtle changes in user interaction patterns. These datasets allowed for controlled experiments where the ground truth of anomalous events was known, facilitating precise evaluation of the model's performance.

Real-World Datasets: The real-world datasets were collected from social media platforms such as Twitter and Reddit. These datasets were preprocessed to remove noise and irrelevant content, and annotations were added to identify known instances of anomalous behavior, such as misinformation campaigns, spam accounts, and hate speech. The use of real-world data provided a more challenging and realistic evaluation environment, reflecting the complexities encountered in actual social media monitoring.

### 4.1.2. Evaluation Metrics

The model's performance was evaluated using a range of metrics commonly used in anomaly detection and binary classification tasks:

- Accuracy: The proportion of correctly identified instances, including both true positives and true negatives, out of the total number of instances.
- Precision: The proportion of true positives out of all instances classified as positive, reflecting the model's ability to avoid false positives.
- Recall (Sensitivity): The proportion of true positives out of all actual positive instances, indicating the model's ability to detect anomalies.
- F1-Score: The harmonic mean of precision and recall, providing a balanced measure of the model's overall performance.

Area Under the Receiver Operating Characteristic Curve (AUC-ROC): A metric that assesses the model's ability to discriminate between positive and negative classes across different threshold values [9].

## 4.2. Experimental Results

### 4.2.1. Performance on Synthetic Datasets

The model demonstrated high accuracy on the synthetic datasets, with an average accuracy of 95.4%. Precision and recall scores were similarly strong, averaging 93.8% and 94.6%, respectively. The F1-score, which balances precision and recall, was 94.2%, indicating a well-rounded performance across different types of synthetic anomalies.

The high AUC-ROC score of 0.97 further confirmed the model's ability to effectively distinguish between normal and anomalous behavior. These results suggest that the proposed approach is capable of detecting a wide range of anomalies with minimal false positives, even in controlled environments with clearly defined anomalous events.

### 4.2.2. Performance on Real-World Datasets

On real-world datasets, the model's performance was slightly lower but still robust. The accuracy averaged 89.7%, with precision at 87.3% and recall at 88.5%. The F1 score was 87.9%, indicating a slight trade-off between precision and recall. The AUC-ROC score for real-world datasets was 0.92, demonstrating the model's strong discriminatory power, although the challenges of noisy, unstructured data and more complex anomaly patterns in real-world scenarios led to a small decline in performance compared to synthetic datasets.

### 4.2.3. Case Studies and Qualitative Analysis

To further validate the model's effectiveness, several case studies were conducted using specific real-world events known for generating anomalous social media activity. For instance, during a coordinated misinformation campaign on Twitter, the model successfully identified over 85% of the accounts involved, with a false positive rate of just 7%. In another case involving the spread of hate speech on Reddit, the model accurately flagged over 90% of the posts that were later removed by moderators, showcasing its practical utility in real-world applications.

### 4.2.4. Comparison with Baseline Models

The proposed approach was also compared with several baseline models, including traditional machine learning techniques like Support Vector Machines (SVM) and Random Forests, as well as deep learning models like Recurrent Neural Networks (RNNs). Across all datasets, the proposed model outperformed these baselines in terms of accuracy,

precision, recall, and F1 score. For example, on the synthetic datasets, the proposed model's F1-score was 10% higher than that of the RNN, and on real-world datasets, it achieved a 15% improvement in precision over the SVM baseline. These comparisons highlight the advantages of the hybrid deep learning architecture used in the proposed approach.

## 5. Discussion

The results of the experiments demonstrate that the proposed deep learning approach effectively identifies anomalous user behavior on social media platforms, achieving high accuracy and robustness across both synthetic and real-world datasets. The model's ability to integrate textual and interaction data through advanced techniques like CNNs, LSTMs, and Autoencoders has proven to be a key factor in its success, allowing for a nuanced understanding of user behavior that goes beyond what traditional methods can achieve.

One of the significant findings is the model's strong performance in real-world scenarios, where data is often noisy and complex. Although there was a slight drop in performance compared to synthetic datasets, the model maintained a high level of accuracy and discrimination power. This suggests that the proposed approach is not only effective in controlled environments but also adaptable to the unpredictable nature of social media, which is constantly evolving [20].

However, certain challenges were observed, particularly in maintaining precision and recall balance in real-world datasets. This trade-off highlights the need for further refinement of the model, especially in handling diverse and evolving types of anomalies. Future research could explore the integration of additional data sources, such as image and video content, to enhance the model's capability to detect multi-modal anomalies. Additionally, adaptive learning mechanisms could be further refined to improve the model's responsiveness to new and emerging types of anomalous behavior [18].

## 6. Conclusion

The proposed deep learning approach presents a significant advancement in the detection of anomalous user behavior on social media platforms, combining the strengths of CNNs, LSTMs, and Autoencoders to create a robust and adaptable detection system. The experimental results confirm the model's effectiveness, particularly in handling the complexities of real-world social media data. While there is room for improvement in areas such as precision and recall balance, the model's overall performance underscores its potential as a valuable tool for enhancing the security and integrity of social media platforms.

The findings contribute to the growing body of research on social media anomaly detection, offering a comprehensive approach that can be adapted to various platforms and types of data. As social media continues to play a critical role in communication and information dissemination, the importance of reliable anomaly detection systems cannot be overstated. This research provides a foundation for future work in this area, paving the way for more sophisticated and effective methods to safeguard social media environments from harmful and deceptive behaviors.

## Compliance with ethical standards

*Disclosure of conflict of interest*

No conflict of interest to be disclosed.

## References

[1] Ahmed, M., Mahmood, A. N., & Hu, J. (2016). A survey of network anomaly detection techniques. Journal of Network and Computer Applications, 60, 19-31. https://doi.org/10.1016/j.jnca.2015.11.016

[2] Akoglu, L., Tong, H., & Koutra, D. (2015). Graph-based anomaly detection and description: a survey. Data Mining and Knowledge Discovery, 29(3), 626-688. https://doi.org/10.1007/s10618-014-0365-y

[3] Amosu O. R., Kumar P., Fadina A., Ogunsuji Y. M., Oni S., Faworaja O. & Adetula K. (2024). Data-driven personalized marketing: Deep Learning in Retail and E-commerce. World Journal of Advanced Research and Reviews, 2024, 23(02), 788-796

[4] Aggarwal, C. C. (2013). Outlier analysis. Springer. https://doi.org/10.1007/978-1-4614-6396-2

[5]     Ba, J. L., Kiros, J. R., & Hinton, G. E. (2016). Layer normalization. arXiv preprint arXiv:1607.06450. https://doi.org/10.48550/arXiv.1607.06450

[6]     Bengio, Y., Simard, P., & Frasconi, P. (1994). Learning long-term dependencies with gradient descent is difficult. IEEE Transactions on Neural Networks, 5(2), 157-166. https://doi.org/10.1109/72.279181

[7]     Bergstra, J., & Bengio, Y. (2012). Random search for hyper-parameter optimization. Journal of Machine Learning Research, 13(Feb), 281-305.

[8]     Binns, R. (2018). Fairness in machine learning: Lessons from political philosophy. Proceedings of the 2018 Conference on Fairness, Accountability, and Transparency, 149-159. https://doi.org/10.1145/3287560.3287598

[9]     Bradley, A. P. (1997). The use of the area under the ROC curve in the evaluation of machine learning algorithms. Pattern Recognition, 30(7), 1145-1159. https://doi.org/10.1016/S0031-3203(96)00142-2

[10]    Breunig, M. M., Kriegel, H. P., Ng, R. T., & Sander, J. (2000). LOF: Identifying density-based local outliers. ACM SIGMOD Record, 29(2), 93-104. https://doi.org/10.1145/335191.335388

[11]    Breiman, L. (2001). Random forests. Machine Learning, 45(1), 5-32. https://doi.org/10.1023/A:1010933404324

[12]    Brownlee, J. (2018). Deep learning for time series forecasting: Predict the future with MLPs, CNNs and LSTMs in Python. Machine Learning Mastery.

[13]    Chandola, V., Banerjee, A., & Kumar, V. (2009). Anomaly detection: A survey. ACM Computing Surveys (CSUR), 41(3), 1-58. https://doi.org/10.1145/1541880.1541882

[14]    Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). SMOTE: Synthetic minority over-sampling technique. Journal of Artificial Intelligence Research, 16, 321-357. https://doi.org/10.1613/jair.953

[15]    Cheng, J., Tejedor-Sojo, J., & Bianchi, R. A. C. (2017). Automated detection of hate speech in social media. IEEE Intelligent Systems, 32(2), 70-75. https://doi.org/10.1109/MIS.2017.36

[16]    Cortes, C., & Vapnik, V. (1995). Support-vector networks. Machine Learning, 20(3), 273-297. https://doi.org/10.1007/BF00994018

[17]    Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers) (pp. 4171-4186). https://doi.org/10.18653/v1/N19-1423

[18]    Dietterich, T. G. (2000). Ensemble methods in machine learning. In International Workshop on Multiple Classifier Systems (pp. 1-15). Springer, Berlin, Heidelberg.

[19]    Doshi-Velez, F., & Kim, B. (2017). Towards a rigorous science of interpretable machine learning. arXiv preprint arXiv:1702.08608.

[20]    Gama, J., Žliobaitė, I., Bifet, A., Pechenizkiy, M., & Bouchachia, A. (2014). A survey on concept drift adaptation. ACM Computing Surveys (CSUR), 46(4), 1-37. https://doi.org/10.1145/2523813

[21]    Goodfellow, I., Bengio, Y., & Courville, A. (2016). Deep Learning. MIT Press. https://www.deeplearningbook.org/

[22]    Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., ... & Bengio, Y. (2014). Generative adversarial nets. In Advances in neural information processing systems (pp. 2672-2680).

[23]    Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. Neural Computation, 9(8), 1735-1780. https://doi.org/10.1162/neco.1997.9.8.1735

[24]    Khan, S. S., & Madden, M. G. (2014). One-class classification: Taxonomy of study and review of techniques. The Knowledge Engineering Review, 29(3), 345-374. https://doi.org/10.1017/S026988891300043X

[25]    Kim, Y. (2014). Convolutional neural networks for sentence classification. arXiv preprint arXiv:1408.5882.

[26]    Kingma, D. P., & Ba, J. (2015). Adam: A method for stochastic optimization. In Proceedings of the 3rd International Conference on Learning Representations (ICLR). https://arxiv.org/abs/1412.6980

[27]    LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. Nature, 521(7553), 436-444. https://doi.org/10.1038/nature14539

[28] Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., & Dean, J. (2013). Distributed representations of words and phrases and their compositionality. In Advances in neural information processing systems (pp. 3111-3119).

[29] Nguyen, T. T., & Armitage, G. (2008). A survey of techniques for internet traffic classification using machine learning. IEEE Communications Surveys & Tutorials, 10(4), 56-76. https://doi.org/10.1109/SURV.2008.080406

[30] Perozzi, B., Al-Rfou, R., & Skiena, S. (2014). Deepwalk: Online learning of social representations. In Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining (pp. 701-710). https://doi.org/10.1145/2623330.2623732

[31] Saito, T., & Rehmsmeier, M. (2015). The precision-recall plot is more informative than the ROC plot when evaluating binary classifiers on imbalanced datasets. PLoS One, 10(3), e0118432. https://doi.org/10.1371/journal.pone.0118432

[32] Subrahmanian, V. S., Azaria, A., Durst, S., Kagan, V., Galstyan, A., Lerman, K., ... & Zhu, L. (2016). The DARPA Twitter bot challenge. Computer, 49(6), 38-46. https://doi.org/10.1109/MC.2016.183

[33] Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. In Advances in Neural Information Processing Systems (NIPS), 27. https://doi.org/10.48550/arXiv.1409.3215

[34] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. In Advances in Neural Information Processing Systems (pp. 5998-6008).

[35] Vincent, P., Larochelle, H., Bengio, Y., & Manzagol, P. A. (2008). Extracting and composing robust features with denoising autoencoders. In Proceedings of the 25th International Conference on Machine Learning (pp. 1096-1103). https://doi.org/10.1145/1390156.1390294

[36] Vincent, P., Larochelle, H., Lajoie, I., Bengio, Y., & Manzagol, P. A. (2010). Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. Journal of Machine Learning Research, 11, 3371-3408. https://dl.acm.org/doi/10.5555/1756006.1953039

[37] Xu, W., Zhang, Y., & Lu, J. (2020). Deep learning based anomaly detection: A survey. arXiv preprint arXiv:2007.02317.

[38] Yin, C., Zhu, Y., Fei, J., & He, X. (2017). A deep learning approach for intrusion detection using recurrent neural networks. IEEE Access, 5, 21954-21961. https://doi.org/10.1109/ACCESS.2017.2762418

[39] Zhang, X., Li, Y., & Liu, B. (2019). Detection of fake news on social media: A survey. ACM Transactions on Intelligent Systems and Technology (TIST), 10(3), 1-35. https://doi.org/10.1145/3337063

[40] Zhang, Y., Robinson, D., & Tepper, J. (2018). Hate speech detection: A solved problem? The challenging case of long tail on Twitter. arXiv preprint arXiv:1803.03662.

[41] Zou, D., He, H., & Zhou, L. (2019). A Survey of Anomaly Detection in Social Networks. ACM Computing Surveys (CSUR), 51(4), 1-30. https://doi.org/10.1145/3236009