

Leveraging Machine Learning for Fraudulent Social Media Profile Detection

Soorya Ramdas, Agnes Neenu N. T.

Department of Computer Science, St. Albert's College, India

E-mails: sooryaramdas19@gmail.com agnesneenu1@gmail.com

Abstract: Fake social media profiles are responsible for various cyber-attacks, spreading fake news, identity theft, business and payment fraud, abuse, and more. This paper aims to explore the potential of Machine Learning in detecting fake social media profiles by employing various Machine Learning algorithms, including the Dummy Classifier, Support Vector Classifier (SVC), Support Vector Classifier (SVC) kernels, Random Forest classifier, Random Forest Regressor, Decision Tree Classifier, Decision Tree Regressor, MultiLayer Perceptron classifier (MLP), MultiLayer Perceptron (MLP) Regressor, Naïve Bayes classifier, and Logistic Regression. For a comprehensive evaluation of the performance and accuracy of different models in detecting fake social media profiles, it is essential to consider confusion matrices, sampling techniques, and various metric calculations. Additionally, incorporating extended computations such as root mean squared error, mean absolute error, mean squared error and cross-validation accuracy can further enhance the overall performance of the models.

Keywords: Dummy Classifier, SVC Classifier, Random Forest classifier, Decision Tree Classifier, MLP (MultiLayer Perceptron) Classifier.

1. Introduction

In today's tech-driven era, digitalization is pervasive. Social media platforms like Facebook, Twitter, and Instagram have become integral to daily life, with a multitude of users engaging in diverse activities, fostering connections, and contributing to both positive and negative aspects [1, 2]. This paper aims to investigate the effectiveness of Machine Learning in identifying fake social media profiles. It explores multiple Machine Learning approaches and evaluates their performance using metrics such as confusion matrices, sampling techniques, and extended computations like root mean squared error, mean absolute error, mean squared error and cross-validation accuracy. As responsible users, caution is crucial to prevent harm, considering the presence of potential fraudsters [3, 4]. Unaware and trusting users on social media are susceptible to cyber-attacks, with crimes like fraud, abuse, phishing, and identity theft perpetrated through fake profiles. The rapid expansion of social media

intensifies the proliferation of such malicious accounts, posing a serious threat to platform security and integrity [5].

Addressing the challenges posed by the widespread presence of fake accounts on social media requires a careful evaluation of machine learning algorithms [6, 7]. The “No Free Lunch” theorem in machine learning implies that there is no universal algorithm excelling in all situations. Hence, the selection of algorithms should be customized to the unique attributes and intricacies of the specific challenge. The structure of this paper involves a range of machine learning algorithms, acknowledging the importance of a nuanced and context-specific approach to guarantee precise and dependable results. Analyzing insights into their limitations and strengths, along with metrics such as accuracy, precision, recall, F1 score, and confusion matrices, can provide a summary of correct and incorrect predictions. The algorithms used in this approach include the Dummy Classifier, MultiLayer Perceptron (MLP) Classifier and Regressor, Logistic Regression, Naïve Bayes Classifier, Decision Tree Classifier and Regressor, Random Forest Classifier and Regressor, Support Vector Classifier (SVC), and Support Vector Classifier Kernel functions. Additionally, expanded computations like mean absolute error, mean squared error, root mean squared error, and cross-validation calculations are performed to estimate the models’ capabilities. Various resampling techniques such as the Synthetic Minority Over sampling TEchnique (SMOTE), ADActive SYNthetic sampling technique (ADASYN), Random Oversampling, and Under-sampling are employed to address dataset class imbalance. This approach aims to present a method to mitigate risks and ensure the security of such platforms by effectively detecting fake profiles through a combination of these Machine Learning techniques.

2. Literature review

Using diverse algorithms is crucial to address the varied features and evolving strategies of fake profiles. Each algorithm has distinct strengths and limitations in detecting specific patterns and behaviors associated with counterfeit profiles. Considering some existing drawbacks, it is evident that traditional or manual identification approaches Regular Expression and Deterministic Finite Automaton approaches, and graph-based methods are inefficient, time-consuming, and subjective, making it impractical to address the escalating number of users [8, 9]. Therefore, leveraging large volumes of data, data-driven methods can be employed to utilize machine-learning algorithms and develop accurate, robust, and reliable techniques. These methods can be validated by comparing them with existing models [10]. Many existing papers concentrate on a specific algorithm with consistently low accuracy. In contrast, this paper emphasizes leveraging commonly used algorithms. One paper proposes an efficient framework for automatically detecting fake profiles using the Random Forest Classifier, achieving 95% accuracy. This study underscores the importance of the Random Forest Classifier Algorithm in addressing the issue of fake profiles. The framework introduces an automated solution for detecting fake profiles, particularly emphasizing the Random Forest Classifier for classification.

This automated approach is practical for online social networks managing a high volume of profiles, making manual examination impractical. The presented framework exhibits an impressive efficacy rate of around 95% in accurately identifying fake profiles using the Random Forest Classifier [11]. Another paper emphasizes the importance of choosing the right algorithm to enhance accuracy in detecting fake social media accounts. Despite prior use of traditional methods, there is a recognized need for precision improvement. The study employs Machine Learning and Natural Language Processing, leveraging their advantages to analyze data patterns linked to fake accounts. Specifically, the researchers opt for the Random Forest tree classifier algorithm, known for its robustness in handling diverse data types and improving accuracy in identifying fake accounts on social media [12].

In a paper is suggested that to ensure dependable predictions about profile authenticity; the research aims to assess the effectiveness of three supervised machine-learning algorithms: Random Forest (RF), Decision Tree (DT-J48), and Naïve Bayes (NB). This assessment aims to offer valuable insights, helping identify the most suitable algorithm(s) to develop resilient systems and address the pervasive issue of fake profiles and associated risks on social media platforms [13]. One of the studies employs various machine learning algorithms, including Naive Bayes, Logistic Regression, Support Vector Machine-Kernel, K-Nearest Neighbor, Boosted Tree, Neural Networks, and Logistic Regression Kernel to analyze datasets containing both fake and legitimate accounts from Facebook and Instagram. The algorithms' effectiveness is evaluated based on their classification accuracy in identifying fake profiles, with SVM achieving the highest accuracy at 97.1% on Fake profile detection datasets. By leveraging machine-learning techniques like SVM, the study demonstrates the effectiveness of these methods in distinguishing between fake and genuine profiles, contributing to the enhancement of security and privacy on Online Social Networks (OSNs) [14].

One of the studies uses the friend-to-follower ratio, a critical attribute readily available on social media profiles, to employ machine learning techniques for detecting fake profiles. With a focus on supervised and unsupervised classifiers such as Naïve Bayes, decision trees, SVM, ANNs, and NLP, the research employs targeted feature sets encompassing attributes like name, chat history, location, friends list, followers, likes, comments, and tagging. Through extensive research, the study achieves substantial improvements in detection accuracy, ranging from 50% to 96%. The careful selection and application of diverse machine learning algorithms significantly contribute to the study's reliable and effective detection of fake profiles [15]. Choosing and employing various Machine Learning (ML) algorithms, such as Support Vector Machine (SVM), Logistic Regression (LR), Random Forest (RF), and k-Nearest Neighbors (KNN), is crucial for assessing their effectiveness in detecting fake Twitter accounts and bots. Each algorithm's unique characteristics allow for a comprehensive evaluation of its performance. One of the papers employs two normalization techniques, Z-Score and Min-Max, to enhance detection accuracy by ensuring uniformly scaled features and preventing bias. Research findings highlight that Random Forest and k-Nearest Neighbors algorithms achieve high accuracy and true positive rates in detecting fake profiles, emphasizing their effectiveness in

addressing the challenge of identifying fraudulent social media accounts. The study's use of diverse ML algorithms and normalization techniques significantly contributes to exploring machine learning's potential for detecting fake profiles, offering valuable insights for accurate results. [16].

In another research [17] is shown that the algorithm plays a crucial role, serving as the cornerstone for detecting misleading fake profiles. Using machine-learning techniques, the study aims to identify and mitigate the risks associated with fake profiles. The dataset undergoes preprocessing with various Python libraries to ensure effective analysis. A comparison model is then used to select the most suitable algorithm based on dataset characteristics. Multiple machine learning algorithms, including Random Forest, Neural Network, and Support Vector Machines, have been evaluated for their efficacy in identifying fake accounts. By assessing their classification performance, researchers determine the most effective algorithm to detect counterfeit profiles, facilitating improved identification on social media platforms [17]. Another paper utilizes Deep Neural Networking and Machine Learning algorithms, including Artificial Neural Networks (ANN), Random Forest, and Support Vector Machine (SVM). The selection is based on their established effectiveness in classification tasks and compatibility with the dataset sourced from GitHub, specifically the Facebook profile Dataset. This dataset, designed for distinguishing between genuine and fake profiles, served as the basis for applying ANN, Random Forest, and SVM algorithms. The results have highlighted SVM's superiority in accuracy, making it a suitable choice for detecting fake profiles in this context [18]. In a relevant paper, the strategic choice of machine learning algorithms is paramount for discerning between counterfeit and legitimate Twitter profiles based on various characteristics. Utilizing algorithms such as neural networks, LSTM (Long Short-Term Memory), XG Boost, and Random Forest, the study effectively classifies genuine Twitter accounts as TFP (True Follower Profile) and E13 (Echobot 13/bot accounts). Simultaneously, it accurately identifies fake accounts as INT (Intentional – Accounts created with the explicit purpose of manipulation), TWT (Tweeterbot – Automated Accounts to perform certain actions), and FSF (Followers Selling Followers). This classification process significantly contributes to determining the authenticity of social media pages. Additionally, the paper delves into the architecture and hyperparameters of the selected algorithms, emphasizing the specific setups used to train models for optimal effectiveness. Post-training, the models yield outcomes where a value of 0 denotes a genuine profile, and a value of 1 indicates a fraudulent profile [19]. The authors of another study have discovered that traditional supervised learning methods may struggle with dynamically changing bot behavior. To address this, they suggest using machine learning techniques, specifically unsupervised learning, to identify fake profiles on Instagram. Utilizing unsupervised learning, the algorithm analyzes a dataset with seventeen metadata features from both genuine and fraudulent accounts, crucial for distinguishing between real and fake profiles. By scrutinizing data for patterns and anomalies, the algorithm can identify common characteristics among fake profiles, allowing it to adapt and detect fraudulent behavior even as it evolves over time [20].

Research paper, machine-learning algorithms are chosen based on their effectiveness in classification tasks. Logistic Regression and Random Forest are pivotal in analyzing Instagram's distinctive features and verifying user account authenticity. Logistic Regression, a widely used algorithm for binary classification, is trained on labeled data to learn patterns and make predictions for unseen profiles. It has been chosen for its ability to establish correlations between different features and the likelihood of an account being fake or genuine. Random Forest, an ensemble-learning algorithm using multiple decision trees, enhances prediction accuracy, especially in handling high-dimensional feature spaces and capturing intricate relationships among Instagram account attributes. Through aggregating predictions from multiple trees, Random Forest delivers robust outcomes in identifying fake profiles [21]. In a paper, AI algorithms have been leveraged to boost the precision and efficiency of the detection process, enhancing the overall success of the study. The researchers have employed decision trees, logistic regression, and support vector machines specifically for detecting fake records. A performance comparison has shown that logistic regression outperforms the other methods, highlighting its significance in accurately identifying fraudulent social media profiles [22]. Another research introduces the Support Vector Machine and Neural Network algorithm, a novel approach for efficiently detecting fraudulent Twitter accounts and bots. By incorporating four techniques for selecting relevant features and reducing dimensions, the algorithm enhances accuracy and efficiency in the detection process. In comparison to other algorithms, it shows promising outcomes, accurately classifying approximately 98% of accounts in the training dataset while utilizing a reduced set of features [23]. The study on fake profile detection emphasizes the significance of selecting appropriate machine learning algorithms, such as Logistic Regression. Logistic Regression operates as a statistical model predicting the probability of a binary outcome based on input variables. By training the algorithm on a dataset containing labelled examples of both genuine and fake profiles, it learns underlying patterns and relationships for classification. This enables automated and scalable identification of fake profiles on social media platforms. Utilizing Logistic Regression and other machine learning algorithms is crucial, allowing researchers to develop precise and efficient methods for detecting and mitigating the risks posed by fake profiles [3].

3. Data source

The data source used for this study has been obtained from **www.github.com**, specifically the datasets `fusers.csv` and `users.csv` [19]. These datasets are vital for investigating challenges related to fake social media profiles, serving as the foundation for validation. They play a crucial role in testing machine learning algorithms, forming the fundamental components for training, testing, and evaluating models in the detection of fake social media profiles [21]. These datasets consist of various attributes such as: "Id" which is for user identification; "name" represents users display name; "screen name" represents the user's unique social media handle; "statuses count" represents Number of statuses posted by the user; "followers count"

represents the number of followers the user has; “friends count” represents the number of accounts the user is following; “favourites count” represents the number of posts the user has marked as favorites; “listed count” represents number of public lists that include the user; “url” represents the user’s provided URL in the profile; “lang” represents the user’s language preference; “time zone” represents the user’s time zone setting; “location” represents the user’s specified location; “default profile” is the Indicator if the user has the default profile; “default profile image” indicates if the user has the default profile image; “geo enable” indicates if the user has enabled geo tagging; “profile image url” represents the Uniform Resource Locator of the user’s profile image; “profile banner url” represents the Uniform Resource Locator of the user’s profile banner; “profile use background image” indicates if the user uses a background image; “profile background image url https” represents the Uniform Resource Locator of the user’s background image (Hypertext Transfer Protocol Secured); “profile text color” represents the color of the text in the user’s profile; “profile image url https” represents the Uniform Resource Locator of the user’s profile image (Hypertext Transfer Protocol Secured); “profile sidebar border color” represents the border color of the user’s sidebar; “profile background tile” Indicates if the user’s background image is tiled; “profile sidebar fill color” represents the fill color of the user’s sidebar; “profile background image url” represents the Uniform Resource Locator of the user’s background image; “profile background color” represents the background color of the user’s profile; “profile link color” represents the color of links in the user’s profile; “utc offset” represents the UTC offset of the user’s time zone; “is translator” indicates if the user is a translator; “follow request sent” indicates if a follow request has been sent; “protected” indicates if the user has a protected account; “verified” indicates if the user’s account is verified; “notifications” indicates if the user has notifications enabled; “description” represents the user’s profile description/bio; “contributors enabled” Indicates if contributors are enabled for the user; “following” indicates if the authenticated user is following this user; “created at” represents the date when the user’s account was created; “timestamp” represents the timestamp associated with the data; “crawled at” represents the timestamp indicating when the data was crawled; “updated” represents the timestamp indicating when the data was last updated. The users.csv dataset, comprising 3475 rows and 42 columns, represents genuine user profiles, offering valuable insights into real user behaviors on social media. This realistic data is crucial for training and validating machine-learning models to accurately identify legitimate accounts. Conversely, the fusers.csv dataset, with 3352 rows and 38 columns, contains information about fake user profiles and plays a vital role in training models to detect fraudulent activities. Both datasets are not human-annotated because all of the entries are taken in automated ways. By analyzing features like unusual behavior and suspicious interaction patterns in fake profiles, machine-learning models learn common indicators of fraudulent accounts [16]. These datasets serve as a baseline for distinguishing between fake and genuine user profiles, allowing researchers to assess model accuracy and performance. Both datasets significantly contribute to evaluating machine-learning models, aiding researchers in selecting effective approaches. Users.csv contains genuine profiles, and fusers.csv has data on fake profiles, enabling

[illegible]

4. Methodology

This step is vital for cleaning the dataset, and removing missing values, irrelevant data, and inconsistencies. It ensures dataset compatibility for training and analysis using Machine Learning algorithms, forming the foundation for accurate fake profile detection in social media. Various data preprocessing techniques, including [13, 16] are four.

- Dealing with missing values in the dataset is crucial as they can lead to inaccurate or biased results and hinder the proper functioning of algorithms. Missing values are managed by replacing them with a default value, often changing all “Not a Number” (NaN) values to zero. The choice of zero replacement depends on the objectives and nature of the data analysis. In this study, since the dataset is mostly numeric in nature, replacing missing values with zero is a reasonable approach with a reasonable impact [19].

4.1.1. Attribute selection

In this study, a subset of features has been selected based on their correlation with the target variable [13]. There are many attributes with NaN and with values not falling under the correlation range. The following features have been chosen for further analysis: “statuses count”, “followers count”, “friends count”, “favorites count”, “listed count”, “url”, “time zone”, “follow request”, “request sent”, “id”, “location”, “following”, “test set1”, “test set2”. The selection of these features has been guided by exploring their correlation with the target variable. By examining the relationships between these features and the target variable, it was determined that they potentially hold relevant information for distinguishing between genuine and fake profiles and can indicate discrepancies or patterns. The presence of a URL in the profile can be an informative feature for differentiating between fake and very social media profiles similarly, the attribute called “time zone” can help uncover anomalies or inconsistencies by capturing the time zone associated with a profile, aiding in fake profile detection. Measures of user activity, such as statuses count, followers count, favorites count, and listed count, with low or high values can also show signals of inauthenticity [23]. Additionally, “test set1” and “test set2” may include features representing additional data sources or attributes related to the profiles and profile information flagged as fake or real, as well as other features related to testing and validating machine learning algorithms [17]. The aim is to focus on the most informative attributes that can aid in the detection of fake social media profiles [24, 25]. So certain fields like “name”, “screen name”, “created at”, “lang”, “default profile”, “default profile image”, “geo enabled”, “profile image url”, “profile banner url”, “profile use background image”, “profile background image url https”, “profile text color”, “profile image url HTTPs”, “profile sidebar border color”, “profile background tile”, “profile sidebar fill color”, “profile background image url”, “profile background-color”, “profile link color”, “utc offset”, “is translator”, “protected”, “verified”, “notifications”, “description”, “contributors enabled”, “updated” have been removed because some of them contained inconsistent, unreliable data and highly correlated information that are not crucial for this study. Also, with a large number of features, it may suffer from the curse of dimensionality. To address this, some features might be excluded to improve model performance.

4.1.2. Dataset merging and standardization

This study employs two datasets: `fusers.csv`, containing information on fake social media profiles, and `users.csv`, containing details about genuine user profiles. After undergoing preprocessing steps and feature selection, the datasets were merged [16]. This merging process enhances the accuracy and reliability of detection models by allowing identification based on various features. Data preprocessing, using the `StandardScaler` module from the `sklearn` library, standardized the combined data by subtracting the mean and scaling to unit variance [23]. This ensures zero mean and unit variance for all features, making them comparable and preventing biases in the model’s performance due to differing scales.

4.2. Train and test split

This study utilized two datasets, one containing information about fake user profiles, and the other about legitimate user profiles [21]. These datasets have been merged into a single cohesive dataset, incorporating relevant labels and features from both categories. To assess model performance, functionality, and capability, as well as to ensure diverse training, the dataset was split into training and testing sets using the sklearn library, following a common 70:30 ratio. This ratio strikes a balance between adequate training data and robust testing, preventing both underfitting and overfitting [13]. Experimentation confirmed that the 70:30 split provided optimal accuracy, benefiting from a larger training set while maintaining reasonable performance and complexity for accurate predictions.

4.3. Machine learning algorithms

The research aims to develop a machine-learning system for detecting fraudulent social media profiles. Exploring various machine-learning algorithms, the study leverages their ability to learn from complex datasets, recognizing patterns and relationships. Focused on detecting fake profiles, the study evaluates the performance of algorithms, including Dummy Classifier, Logistic Regression Classifier, Naive Bayes Classifier, Multilayer Perceptron Classifier, Decision Tree Classifier, Random Forest Classifier, Support Vector Classifier.

Dummy Classifier. One of the main algorithms utilized in this research is the Dummy Classifier, which plays a vital role in assessing the performance of more advanced models in the detection of fake social media profiles. This simple machine learning algorithm serves as a reference or baseline to evaluate and compare the performance of more sophisticated models. It operates with simple rules based on the class distribution of the training data and assigns labels accordingly. It does not learn from the data or consider data patterns but instead provides class distribution and performance insights through random guessing. Evaluation metrics such as F1 score, Precision, Recall, and Accuracy, along with the confusion matrix, are used to assess the algorithm's performance in correctly identifying fake and genuine profiles.

Logistic Regression. Logistic regression, a key algorithm in this research, is crucial for the objective of employing machine learning to detect fake social media profiles [16, 20]. It models the relationship between input variables and binary outcomes, effectively predicting whether a profile is fake or genuine. The algorithm's simplicity and interpretability make it valuable, but it may face challenges in capturing complex interactions and nonlinear patterns, particularly in the task of identifying fake profiles. Despite these limitations, logistic regression is employed as part of the research's machine-learning algorithms, aiming to achieve the specific objective of detecting fake social media profiles [14].

Gaussian Naive Bayes Classifier. The Gaussian Naive Bayes algorithm is pivotal for detecting fake social media profiles in machine learning [15, 20]. Using Bayes' theorem and prior knowledge, it estimates the likelihood of a profile being fake, assuming feature independence. Fake profile detection identifies deceptive profiles by selecting the class with the highest probability. Known for its speed, effectiveness, and simplicity, Gaussian Naive Bayes is a popular choice for

classifying fake social media profiles, particularly in datasets with numerous features and high-dimensional data [3, 13, 14, 27].

Multilayer Perceptron Classifier. To detect fake social media profiles, powerful feedforward neural network algorithms are applied to learn complex nonlinear relationships between output classes and input variables. The research aims to leverage these algorithms for accurate classification, contributing to the field of fake profile detection using machine learning. The network structure includes input, hidden, and output layers, utilizing backpropagation to minimize prediction errors. Techniques like hyperparameter tuning and regularization are implemented to optimize performance and prevent overfitting. Challenges such as training cost and limited data availability are addressed through preprocessing techniques like dimensionality reduction and feature scaling, enhancing the neural network's efficiency in identifying fake profiles [20].

Multilayer Perceptron Regressor. This powerful neural network algorithm excels in modeling complex input-output relationships, particularly in regression tasks, making it ideal for identifying fraudulent profiles. Leveraging its ability to understand nonlinear relationships, the algorithm effectively captures intricate patterns and irregularities in data associated with fake profiles. Its proficiency in handling noisy and incomplete social media data aligns with the typical characteristics of such datasets. Operating through activation functions in hidden layers, the algorithm adjusts weights using backpropagation, minimizing errors until predicted and actual outputs align. This iterative process aids in detecting fake profiles by estimating the probability based on available input data [20].

Decision Tree Classifier. The decision tree, a widely used method for classification and regression tasks, proves essential in analyzing social media profile data to identify fraudulent accounts [13, 28]. By creating distinct regions based on informative splits from profile features, decision trees offer interpretability and examine specific attributes influencing the classification of fake profiles. They enhance generalization, reduce overfitting, and facilitate accurate model construction. However, potential biases towards features with large values require careful consideration in feature selection and preprocessing. Dealing with imbalanced classes in real-world social media datasets may pose challenges, but employing ensemble methods can enhance decision trees' performance in handling imbalanced data and improving accuracy [15, 20].

Decision Tree Regressor. Decision trees are used to analyze attributes of social media profiles, such as posting behavior and information consistency. Their versatility captures nonlinear relationships, aiding in understanding the traits of fake profiles. The input feature is split based on a criterion minimizing mean squared error, creating subsets fitted with linear regression models. This iterative process forms a tree-like structure, with root and leaf nodes representing starting points and predictions. Effective for continuous variables, decision trees handle both categorical and numerical inputs. They excel at capturing nonlinear relationships, identifying patterns, and generating predictions based on distinctive profile attributes.

Random Forest Classifier. Ensemble methods, employing a combination of decision trees based on different feature subsets, enhance the accuracy and efficacy

of detecting fraudulent profiles [16, 21, 29, 30]. This technique uncovers intricate patterns, offers adaptability, and reduces overfitting by utilizing multiple trees. Integrated into the research on fake profile detection, ensemble methods provide precise and dependable predictions. Their ability to handle datasets with numerous variables and fine-tuning through parameter adjustments makes them valuable in achieving optimal performance. The ensemble of trees utilizes a majority voting scheme to effectively predict the final class label for a new data point [3, 11, 13, 18, 19].

Random Forest Regressor. In this ensemble technique, multiple decision trees are built using randomly selected features, and their predictions are combined by averaging to determine the legitimacy of a social media profile. This method enhances our research by being robust against outliers, revealing non-linear relationships, and offering insights into the importance of different features. Despite its advantages, caution is needed to prevent overfitting risks, particularly when dealing with extensive datasets in social media analysis. The resilience and interpretability of this technique contribute to its effectiveness in detecting fake social media profiles.

Support Vector Classifier. This supervised learning algorithm, Support Vector Classifier (SVC), is designed for binary class separation, distinguishing between genuine and fake profiles. Widely applied in domains like bioinformatics and text/image classification, SVC identifies support vectors crucial for defining decision boundaries [3, 15, 21]. It optimizes a hyperplane to maximize the margin between classes, aiming for effective separation with minimal errors. For detecting fake social media profiles, SVC utilizes extracted features to learn patterns indicative of fraud, analyzing attributes like posting behavior and account activity. It offers flexibility with kernel functions, transforming data for handling complex relationships and non-linear separability. Despite its strengths, SVC's high computational complexity presents challenges with large datasets. Efficient implementation and scalability considerations are crucial for overcoming this limitation in the context of fake profile detection.

Support Vector Classifier Kernel. To enhance the detection of fake profiles, Support Vector Classification (SVC) utilizes kernel functions to transform input data into a higher-dimensional feature space, enabling the identification of intricate patterns and nonlinear relationships [16, 18]. These functions play a vital role in mapping the data, facilitating the separation of profiles into distinct classes (fake and genuine) [14]. By computing inner products between transformed feature vectors, kernel functions empower SVCs to uncover complex patterns, contributing to accurate fake profile detection [23]. The selection of the best kernel function involves experimentation, as it depends on data characteristics and the nature of the problem. Various types of kernel functions provide flexible approaches for handling different data patterns, enhancing the overall detection process [17].

Linear Kernel Function. This function is useful when genuine and fake profiles display linearly separable patterns, allowing for efficient separation using a straight line or hyperplane without high-dimensional transformations. Employed in Support Vector Machines and Support Vector Classification, it computes the dot

product of feature vectors for linear separation. Its benefits include suitability for large datasets, high efficiency, and low memory requirements. The function aims to find an optimal hyperplane to maximize class margins but faces challenges with non-linearly separable data.

Polynomial Kernel Function. The Optimal Hyperplane is found by mapping data into a higher-dimensional space, beneficial for handling non-linearly separable and polynomial-structured data. Careful selection of kernel parameters and degree is crucial to prevent issues like overfitting or underfitting. These functions are particularly useful when dealing with fake profiles exhibiting intricate or nonlinear patterns not easily discernible in the original feature space.

Radial Basis Kernel Function. This function, utilized in both Support Vector Machine and Support Vector Classification, is crucial for handling non-linearly separable data. By mapping inputs into an infinite-dimensional space, it enables the use of non-linear decision boundaries. Its benefits include capturing complex data, facilitating clustering, regression, and classification, while drawbacks involve high cost, large memory requirements, and challenges in tuning the right gamma hyperparameter for controlling decision boundary shape [26].

Sigmoid Kernel Function. The Sigmoid Kernel Function is crucial in neural networks, mapping inputs to higher dimensions. While less common in Support Vector Machines (SVM) and Support Vector Classification (SVC), it's adept at handling two-feature vectors. Ranging between -1 and 1 , it's suitable for nonlinear classification, particularly with challenging data separations. Utilizing this function empowers SVMs/SVCs to adeptly detect fake profiles with complex features by capturing nonlinear decision boundaries.

4.4. Model evaluation with different performance metrics

After training various models, a comprehensive evaluation assesses their accuracy and effectiveness, revealing strengths and weaknesses [4, 23, 9, 10, 15]. These insights are pivotal for identifying areas of improvement and guiding enhancements to better detect fake profiles. The research aims to pinpoint top-performing models, guiding ongoing improvements and strengthening overall detection capabilities [11, 13].

- **Accuracy.** Accuracy is a metric that gauges the model's overall prediction performance by measuring the ratio of correctly classified instances to the total instances in the dataset. A higher accuracy score indicates a more effective distinction between fake and genuine profiles as predicted by the model.

$\text{Accuracy} = (\text{Number of correct predictions}) \div (\text{Total Number of predictions})$.

- **Precision.** Precision measures the accuracy of identifying true positives among instances classified as positive. It is calculated by dividing true positives by the sum of false positives and true positives. A high precision score signifies accurate identification of relevant instances with minimal false positives, crucial for detecting fake profiles.

$\text{Precision} = \text{True Positives} \div (\text{True Positives} + \text{False Positives})$.

- **Recall.** This metric gauges accurate identification of a specific class by dividing true positives by the total actual positives. A high value signals the model's

proficiency in identifying positive instances, minimizing false negatives, and ensuring fake profiles are not overlooked.

$$\text{Recall} = \text{True Positives} \div (\text{True Positives} + \text{False Negatives}).$$

- **F1 Score.** The F1 score, a harmonic mean of precision and recall gauges their balance with values between 0 and 1. Calculated as

$$2 \times (\text{Precision} \times \text{Recall}) \div (\text{Precision} + \text{Recall}),$$

high scores signify accurate positive identifications, while low scores signal a trade-off between precision and recall. Ideal for imbalanced datasets, it addresses class imbalance by excelling in scenarios with notable positive-negative instance disparities,

$$\text{F1 Score} = (2 \times \text{Precision} \times \text{Recall}) \div (\text{Precision} + \text{Recall}).$$

Table 1 represents the Evaluation Metrics of different algorithms.

Table 1. Evaluation metrics

Machine Learning Algorithms	Accuracy	Precision	Recall	F1 Score
Dummy Classifier	0.5112	-	-	-
Logistic Regression	0.9819	0.9839	0.9790	0.9814
Naive Bayes Classifier	0.9677	0.9905	0.9430	0.9662
MLP (Multilayer perceptron) Classifier	0.9843	0.9859	0.9820	0.9839
Decision Tree Classifier	0.9838	0.984	0.9830	0.9835
Decision Tree Regressor	0.9814	0.9791	0.9830	0.9810
Random Forest Classifier	0.9902	0.9791	0.9850	0.9899
Support Vector Classifier	0.9863	0.9919	0.9800	0.9859
SVC Linear Kernel	0.9838	0.9849	0.9820	0.9834
SVC Polynomial Kernel	0.9838	0.9859	0.9810	0.9834
SVC Radial basis function Kernel	0.9863	0.9919	0.9800	0.9859
SVC sigmoid Kernel	0.9619	0.9520	0.9710	0.9614

- **Mean Absolute Error.** This metric calculates the average absolute difference between actual and predicted values, evaluating regression model performance. A lower score signifies greater accuracy and fit. Notably, it shares the unit of the data, simplifying interpretation, and is less affected by outliers compared to other regression metrics.

Mean absolute error = $(1/n) \times (\sum |y - \hat{y}|)$, y represents actual value, \hat{y} is the value predicted, and n is the sample's total number.

- **Mean Squared Error.** This metric measures the average squared difference between actual and predicted values, with lower scores indicating better model fit and prediction accuracy. It is less interpretable than the Mean Absolute Error but more sensitive to outliers due to the squared data unit.

Mean squared error = $(1/n) \times (\sum (y - \hat{y})^2)$, y represents the actual value, \hat{y} is the value predicted, and n is the sample's total number.

- **Root Mean Squared Error.** The score is determined by averaging the discrepancies between actual and predicted values, then taking the square root of the mean squared error. This metric, sensitive to outliers, yields a score with the same units as the data. A lower score signifies a superior fit and greater accuracy in the model's predictions.

Root Mean Squared Error = $\text{sqrt}((1/n) \times (\sum (y - \hat{y})^2))$, where y is the actual value, \hat{y} is the value predicted, and n is the sample's total number. Table 2 represents regression metrics.

Table 2. Regression metrics

Machine Learning Algorithms	Mean absolute Error	Mean Squared Error	Root Mean squared Error
Logistic Regression	0.0361	0.0722	0.2688
Naive Bayes Classifier	0.0644	0.1289	0.3590
MLP (Multilayer perceptron) Classifier	0.0312	0.0625	0.25
MLP Regressor	0.0292	0.0585	0.2420
Decision Tree Classifier	0.0322	0.0644	0.2538
Decision Tree Regressor	0.0371	0.0742	0.2724
Random Forest Classifier	0.0195	0.0390	0.1976
Random Forest Regressor	0.1115	0.0497	0.2229
Support Vector Classifier	0.0273	0.0546	0.2338
SVC Linear Kernel	0.0322	0.0644	0.2538
SVC Polynomial Kernel	0.0322	0.0644	0.2538
SVC rbf Kernel	0.0273	0.0546	0.2338
SVC sigmoid Kernel	0.0761	0.1523	0.3903

- **Confusion Matrix.** Confusion matrix serves as a valuable tool for assessing the performance of machine learning models. Analyzing true positives assesses how well the algorithm detects fake profiles, indicating its effectiveness in spotting deceptive accounts. True negatives gauge the algorithm's accuracy in identifying genuine profiles, demonstrating its ability to distinguish real users from fake ones. False positives highlight potential areas of improvement, where the algorithm may be overly sensitive, impacting user trust. Conversely, false negatives reveal instances where the algorithm overlooks fake profiles, emphasizing the need for enhancements to minimize the risk of missing fraudulent accounts. According to this study, the Random Forest Classifier had the highest accuracy and the confusion Matrix of the model with the highest accuracy is presented in Fig 2.

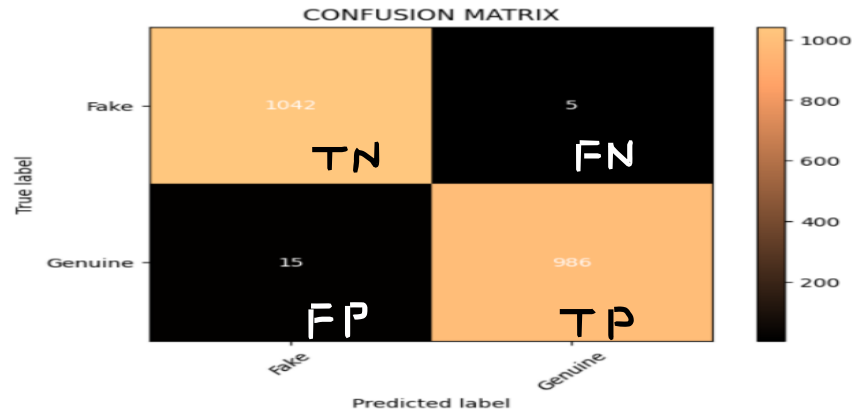


Fig. 2. Confusion matrix of Random Forest Classifier

- **Cross Validation Accuracy.** Cross-validation techniques ensure a robust estimation of a model's performance, offering a comprehensive understanding of its generalization ability. In fake profile detection, the dataset is divided into k subsets. The model undergoes training on one subset and testing on the remaining folds, a process repeated multiple times. The average accuracy across folds determines

overall performance, with K-Fold cross-validation being the most common. Equal-sized subsets are drawn from the dataset, with the model trained on $k - 1$ subsets and tested on the remaining one, repeated k times. Each iteration's performance metrics are averaged to assess the algorithm's effectiveness in detecting fake social media profiles. The cross-validation accuracies of different algorithms used in this study are given in Table 3.

Table 3. Cross validation accuracy

Machine Learning Algorithms	Cross validation accuracy
Logistic Regression	0.9968
Naive Bayes Classifier	0.9947
MLP (Multilayer perceptron) Classifier	0.9980
MLP Regressor	0.9978
Decision Tree Classifier	0.9845
Decision Tree Regressor	0.9841
Random Forest Classifier	0.9987
Random Forest Regressor	0.9976
SVC Classifier	0.9978
SVC Classifier Linear Kernel	0.9964
SVC Classifier Polynomial Kernel	0.9968
SVC Classifier Radial basis function Kernel	0.9978
SVC sigmoid Kernel	0.9843

4.5. Post-processing using Sampling Techniques

Table 4. Accuracy of different models

Machine Learning Algorithms	Random under-sampling accuracy	Random over sampling accuracy	SMOTE accuracy	ADASYN accuracy
Logistic Regression	0.9765	0.9765	0.9760	0.9750
Naive Bayes Classifier	0.9780	0.9780	0.9780	0.9765
MLP (Multilayer perceptron) Classifier	0.9833	0.9858	0.9868	0.9863
Decision Tree Classifier	0.9843	0.9838	0.9824	0.9775
Random Forest Classifier	0.9892	0.9892	0.9887	0.9892
SVC Classifier	0.9838	0.9838	0.9843	0.9838

Some of the Sampling techniques have been implemented to address class imbalance in the dataset [31].

- **Random Sampling Technique.** This technique addresses class imbalance in datasets, where one class has fewer instances than the other, leading to biased modeling. It balances class distribution by randomly removing instances from the majority class until the desired ratio is achieved. Applied to detecting fake social media profiles, it ensures equal representation of genuine and fake profiles, reducing bias and improving overall performance.

- **Random Over Sampling Technique.** This technique tackles dataset class imbalance by randomly duplicating minority class samples. It creates a balanced distribution among classes, fostering meaningful pattern learning from both genuine and fake profiles. In identifying fake social media profiles, this approach enhances the representation of the minority class (fake profiles), improving the model's understanding of their characteristics.

- **Synthetic Minority Over Sampling Technique (SMOTE)** or Synthetic Minority Over-sampling Technique, addresses dataset class imbalance by generating synthetic samples similar to the minority class. It works by identifying k nearest neighbors in the feature space, randomly selecting neighbors to create new samples along the connecting line, and repeating until the desired level of minority class oversampling is achieved. This technique is particularly useful for small datasets or when the minority class is underrepresented. In the context of detecting fake social media profiles, SMOTE can be employed to increase the number of fake profiles, enhancing the model's ability to learn and improve detection capabilities.

- **Adaptive Synthetic Sampling Technique (ADASYN)** is an extension of SMOTE, tackles class imbalance by adapting synthetic sample generation to the class distribution. It focuses on low-density areas of the minority class, prioritizing harder-to-learn instances. Particularly beneficial for detecting fake social media profiles, ADASYN allows the model to concentrate on challenging cases within the minority class, enhancing representation through synthetic samples.

5. Discussion

The study highlights the Random Forest Classifier's effectiveness in detecting fake social media profiles, showcasing high accuracy across diverse evaluation metrics. Outperforming other algorithms in precision, recall, and F1 score, it demonstrated consistency and resilience. Regression metrics revealed lower error rates, emphasizing favorable results. Confusion Matrix analysis provided insights into strengths and limitations, confirming accurate classification of fake profiles. Cross-validation reinforced reliability, with the classifier consistently achieving high accuracy on unseen data. The algorithm's advantageous features, including ensemble learning and effective handling of outliers, contribute to its superior performance in capturing intricate patterns and differences in detecting fake social media profiles.

6. Result

The study highlights the **Random Forest Classifier's** consistent and robust performance in detecting fake profiles, achieving the highest accuracy (0.9902343) among evaluated algorithms. Various metrics, including Accuracy, Precision, Recall, and F1 Score, were considered, with the Random Forest Classifier consistently outperforming others. Table 1 summarizes the Evaluation Metrics for various algorithms used to detect fake profiles in this study, providing valuable evidence and insights into the effectiveness of the Random Forest Classifier.

Various algorithms were assessed using regression metrics, including Mean Absolute Error, Mean Squared Error, and Root Mean Squared Error, as part of a comprehensive evaluation process. The Random Forest Classifier demonstrated superior performance with lower errors compared to other algorithms, specifically achieving a Mean Absolute Error of 0.019531, Mean Squared Error of 0.03906, and Root Mean Squared Error of 0.1976423 (as presented in Table 2). This highlights the algorithm's effectiveness in minimizing discrepancies and enhancing predictive

accuracy within the study's evaluation framework. The study found that the Random Forest Classifier achieved the highest accuracy, as illustrated by the corresponding confusion matrix in Fig. 1. The Confusion Matrix, analyzing true positives, false positives, true negatives, and false negatives, serves as a crucial tool for evaluating the model's performance in detecting fake social media profiles. The Random Forest Classifier's highest accuracy underscores its effectiveness in distinguishing between authentic and fake profiles. In this study, cross-validation accuracies were calculated using 10-fold cross-validation. This method averages accuracies across folds, offering a more reliable performance estimate and reducing variance compared to a single test-train split. Table 3 displays the cross-validation accuracies for the various algorithms assessed in the study. The Random Forest Classifier achieved the highest Cross Validation Accuracy at 0.9987, surpassing other algorithms. Validation with various sampling techniques consistently showed its superior accuracy of 0.9893, indicating strong generalization to unseen data. Table 4 represents the accuracies different machine learning algorithms and sampling methods.

7. Conclusion

The proposed study presents a comprehensive approach that utilizes multiple machine-learning algorithms. Each algorithm is evaluated using different performance metrics, and the accuracy of the system is significantly improved through the implementation of various sampling techniques. The study demonstrates that the performance of each algorithm varies, and the choice of the most suitable algorithm depends on the characteristics of the dataset and the main objective. In conclusion, this study highlights the potential of machine learning algorithms in addressing the challenges associated with fake social media profile detection. It provides valuable insights for enhancing the integrity and security of online platforms by identifying and removing fake profiles. The findings of this study can serve as a baseline for future advancements in this field.

This paper highlights future improvement prospects by incorporating deep learning, such as Convolutional and Recurrent Neural Networks, for enhanced image and text classification. Adding Natural Language Processing aids language analysis, identifying deceptive language in descriptions and messages to detect fake profiles. Social Network Analysis further contributes by scrutinizing behavior, activities, and network structures of fake profiles. Expanding this system to detect diverse fraudulent activities in social media content promises to uphold the integrity and security of online platforms.

References

1. Kemp, S. Digital 2023 Global Overview Report – Reports – Datareportal – Global Digital Insights, DataReportal.
<https://datareportal.com/reports/tag/Digital+2023+Global+Overview+Report>
2. Dean, B. How Many People Use Social Media in 2023? (65+ Statistics), Backlinko (Accessed 22 June 2023).
<https://backlinko.com/social-media-users>

3. Ramalingam, D., V. Chinnaiah. Fake Profile Detection Techniques in Large-Scale Online Social Networks: A Comprehensive Review. – Computers and Electrical Engineering, Vol. **65**, 2018, pp. 165-177. DOI: 10.1016/j.compeleceng.2017.05.020.
4. Goyal, B., N. S. Gill, P. Gulia. Detection of Fake Profiles on Online Social Media. – In: Proc. of the Strategy National Conference on Computational Intelligence and Data Science (NCCIDS'23). MDU Rohtak, 2023.
https://www.researchgate.net/publication/369643807_Detection_of_Fake_Profiles_on_Online_Social_Media_A_Strategy
5. Singh, N., T. Sharma, A. Thakral, T. Choudhury. Detection of Fake Profile in Online Social Networks Using Machine Learning. – In: Proc. of International Conference on Advances in Computing and Communication Engineering (ICACCE'18), Paris, France, 2018, pp. 231-234. DOI: 10.1109/ICACCE.2018.8441713.
6. Nikhitha, K. V., K. Bhavya, D. U. Nandini. Fake Account Detection on Social Media Using Random Forest Classifier. – In: Proc. of 7th International Conference on Intelligent Computing and Control Systems (ICICCS'23), Madurai, India, 2023, pp. 806-811. DOI: 10.1109/ICICCS56967.2023.10142841.
7. Ritchie, J. N. A., et al. Scams Starting on Social Media Proliferate in Early 2020. Federal Trade Commission, 2022.
<https://www.ftc.gov/news-events/data-visualizations/data-spotlight/2020/10/scams-starting-social-media-proliferate-early-2020>.
8. Spoorthy, A. S., S. Sinha. Trust Based Fake Node Identification in Social Networking Sites. – IOP Conference Series: Materials Science and Engineering, Vol. **1123**, 2021, No 1, p. 012036. DOI:10.1088/1757-899x/1123/1/012036.
9. Meligy, A., M. H. Ibrahim, F. M. Torky. Identity Verification Mechanism for Detecting Fake Profiles in Online Social Networks. – International Journal of Computer Network and Information Security, Vol. **9**, 2017, No 1, pp. 31-39. DOI:10.5815/ijcnis.2017.01.04.
10. Sheikh, S. An Efficient Method for Detection of Fake Accounts on the Instagram Platform. – Revue d'Intelligence Artificielle, Vol. **34**, 2020, No 4, pp. 429-436. DOI:10.18280/ria.340407.
11. Reddy, K. D. Fake Profile Identification Using Machine Learning. – International J. of Scientific Research in Science Engineering, 2020 [Preprint].
12. Latha, P., et al. Fake Profile Identification in Social Network Using Machine Learning and NLP. – In: Proc. of International Conference on Communication, Computing and Internet of Things (IC3IoT'22), 2022, [Preprint]. DOI: 10.1109/ic3iot53935.2022.9767958.
13. Elyusufi, Y., Z. Elyusufi, M. A. Kbir. Social Networks Fake Profiles Detection Using Machine Learning Algorithms. – Innovations in Smart Cities Applications Edition 3, 2020, pp. 30-40. DOI:10.1007/978-3-030-37629-1_3.
14. Mughaid, A., I. Obeidat, E. Abu Elsou, A. Alnajjar et al. A Novel Machine Learning and Face Recognition Technique for Fake Accounts Detection System on Cyber Social Networks. – Multimedia Tools and Applications, Vol. **82**, 2023, pp. 26353-26378. DOI: 10.1007/s11042-023-14347-8.
15. Patel, K., S. Agrahari, S. Srivastava. Survey on Fake Profile Detection on Social Sites by Using Machine Learning Algorithm. – In: Proc. of 8th International Conference on Technologies and Optimization (Trends and Future Directions) (ICRITO'20), 2020 [Preprint]. DOI:10.1109/icrito48877.2020.9197935.
16. Kondeti, P., L. P. Yerramreddy, A. Pradhan, G. Swain. Fake Account Detection Using Machine Learning. – In: V. Suma, N. Bouhmala, H. Wang, Eds. Evolutionary Computing and Mobile Sustainable Networks. – Lecture Notes on Data Engineering and Communications Technologies, Vol. **53**, Springer, Singapore, 2021.
https://doi.org/10.1007/978-981-15-5258-8_73
17. Rao, K. S., S. Gutha, B. D. Raju. Detecting Fake Account on Social Media Using Machine Learning Algorithms. – International Journal of Control and Automation, Vol. **13**, 2020, pp. 95-100.
18. Shreya, K., A. Kothapelly, D. V. H. Shanmugasundaram. Identification of Fake Accounts in Social Media Using Machine Learning. – In: Proc. of 4th International Conference on Emerging Research in Electronics, Computer Science and Technology (ICERECT'22), Mandya, India, 2022, pp. 1-4. DOI: 10.1109/ICERECT56837.2022.10060194.
19. Harish, K., R. Naveen Kumar, Dr. J. Briso Becky Bell. Fake Profile Detection Using Machine Learning. – International Journal of Scientific Research in Science, Engineering and Technology, 2023, pp. 719-725. DOI:10.32628/ijrsrset2310264.
20. Munoz, S. D., P. G. E. Pinto. A Dataset for the Detection of Fake Profiles on Social Networking Services. – In: Proc. of International Conference on Computational Science and Computational Intelligence (CSCI'20), 2020 [Preprint]. DOI:10.1109/csci51800.2020.00046.
21. Mesram, P., B. Karbikar. Automatic Detection of Fake Profile Using Machine Learning on Instagram. – International Journal of Scientific Research in Science and Technology, 2021 pp. 117-127. DOI: 10.32628/ijrsrset218330.

22. Aydin, İ., M. Sevi, M. U. Salur. Detection of Fake Twitter Accounts with Machine Learning Algorithms. – In: Proc. of International Conference on Artificial Intelligence and Data Processing (IDAP'18), Malatya, Turkey, 2018, pp. 1-4. DOI: 10.1109/IDAP.2018.8620830.
23. Khaled, S., N. El-Tazi, H. M. O. Mokhtar. Detecting Fake Accounts on Social Media. – In: Proc. of IEEE International Conference on Big Data (Big Data), Seattle, WA, USA, 2018, pp. 3672-3681. DOI: 10.1109/BigData.2018.8621913.
24. Akhlat, Y., et al. A New Noisy Random Forest-Based Method for Feature Selection. – Cybernetics and Information Technologies, Vol. **21**, 2021, No 2, pp. 10-28.
25. Venkatesh, B., J. Anuradha. A Review of Feature Selection and Its Methods. – Cybernetics and Information Technologies, Vol. **19**, 2019, No 1, pp. 3-26.
26. Jyothi, M. A., T. Sridevi, K. Rajani, U. Channabasava, S. Bethu. Rapport of Counterfeit Profiles in Social Networking Using Machine Learning Techniques. – In: Proc. of 3rd International Conference on Intelligent Communication and Computational Techniques (ICCT'23), Jaipur, India, 2023, pp. 1-7. DOI: 10.1109/ICCT56969.2023.10076017.
27. Zaman, B., et al. An Indonesian Hoax News Detection System Using Reader Feedback and Naïve Bayes Algorithm. – Cybernetics and Information Technologies, Vol. **20**, 2020, No 1, pp. 82-94.
28. Patil, D. R., J. B. Patil. Malicious URL's Detection Using Decision Tree Classifiers and Majority Voting Technique. – Cybernetics and Information Technologies, Vol. **18**, 2018, No 1, pp. 11-29.
29. Chakraborty, P., et al. Fake Profile Detection Using Machine Learning Techniques. – Journal of Computer and Communications, Vol. **10**, 2022, No 10, pp. 74-87. DOI: 10.4236/jcc.2022.1010006.
30. Huanrui, H. New Mixed Kernel Functions of SVM Used in Pattern Recognition. – Cybernetics and Information Technologies, Vol. **16**, 2016, No 5, pp. 5-14.
31. Matthews, L. M., H. Seetha. On Improving the Classification of Imbalanced Data. – Cybernetics and Information Technologies, Vol. **17**, 2017, No 1, pp. 45-62.

Received: 29.05.2023; Second Version: 04.12.2023; Accepted: 08.01.2024