

A PROJECT REPORT ON

## **Fake Instagram Profile Identification and Classification using Machine Learning**

SUBMITTED TO THE SAVITRIBAI PHULE PUNE UNIVERSITY, PUNE  
IN THE PARTIAL FULFILLMENT OF THE REQUIREMENTS  
FOR THE AWARD OF THE DEGREE

OF

**BACHELOR OF ENGINEERING (COMPUTER ENGINEERING)**

**SUBMITTED BY**

**STUDENT NAME**

Apurv Kamalakar Badave  
Aniruddha Shivaji Lalge  
Nikhil Lala Saheb Elajale  
Pankaj Dayaram Godara

**SEAT NO**

B190364217  
B190364389  
B190364283  
B190364304



**Sinhgad Institutes**

**DEPARTMENT OF COMPUTER ENGINEERING**

**STES'S SMT. KASHIBAINA VALE COLLEGE OF ENGINEERING**

**VADGAON BK, OFF SINHGAD ROAD, PUNE 411041**

**SAVITRIBAI PHULE PUNE UNIVERSITY**

**2023-2024**



**Sinhgad Institutes**

## **CERTIFICATE**

This is to certify that the project report entitles

### **“Fake Instagram Profile Identification and Classification using Machine Learning”**

Submitted by

<b>Apurv Kamalakar Badave</b>	<b>Seat No: B190364217</b>
<b>Aniruddha Shivaji Lalge</b>	<b>Seat No: B190364389</b>
<b>Nikhil Lala Saheb Elajale</b>	<b>Seat No: B190364283</b>
<b>Pankaj Dayaram Godara</b>	<b>Seat No: B190364304</b>

is a bonafide work of this institute and the work has been carried out by them under the supervision of **Prof. Priyanka Kinage** and it is approved for the partial fulfilment of the requirement of Savitribai Phule Pune University, for the award of the degree of **Bachelor of Engineering** (Computer Engineering)

**(Prof. Priyanka Kinage)**

Internal Guide

Department of Computer Engineering

**(Prof. R. H. Borhade)**

Head,

Department of Computer Engineering

**(Dr. A. V. Deshpande)**

Principal,

Smt. Kashibai Navale College of Engineering Pune – 41

## **ACKNOWLEDGEMENT**

It gives us great pleasure in presenting the preliminary project report on '**Fake Instagram Profile Identification and Classification using Machine Learning**'. We would like to take this opportunity to thank our internal guide **Prof. Priyanka Kinage** for giving us all the help and guidance we needed. We are really grateful to her for their kind support. Their valuable suggestions were very helpful.

We are also grateful to **Prof. R. H. Borhade**, Head of Computer Engineering Department, SKNCOE for his indispensable support, suggestions.

We would like to convey our gratitude to **Dr. A. V. Deshpande (Principal)** all the teaching and non-teaching staff members of the Computer Engineering Department who gave us the freedom to explore and guided us the right way, also our friends and families for their valuable suggestions and support.

<b>Apurv Kamalakar Badave</b>	<b>(B190364217)</b>
<b>Aniruddha Shivaji Lalge</b>	<b>(B190364389)</b>
<b>Nikhil Lalasaheb Elajale</b>	<b>(B190364283)</b>
<b>Pankaj Dayaram Godara</b>	<b>(B190364304)</b>

## **ABSTRACT**

Social media platforms have become integral to modern communication, enabling users to connect, share, and engage in various activities. However, the rise of fake profiles on platforms like Instagram possess significant challenges related to user privacy, security, and trust. This work presents a novel approach to identify and classify fake Instagram profiles using machine learning techniques. The findings of this research contribute to the ongoing efforts to combat the proliferation of fake profiles on Instagram and other social media platforms. By leveraging machine learning techniques and a comprehensive feature set, the proposed model demonstrates promising results in identifying and classifying fake profiles, thereby promoting a safer and more trustworthy online environment. This research opens avenues for further exploration, including the integration of real-time data streams and the adaptation of the model to other social media platforms.

**Keywords:** Profile identification, User authentication, Data preprocessing, Model training, Online security, Machine learning.

# INDEX

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Background . . . . .	2
1.2	Motivation . . . . .	2
1.3	Area Of Project . . . . .	2
1.4	Technical Keywords . . . . .	3
1.5	Objectives . . . . .	3
<b>2</b>	<b>Literature Survey</b>	<b>4</b>
2.1	Study Of Research Paper . . . . .	5
<b>3</b>	<b>Software Requirement Specification</b>	<b>13</b>
3.0.1	Problem Statement.....	14
3.0.2	Assumptions and Dependencies.....	14
3.1	FUNCTIONAL REQUIREMENTS.....	14
3.1.1	System Feature1(Functional Requirement).....	14
3.1.2	System Feature2(Functional Requirement).....	15
3.2	EXTERNAL INTERFACE REQUIREMENT.....	15
3.2.1	User Interface.....	15
3.2.2	Hardware Interfaces.....	15
3.2.3	Software Interfaces .....	16
3.3	NON FUNCTIONAL REQUIREMENT .....	16
3.3.1	Performance Requirements.....	16
3.3.2	Safety Requirement .....	16
3.3.3	Software Quality Attributes.....	16
3.4	SYSTEM REQUIREMENTS .....	17
3.4.1	Database Requirements .....	17

3.4.2	Software Requirements(Platform Choice) .....	17
3.4.3	Hardware Requirements.....	17
3.5	ANALYSIS MODEL: SDLC MODEL TO BE APPLIED.....	18
3.6	SYSTEM IMPLEMENTATION PLAN.....	19
<b>4</b>	<b>System Design</b>	<b>20</b>
4.1	SYSTEM ARCHITECTURE .....	21
4.1.1	Module .....	22
4.1.2	Data Flow Diagram.....	22
4.2	UML DIAGRAMS .....	24
<b>5</b>	<b>Software Information</b>	<b>27</b>
5.1	Python.....	28
5.2	Anaconda.....	29
5.3	SPYDER.....	31
5.1.1	Features.....	31
5.1.2	Other Specifications.....	32
<b>6</b>	<b>Software Testing</b>	<b>33</b>
6.1	Introduction.....	34
6.2	TYPES OF TESTING USED .....	34
6.2.1	Unit Testing .....	34
6.2.2	Integration Testing.....	34
6.2.3	White Box & Black Box Testing.....	35
6.3	Test Cases & Test Result .....	36
<b>7</b>	<b>Result</b>	<b>38</b>
7.1	Login Page .....	39
7.2	Registration Page .....	39
7.3	GUI Main.....	40
7.4	Outputs .....	41

<b>8 Conclusion</b>	<b>45</b>
<b>9 References</b>	<b>47</b>
<b>10 Appendix</b>	<b>50</b>
10.1 Appendix A	50
10.2 Appendix B	51
10.3 Appendix C	66

# List of Figures

4.1	System Architecture.....	21
4.2	Data Flow(0) diagram.....	23
4.3	Data Flow(1) diagram.....	23
4.4	Data Flow(2) diagram.....	23
4.5	Class Diagram.....	24
4.6	Usecase Diagram .....	25
4.7	Activity Diagram .....	25
4.8	Sequence Diagram .....	26
6.1	GUI Testing .....	36
6.2	Registration Test Case .....	36
6.3	Login Test Case .....	37
7.1	Login Page .....	39
7.2	Registration Page .....	40
7.3	GUI Main .....	41
7.4.1	Output 1 Page.....	42
7.4.1	Output 2 Page.....	42
7.4.2	Output 3 Page.....	43
7.4.2	Output 4 Page.....	44
7.4.2	Output 5 Page.....	44

# **CHAPTER 1**

## **INTRODUCTION**

## **1.1 BACKGROUND**

Fake Instagram profiles can range from automated bots posting spam to sophisticated imposters attempting to deceive genuine users for financial gain, social manipulation, or other illicit activities. Traditional methods of manual inspection and reporting are insufficient to handle the sheer volume of profiles and interactions, necessitating the use of advanced technological solutions. Machine learning has emerged as a powerful tool in addressing the issue of fake profiles on social media platforms. By harnessing the computational power of machine learning algorithms, it is possible to automatically identify and classify fake profiles based on distinctive patterns and characteristics. The combination of the growing influence of social media, the challenges posed by fake profiles, and the advancements in machine learning techniques has led to the development of solutions aimed at identifying and classifying these profiles. This research addresses the need for a safer and more trustworthy online environment by proposing a comprehensive approach to tackle the issue of fake Instagram profiles using machine learning.

## **1.2 MOTIVATION**

The research on “Fake Instagram Profile Identification and Classification using Machine Learning” is relevant due to its potential to address critical issues related to user trust, online safety, and platform integrity. By leveraging the power of machine learning, this research offers practical solutions that align with the needs of the digital age.

## **1.3 AREA OF PROJECT**

1. Machine Learning

## **1.4 TECHNICAL KEYWORDS**

1. Machine Learning
2. SVM/RF/DT Algorithm
3. Pre-processing
4. Feature Extraction

## **1.5 OBJECTIVES**

The research aims to contribute to a safer online environment with enhanced user experiences and support the ongoing efforts of social media platforms to combat fake profiles.

Identify and extract relevant features from the collected Instagram profiles and associated content. The research encompass developing an effective machine learning based model for identifying and classifying fake Instagram profiles.

# **CHAPTER 2**

# **LITERATURE SURVEY**

The research encompass developing an effective machine learning-based model for identifying and classifying fake Instagram profiles.

## **2.1 STUDY OF RESEARCH PAPER**

### **1. Paper Name: PREDICTION OF FAKE INSTAGRAM PROFILES USING MACHINE LEARNING**

**Author:** Anupriya1, V. Sowmiya, Dr. G. Devika.

#### **Abstract:**

The majority of people now use social networking sites as part of their everyday lives. Every day, a vast number of people build profiles on social networking sites and connect with others, regardless of their place or time. False identities play an important role in advanced persisted threats and are also involved in other malicious activities. Users of social networking sites not only profit from them, but they also face security concerns about their personal details. To assess who is promoting threats in social networks, we must first identify the user's social network profiles. It is necessary to differentiate between genuine and fake accounts on social media based on the classification. Detecting fake accounts on social media has historically focused on a number of classification methods. However, it is possible to boost the accuracy of fake profile identification in social media. Machine Learning and technology is used in the proposed work to increase the percentage of fake profile prediction.

**2. Paper Name:** DETECTION OF FAKE ACCOUNTS IN INSTAGRAM USING  
MACHINE LEARNING.

**Author:** Ananya Dev, Hamsashree Reddy, Manjistha Sinha.

**Abstract:**

With the advent of the Internet and social media, while hundreds of people have benefits from the vast sources of information available, there has been an enormous increase in the rise of cybercrimes, particularly targeted towards women. According to a 2019 report in the [4] Economics Times, India has witnessed a 457 rise in cybercrime in the five year span between 2011 and 2016. Most speculate that this is due to impact of social media such as Facebook, Instagram and Twitter on our daily lives. While these definitely help in creating a sound social network, creation of user accounts in these sites usually needs just an email id. A real life person can create multiple fake IDs and hence impostors can easily be made. Unlike the real world scenario where multiple rules and regulations are imposed to identify oneself in a unique manner (for example while issuing one's passport or driver's license), in the virtual world of social media, admission does not require any such checks. In this paper, we study the different accounts of Instagram, in particular and try to assess an account as fake or real using Machine Learning techniques namely Logistic Regression and Random Forest Algorithm.

**3. Paper Name:** SURVEY ON FAKE PROFILE DETECTION ON SOCIAL SITES BY USING MACHINELEARNING ALGORITHM

**Author:** Kumud Patel, Sudhanshu Agrahari, Saijshree Srivastava

**Abstract:**

To avoid the spam message, malicious and cyber bullies activities which are mostly done by the fake profile. These activities challenge the privacy policies of the social network communities. These fake profiles are responsible for spread false information on social communities. To identify the fake profile, duplicate, spam and bots account there is much research work done in this area. By using a machine-learning algorithm, most of the fake accounts detected successfully. This paper represents the review of Fake Profile Detection on Social Site by Using Machine Learning.

**4. Paper Name: FAKE ACCOUNTS DETECTION ON SOCIAL MEDIA  
(INSTAGRAM AND TWITTER)**

**Author:** Dr.P.V. Kumar, Shanthi Vardhan, Y. Kavya, K. Badri Singh.

**Abstract:**

Online Social Networks (OSNs) have grown in popularity among today's youth, having an effect on their social life and motivating them to sign up for various social media platforms. Social media sites offer the required tools for a range of tasks, including news generation, Fake accounts have grown to be a serious issue with the growth of social media, endangering user security and platform integrity. In this work, we investigate how well machine learning (ML) algorithms identify fake accounts on social media sites like Twitter and Instagram. In order to train ML models for spotting fake accounts, we examine user behavior and account attributes, extracting parameters like the number of followers, activity level, and posting behavior. To preprocess the data and use different ML techniques, such Random Forest, Support Vector Machines, and XG boost, to categorize and identify bogus accounts, we employ Python packages. The findings demonstrate that ML algorithms can accurately detect patterns and abnormalities suggestive of fake accounts and achieve high precision in fake account detection.

## **5. Paper Name: INSTAGRAM FAKE AND AUTOMATED ACCOUNT DETECTION**

**Author:** Fatih Cagatay Akyon, M. Esat Kalfaoglu.

### **Abstract:**

Fake engagement is one of the significant problems in Online Social Networks (OSNs) which is used to increase the popularity of an account in an inorganic manner. The detection of fake engagement is crucial because it leads to loss of money for businesses, wrong audience targeting in advertising, wrong product predictions systems, and unhealthy social network environment. This study is related with the detection of fake and automated accounts which leads to fake engagement on Instagram. As far as we know, there is no publicly available dataset for fake and automated accounts. For this purpose, two dataset have been generated for the detection of fake and automated accounts. For the detection of these accounts, machine learning algorithms like Naive Bayes, logistic regression, support vector machines and neural networks are applied. Additionally, for the detection of automated accounts, cost sensitive genetic algorithm is applied because of the unnatural bias in the dataset. To deal with the unevenness problem in the fake dataset, Smotenc algorithm is implemented. For the automated and fake account detection problem, 86 and 96 are obtained, respectively.

## **6. Paper Name: Fake Profile Detection Using Machine Learning Techniques**

**Author:** Partha Chakraborty, Mahim Musharof Shazan, Mahamudul Nahid,  
Md.Kaysar Ahmed, Prince Chandra Talukder

### **Abstract:**

Our lives are significantly impacted by social media platforms such as Facebook, Twitter, Instagram, LinkedIn, and others. People are actively participating in it the world over. However, it also has to deal with the issue of bogus profiles. False accounts are frequently created by humans, bots, or computers. They are used to disseminate rumors and engage in illicit activities like identity theft and phishing. So, in this project, the author'll talk about a detection model that uses a variety of machine learning techniques to distinguish between fake and real Twitter profiles based on attributes like follower and friend counts, status updates, and more. The author used the dataset of Twitter profiles, separating real accounts into TFP and E13and false accounts into INT, TWT, and FSF. Here, the author discusses LSTM, XG Boost, Random Forest, and Neural Networks. The key characteristics are chosen to assess a social media profile's authenticity. Hyperparameters and the architectureare also covered. Finally, results are produced after training the models. The output is therefore 0 for genuine profiles and 1 for false profiles. When a phony profile is discovered, it can be disabled or destroyed so that cyber security problems can be prevented.

**7. Paper Name:** PREDICTION OF FAKE INSTAGRAM PROFILES USING MACHINE LEARNING

**Author:** S. Saranya Shree, C. Subhiksha, R. Subhashini

**Abstract:**

The majority of people now use social networking sites as part of their everyday lives. Every day, a vast number of people build profiles on social networking sites and connect with others, regardless of their place or time. Users of social networking sites not only profit from them, but they also face security concerns about their personal details. To assess who is promoting threats in social networks, we must first identify the users' social network profiles. We may differentiate between genuine and false accounts on social media based on the classification. Detecting false accounts on social media has historically focused on a number of classification methods. However, we must boost the accuracy of fake profile identification in social media. We suggest machine learning and natural language processing (NLP) in this paper to increase the percentage of fake profile prediction. The Support Vector Machine (SVM) and the Naive Bayes algorithm are two algorithms that can be used.

**8. Paper Name:** Fake Profile Identification using Machine Learning

**Author:** Samala Durga, Prasad Reddy

**Abstract:**

In the present generation, the social life of everyone has become associated with on-line social networks. These sites have made a drastic change in the way we pursue our social life. Making friends and keeping in contact with them and their updates has become easier. But with their rapid growth, many problems like fake profiles, on-line impersonation have also grown. There are no feasible solutions exist to control these problems. In this paper, I came up with a framework with which the automatic identification of fake profiles is possible and is efficient. This framework uses classification techniques like Random Forest Classifier to classify the profiles into fake or genuine classes. As this is an automatic detection method, it can be applied easily by online social networks that have millions of profiles whose profiles cannot be examined manually.

# **CHAPTER 3**

## **SOFTWARE REQUIREMENT SPECIFICATION**

### **3.0.1 PROBLEM STATEMENT**

To develop a windows based model to identify and classify fake Instagram profiles using machine learning algorithms such as SVM, Random Forest and Decision Tree algorithms.

### **3.0.2 ASSUMPTIONS AND DEPENDENCIES**

1. User must install the Python on his pc.
2. User has to install the Spyder on his pc.
3. User has to login to the system.

## **3.1 FUNCTIONAL REQUIREMENTS**

### **3.1.1 System Feature1 (Functional Requirement)**

**Admin:** Admin module will be on web module. Admin will verify user information and allow or reject to user. Load the Data set.

**User:** User registers into system with personal information. Automatically user verification request send to admin. After verification user can login into system.

### **3.1.2 System Feature2 (Functional Requirement)**

**System:** By using RF algorithm, enhance Machine Learning techniques like Feature Selection, Handling Missing Data & Improving Stability and Accuracy.

## **3.2 EXTERNAL INTERFACE REQUIREMENT**

### **3.2.1 User Interface**

Application of fake Instagram Profile Identification and Classification using Machine Learning.

### **3.2.2 Hardware Interfaces:**

- RAM: 8GB

As we are using Machine Learning Algorithm, various high level Libraries and data loading should be fast hence minimum laptop RAM required is 8 GB.

- Hard Disk: 40 GB

- Processor: Intel i5 Processor

- IDE: Spyder

Spyder IDE the Integrated Development Environment to be used.

- Coding Language: Python Version 3.5

Highly specified programming language for Machine Learning because of availability of High Performance Libraries.

- Operating System: Windows 11

### **3.2.3 Software Interfaces**

- Operating System: Windows 11
- IDE: Spyder
- Programming Language: Python

## **3.3 NON FUNCTIONAL REQUIREMENT**

### **3.3.1 Performance Requirements**

The performance of the functions and every module must be well. The overall performance of the software will enable the users to work decently. Performance of encryption of data should be fast. Performance of the providing virtual environment should be fast.

### **3.3.2 Safety Requirement**

The application is designed in modules where errors can be detected and fixed easily. This makes it easier to install and update new functionality if required.

### **3.3.3 Software Quality Attributes**

Our software has many quality attributes that are given below:

- Adaptability: This software is adaptable by all users.
- Availability: This software is freely available to all users. The availability of the software is easy for everyone.
- Maintainability: After the deployment of the project if any error occurs then it can be easily maintained by the software developer.
- Reliability: The performance of the software is better which will increase the reliability of the Software.
- User Friendliness: Since, the software is a GUI application; the output generated is much user friendly in its behavior.
- Integrity: Integrity refers to the extent to which access to software or data.
- Security: Users are authenticated using many security phases so reliable security is provided.
- Testability: The software will be tested considering all the aspects.

## **3.4 SYSTEM REQUIREMENTS**

### **3.4.1 Database Requirements**

SQLITE

### **3.4.2 Software Requirements (Platform Choice)**

- Operating system: Windows 7 or more.
- Coding Language: Python
- IDE: Spyder

### **3.4.3 Hardware Requirements**

- System: Intel I3 Processor and above.
- Hard Disk: 20 GB
- RAM: 8GB

### **3.5 ANALYSIS MODEL: SDLC MODEL TO BE APPLIED**

The software development cycle is a combination of different phases such as designing, implementing and deploying the project. These different phases of the software development model are described in this section. The SDLC model for the project development can be understood using the following figure. The chosen SDLC model is the waterfall model which is easy to follow and fits bests for the implementation of this project.

**Requirements Analysis:** At this stage, the business requirements, definitions of usecases are studied and respective documentations are generated.

**Design:** In this stage, the designs of the data models will be defined and different data preparation and analysis will be carried out.

**Implementation:** The actual development of the model will be carried out in this stage. Based on the data model designs and requirements from previous stages, appropriate algorithms, mathematical models and design patterns will be used to develop the agent's backend and front-end components.

**Testing:** The developed model based on the previous stages will be tested in this stage. Various validation tests will be carried out over the trained model.

**Deployment:** After the model is validated for its accuracy scores its ready to be deployed or used in simulated scenarios.

**Maintenance:** During the use of the developed solution various inputs/scenarios will been countered by the model which might affect the models overall accuracy. Or with passing time the model might not fit the new business requirements. Thus, the model must be maintained often to keep its desired state of operation.

### **3.6 SYSTEM IMPLEMENTATION PLAN**

The System Implementation plan table, shows the overall schedule of tasks compilation and time duration required for each task.

<b>Sr. No.</b>	<b>Name/Title</b>	<b>Start Date</b>	<b>End Date</b>
1	Preliminary Survey	22/08/2023	29/08/2023
2	Introduction and Problem Statement	12/09/2023	26/09/2023
3	Literature Survey	04/10/2023	10/10/2023
4	Software Requirement And Specification	11/10/2023	18/10/2023
5	System Design	19/10/2023	02/11/2023
6	Partial Report Submission	03/11/2023	07/11/2023
7	Architecture Design	19/01/2024	08/02/2024
8	Implementation	19/02/2024	03/03/2024
9	Deployment	10/03/2024	20/03/2024
10	Testing	01/04/2024	15/04/2024
11	Paper Publish	16/04/2024	21/04/2024
12	Report Submission	24/04/2024	03/05/2024

# **CHAPTER 4**

# **SYSTEM DESIGN**

#### 4.1 SYSTEM ARCHITECTURE:

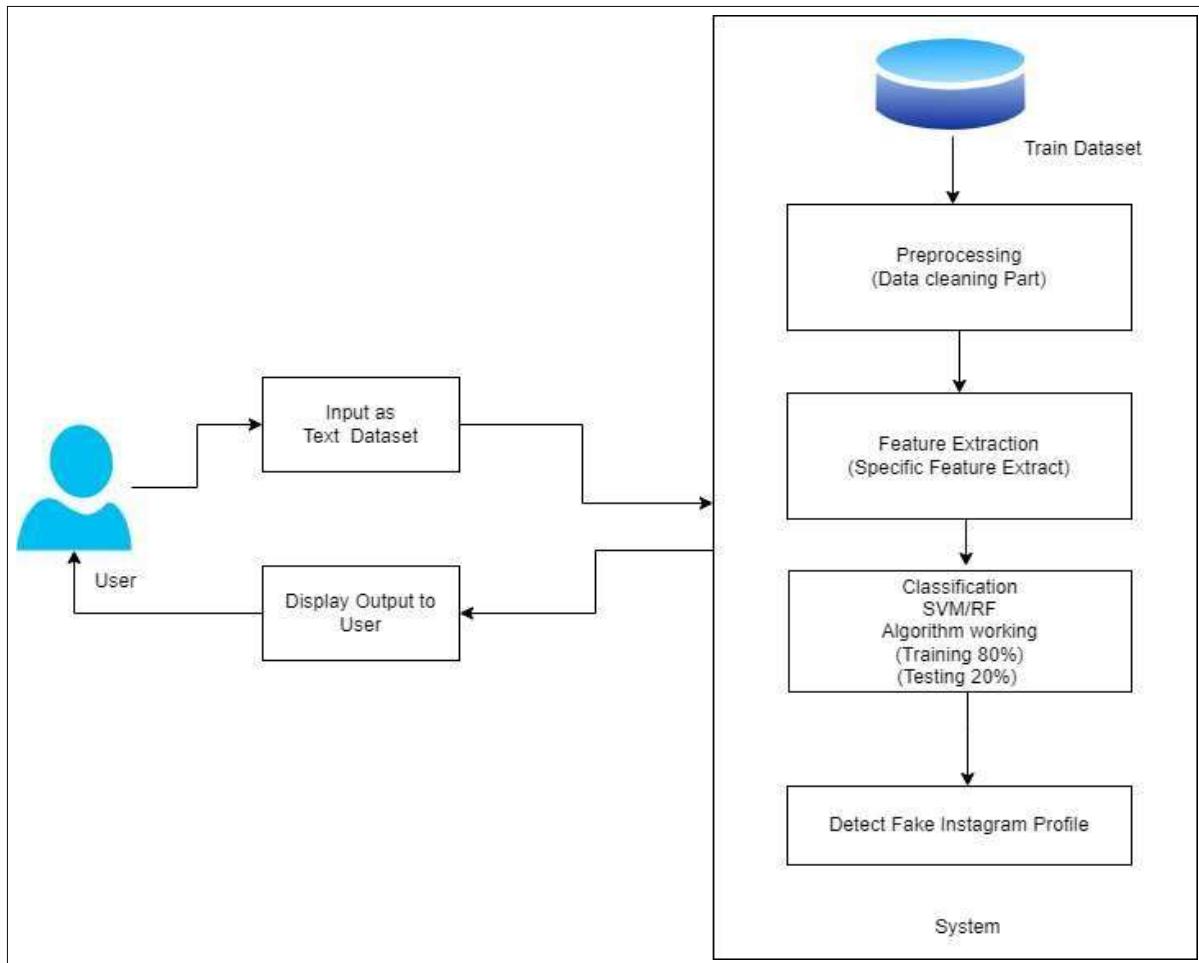


Figure 4.1: System Architecture

#### **4.1.1 Module**

- **Admin**

In this module, the Admin has to log in by using valid user name and password. After login successful he can do some operations such as View All Users and authorize.

- **View and Authorize Users**

In this module, the admin can view the list of users who all registered. In this, the admin can view the user's details such as, username, email, address and admin authorizes the users.

- **View Charts Results**

View All Products Search Ratio, View All Keywords, Search Results, View All Product Review Rank Results.

- **End User**

In this module, there are n numbers of users are present. User should register before doing any operations. Once user registers, their details will best or to the database. After registration successful, he has to login by using authorized user name and password. Once Login is successful user will do some operations like Manage Account.

#### **4.1.2 Data Flow Diagram**

In Data Flow Diagram, we show that flow of data in our system in DFD0 we show that base DFD in which rectangle presents input as well as output and circle shows our system.

In DFD1 we show actual input and actual output of system input of our system is text or image and output is rumor detected likewise in DFD 2 we present operation of user as well as admin.



Figure 4.2: Data Flow (0) diagram

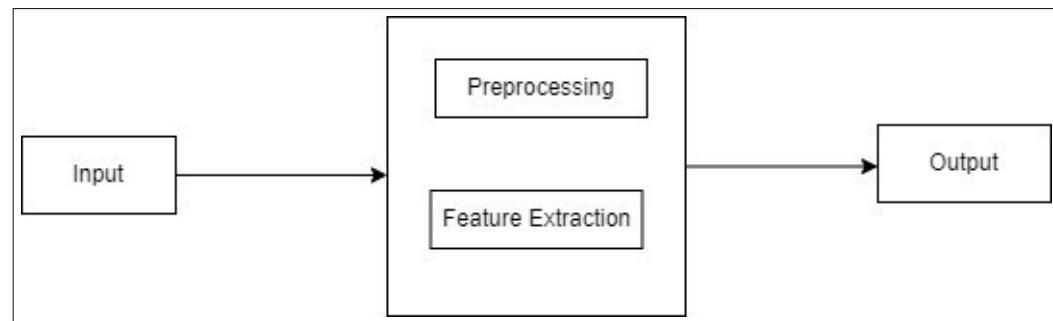


Figure 4.3: Data Flow (1) diagram

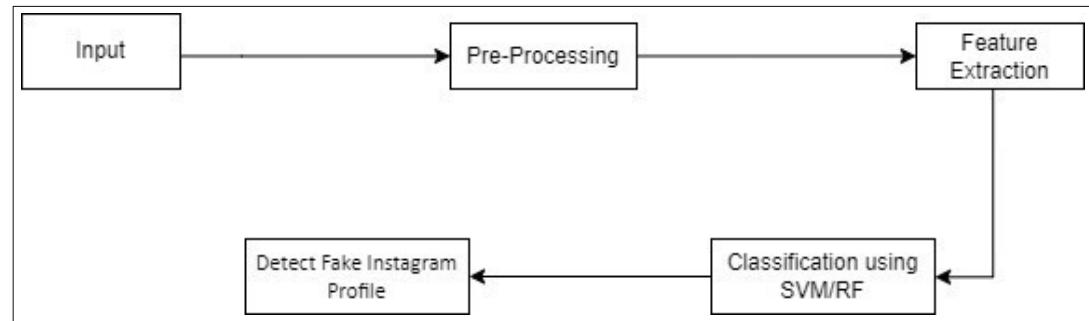


Figure 4.4: Data Flow (2) diagram

## 4.2 UML DIAGRAMS

Unified Modeling Language is a standard language for writing software blueprints. The UML may be used to visualize, specify, construct and document the artifacts of a software intensive system. UML is process independent, although optimally it should be used in process that is use case driven, architecture-centric, iterative and incremental. The Numbers of UML Diagram are available:

- Class Diagram
- Use case Diagram
- Activity Diagram
- Sequence Diagram

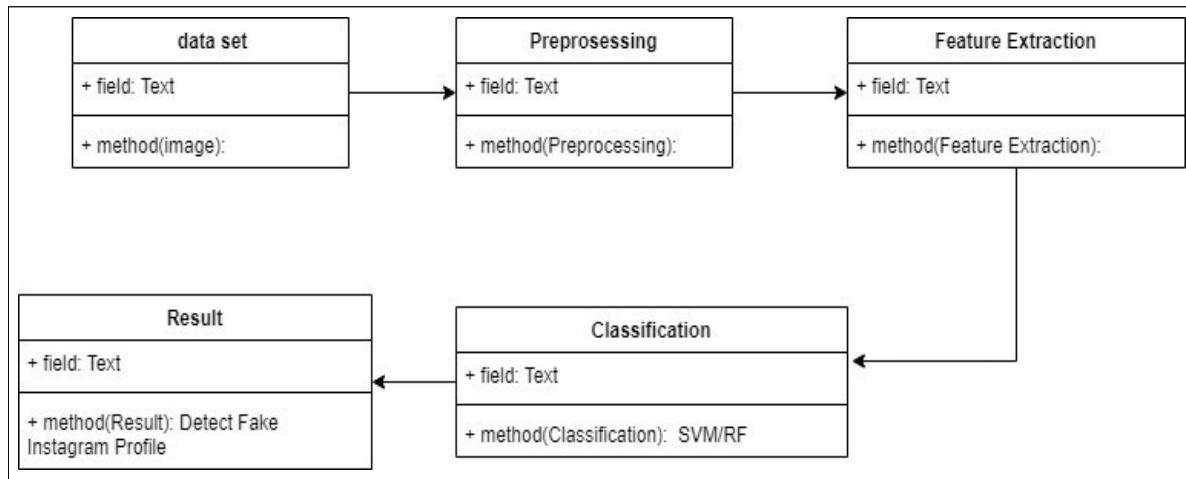


Figure 4.5: Class Diagram

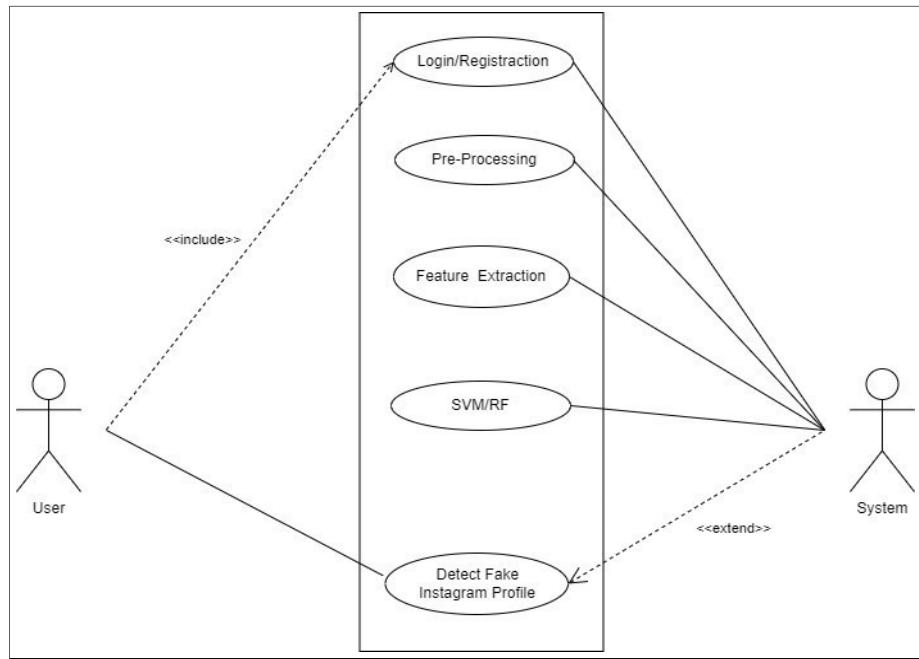


Figure 4.6: Usecase Diagram

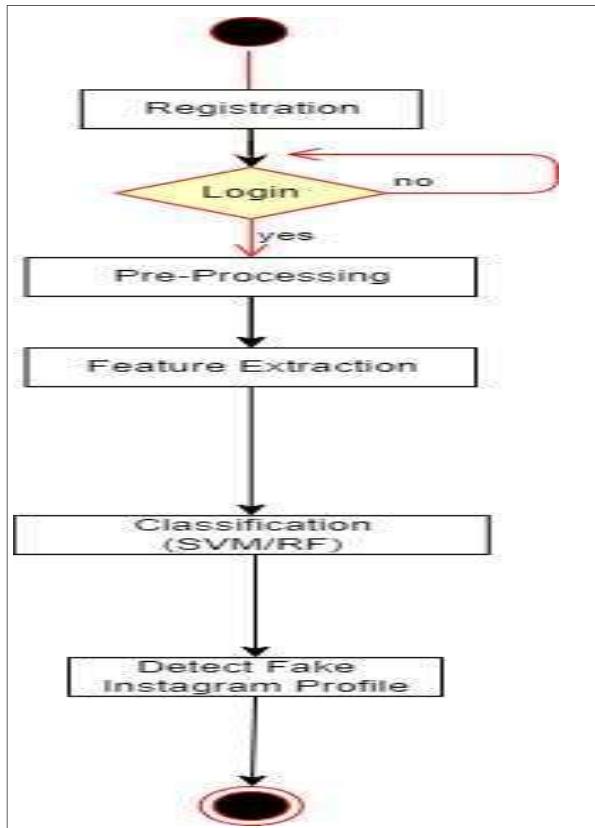


Figure 4.7: Activity Diagram

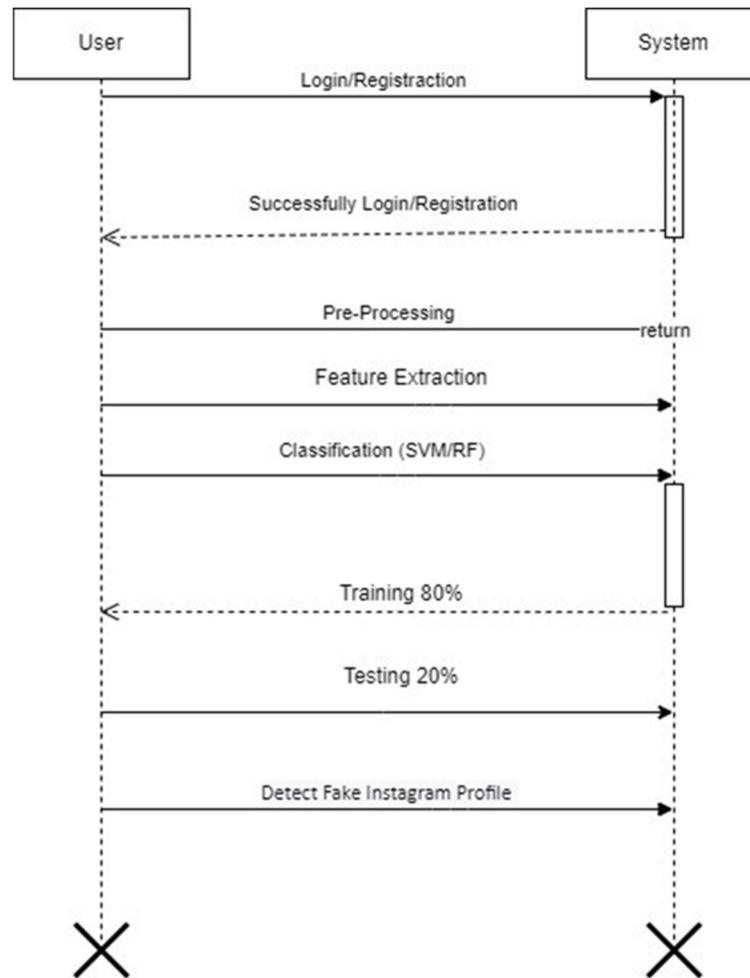


Figure 4.8: Sequence Diagram

# **CHAPTER 5**

# **SOFTWARE INFORMATION**

## 5.1 Python

Python is an interpreted, high-level and general-purpose programming language. Created by Guido van Rossum and first released in 1991, Python's design philosophy emphasizes code readability with its notable use of significant whitespace. Its language constructs and object-oriented approach aim to help programmers write clear, logical code for small and large-scale projects.

Python is dynamically typed and garbage-collected. It supports multiple programming paradigms, including structured (particularly, procedural), object-oriented, and functional programming. Python is often described as a “batteries included” language due to its comprehensive standard library.

Python was created in the late 1980s as a successor to the ABC language. Python 2.0, released in 2000, introduced features like list comprehensions and a garbage collection system with reference counting.

Python 3.0, released in 2008, was a major revision of the language that is not completely backward-compatible, and much Python 2 code does not run unmodified on Python 3.

The Python 2 language was officially discontinued in 2020 (first planned for 2015), and therefore the last Python 2 release.”[30] No more security patches or other improvements will be released for it. With Python 2’s end-of-life, only Python 3.6.x and later are supported.

## 5.2 Anaconda

Anaconda is a free and open-source distribution of the Python and R programming languages for scientific computing (data science, machine learning applications, large-scale data processing, predictive analytics, etc.), that aims to simplify package management and deployment. The distribution includes data-science packages suitable for Windows, Linux, and macOS. It is developed and maintained by Anaconda, Inc., which was founded by Peter Wang and Travis Oliphant in 2012. As an Anaconda, Inc. product, it is also known as Anaconda Distribution or Anaconda Individual Edition, while other products from the company are Anaconda Team Edition and Anaconda Enterprise Edition, both of which are not free.

Package versions in Anaconda are managed by the package management system conda. This package manager was spun out as a separate open-source package as it ended up being useful on its own and for other things than Python. There is also a small, bootstrap version of Anaconda called Miniconda, which includes only conda, Python, the packages they depend on, and a small number of other packages.

Anaconda distribution comes with over 250 packages automatically installed, and over 7,500 additional open-source packages can be installed from PyPI as well as the conda package and virtual environment manager. It also includes a GUI, Anaconda Navigator, as a graphical alternative to the command line interface (CLI). The big difference between conda and the pip package manager is in how package dependencies are managed, which is a significant challenge for Python data science and the reason conda exists.

When pip installs a package, it automatically installs any dependent Python packages without checking if these conflict with previously installed packages. It will install a package and any of its dependencies regardless of the state of the existing installation. Because of this, a user with a working installation of, for example, Google Tensorflow, can find that it stops working having used pip to install a different package that requires a different version of the dependent numpy library than the one used by Tensorflow. In some cases, the package may appear to work but produces different results in detail.

In contrast, conda analyses the current environment including everything currently installed and together with any version limitations specified (e.g. the user may wish to have Tensorflow version 2.0 or higher), works out how to install a compatible set of dependencies, and shows a warning if this cannot be done.

Open source packages can be individually installed from the Anaconda repository, Anaconda Cloud ([anaconda.org](http://anaconda.org)), or the user's own private repository or mirror, using the conda install command. Anaconda, Inc. compiles and builds the packages available in the Anaconda repository itself, and provides binaries for Windows 32/64bit, Linux 64 bit and MacOS 64-bit. Anything available on PyPI may be installed into a conda environment using pip and conda will keep track of what it has installed itself and what pip has installed. Custom packages can be made using the conda build command, and can be shared with others by uploading them to Anaconda Cloud, PyPI or other repositories.

The default installation of Anaconda2 includes Python 2.7 and Anaconda3 includes Python 3.7. However, it is possible to create new environments that include any version of Python packaged with conda.

## 5.3 Spyder

Spyder is a powerful scientific environment written in Python, for Python, and designed by and for scientists, engineers and data analysts. It offers a unique combination of the advanced editing, analysis, debugging, and profiling functionality of a comprehensive development tool with the data exploration, interactive execution, deep inspection, and beautiful visualization capabilities of a scientific package.

Beyond its many built-in features, its abilities can be extended even further via its plugin system and API. Furthermore, Spyder can also be used as a PyQt5 extension library, allowing you to build upon its functionality and embed its components, such as the interactive console, in your own software.

### 5.3.1 Features

- **Editor**

Works efficiently in a multi-language editor with a function/class browser, real-time code analysis tools (pyflakes, pylint, and pycodestyle), automatic code completion (jedi and rope), horizontal/vertical splitting, and go-to-definition.

- **Interactive Console**

Harness the power of as many IPython consoles as you like with full workspace and debugging support, all within the flexibility of a full GUI interface. Instantly run your code by line, cell, or file, and render plots right inline with the output or in interactive windows.

- **Documentation Viewer**

Render documentation in real-time with Sphinx for any class or function, whether external or user-created, from either the Editor or a Console. It is very useful for documenting viewing.

- **Variable Explorer**

Inspect any variables, functions or objects created during your session. Editing and interaction is supported with many common types, including numeric/strings/bools, Python lists/tuples/dictionaries, dates/timedeltas, Numpy arrays, Pandas in- dex/series/dataframes, PIL/Pillow images, and more.

- **Development Tools**

Examine your code with the static analyzer, trace its execution with the interactive debugger, and unleash its performance with the profiler. Keep things organized with project support and a built-in file explorer, and use find in files to search across entire projects with full regex support.

### 5.3.2 Other Specifications

- **Advantages**

Machine learning models can efficiently process large volumes of data, making it feasible to analyze numerous profiles and interactions in real-time. Machine learning enables the automation of fake profile identification, reducing the reliance on manual inspections and user reporting. The proposed approach incorporates both content-based and network-based attributes for classification. Machine learning algorithms are capable of learning intricate patterns that distinguish genuine from fake profiles.

- **Limitations**

The effectiveness of machine learning models heavily relies on the quality and quantity of available data. Developing effective algorithms to capture both content-based and network based attributes accurately requires careful consideration and domain expertise

# **CHAPTER 6**

# **SOFTWARE TESTING**

## **6.1 INTRODUCTION**

Software testing, depending on the testing method employed, can be implemented at any time in the development process. However, most of the test effort occurs after the requirements have been defined and the coding process has been completed. As such, the methodology of the test is governed by the software development methodology adopted. Different software development models will focus the test effort at different points in the development process. Newer development models, such as Agile, often employ test driven development and place an increased portion of the testing in the hands of the developer, before it reaches a formal team of testers.

## **6.2 TYPES OF TESTING USED**

Involves various testing strategies used for the project.

### **6.2.1 UNIT TESTING**

It is the testing of individual software units of the application. It is done after the completion of an individual unit before integration. Unit testing involves the design of test cases that validate that the internal program logic is functioning properly, and that program inputs produce valid outputs. All decision branches and internal code flow should be validated. Unit tests perform basic tests at component level and test a specific business process, application, and system configuration.

### **6.2.2 INTEGRATION TESTING**

Integration tests are designed to test integrated software components to determine if they actually run as one program. Testing is event driven and is more concerned with the basic outcome of screens or fields.

Integration tests demonstrate that although the components were individually satisfactory, as shown. Integration testing is specifically aimed at exposing the problems that arise from the combination of components.

### **6.2.3 WHITE-BOX TESTING & BLACK-BOX TESTING**

Software testing methods are traditionally divided into white and black-box testing. These two approaches are used to describe the point of view that a test engineer takes while designing test cases:

#### **1. White-Box Testing**

In white-box testing an internal perspective of the system, as well as programming skills, are used to design test cases.

#### **2. Black-Box Testing**

Black-box testing treats the software as a examining functionality without any knowledge of internal implementation. The testers are only aware of what the software is supposed to do, not how it does it.

### 6.3 TEST CASES TEST RESULTS

Test Case ID	Test Case	Test Case I/P	Actual Result	Expected Result	Test case criteria(P/F)
001	Store Xml File	Xml file	Xml file store	Error Should come	P
002	Parse the xml file for conversion	parsing	File get parse	Accept	P
003	Attribute identification	Check individual Attribute	Identify Attributes	Accepted	P
004	Weight Analysis	Check Weight	Analyze Weight of individual Attribute	Accepted	P
005	Tree formation	Form them-Tree	Formation	Accepted	P
006	Cluster Evaluation	Check Evaluation	Should check Cluster	Accepted	P
007	Algorithm Performance	Check Evaluation	Should work Algorithm Properly	Accepted	P
008	Query Formation	Check Query Correction	Should check Query	Accepted	P

Figure 6.1: GUI TESTING

Test Case ID	Test Case	Test Case I/P	Actual Result	Expected Result	Test case criteria(P/F)
001	Enter the number in username, middle name, last name field	Number	Error Comes	Error Should Comes	P
001	Enter the character in username, middle name, last name field	Character	Accept	Accept	P
002	Enter the invalid email id format in email id field	Kkgmail.com	Error comes	Error Should Comes	P
002	Enter the valid email id format in email id field	kk@gmail.com	Accept	Accept	P
003	Enter the invalid digit no in phone no field	99999	Error comes	Error Should Comes	P
003	Enter the 10 digit no in phone no field	9999999999	Accept	Accept	P

Figure 6.2: Registration Test Case

Test Case ID	Test Case	Test Case I/P	Actual Result	Expected Result	Test case criteria(P/F)
001	Enter The Wrong username or password click on submit button	Username or password	Error comes	Error Should come	P
002	Enter the correct username and password click on submit button	Username and password	Accept	Accept	P

**Figure 6.3: Login Test Case**

## **CHAPTER 7**

### **RESULT**

## 7.1 Login Page

The login page in the project allows users to access the system where machine learning algorithms analyze and classify Instagram profiles, distinguishing between fake and genuine accounts based on various data inputs.

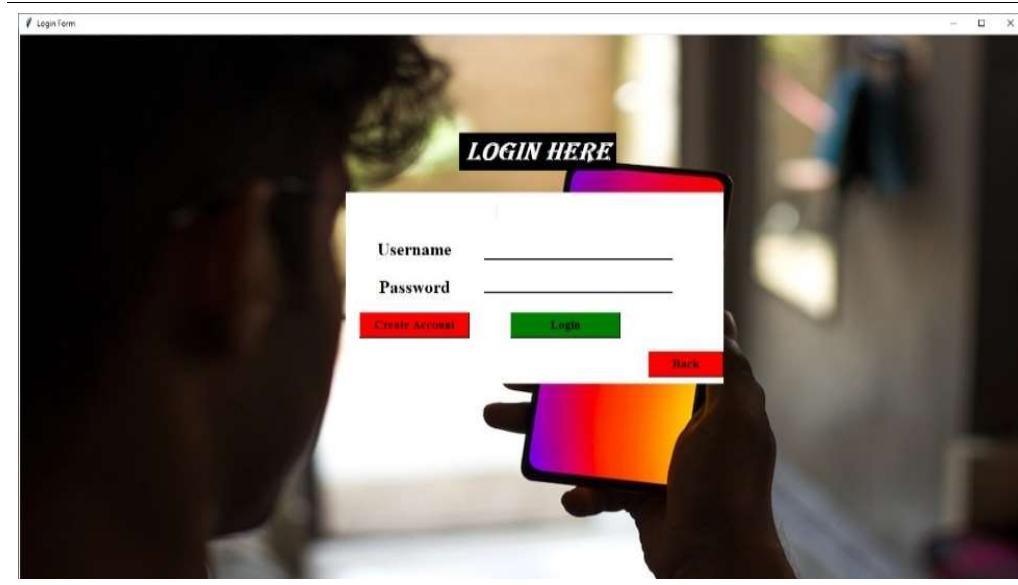
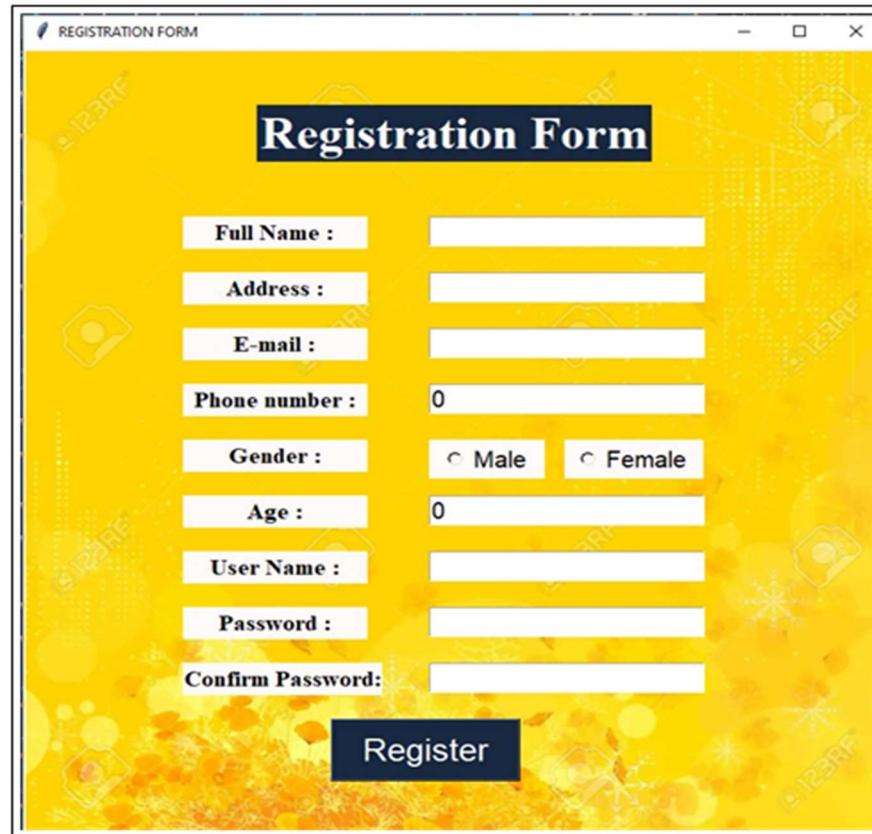


Figure 7.1: Login Page

## 7.2 Registration Page

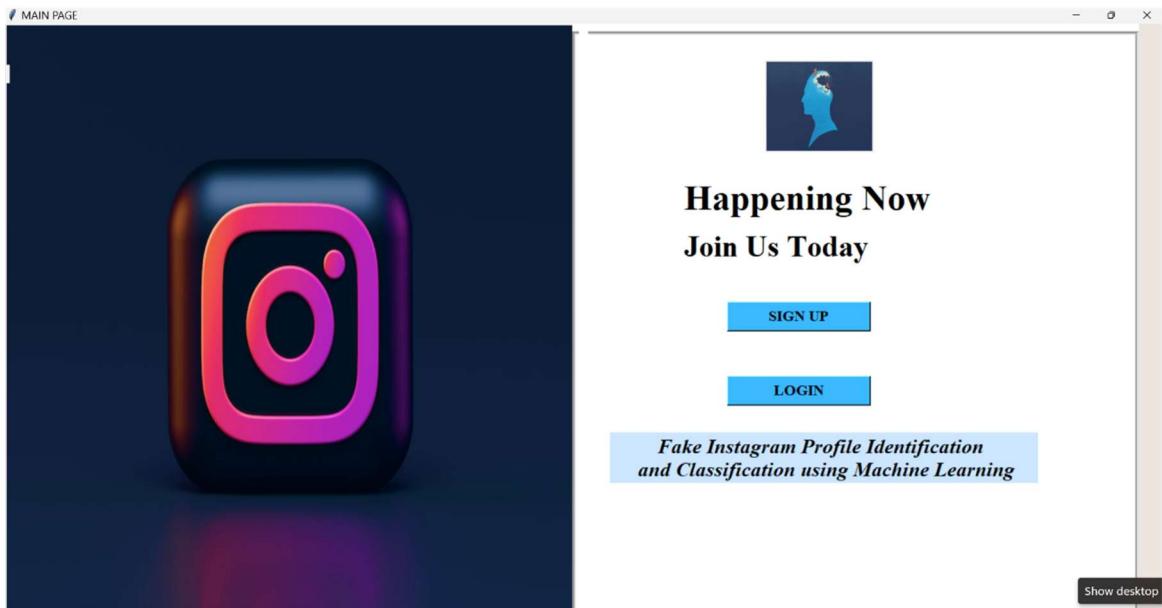
The registration page captures essential user details such as name, email, password, and sometimes additional verification information. This data is used to authenticate users, store their preferences, and secure their interactions within the system. After registration, users can access the machine learning tools that analyze and classify Instagram profiles to detect and identify fake accounts. The process ensures that only authorized users can utilize the system's capabilities, enhancing security and user management.



**Figure 7.2: Registration Page**

### 7.3 GUI Main

The GUI Main page is the primary interface of the Fake Instagram Profile Identification and Classification system. It features a user-friendly dashboard where users can Sign Up or Log in. This central hub ensures easy interaction with the system's machine learning features, providing a seamless user experience for effective profile analysis.



**Figure 7.3: GUI Main**

## 7.4 Outputs

### 7.4.1 Page1 & Page2

These pages contain various attributes like followers, following, username etc. of an Instagram profile whose authenticity to be known. Also based upon these attributes and classification using ML algorithms it displays messages like **Account\_Not\_Fake** (if it's not fake) & **Fake\_Account** (if it's fake).

Instagram Fake Profile Detection	
<b>profilepic</b>	1
<b>numsLengthusername</b>	0.27
<b>fullnamewords</b>	1
<b>numsLengthfullname</b>	0
<b>nameUsername</b>	0
<b>descriptionlength</b>	0
<b>private</b>	0
<b>posts</b>	0
<b>followers</b>	45
<b>follows</b>	64
<b>Submit</b>	
<b>Account_Not_Fake</b>	

Figure 7.4.1: Output Page1

Instagram Fake Profile Detection	
<b>profilepic</b>	1
<b>numsLengthusername</b>	0.33
<b>fullnamewords</b>	1
<b>numsLengthfullname</b>	0.33
<b>nameUsername</b>	1
<b>descriptionlength</b>	30
<b>private</b>	1
<b>posts</b>	35
<b>followers</b>	488
<b>follows</b>	604
<b>Submit</b>	
<b>Fake_Account</b>	

Figure 7.4.1: Output Page2

#### 7.4.2 Page3, Page4 & Page5

These 3 pages include detailed performance metrics for three algorithms: SVM (Support Vector Machine), RF (Random Forest), and DT (Decision Tree). These metrics comprise precision, recall, F-score, and support, displayed to help users evaluate the accuracy and effectiveness of each algorithm. Also there is an option to check authenticity of the profile manually by clicking on “check” button.

Accuracies of the algorithms used:

1. Support Vector Machine: 95.83%
2. Random Forest: 97.22%
3. Decision Tree: 95.83%



**Figure 7.4.2: Output Page3**

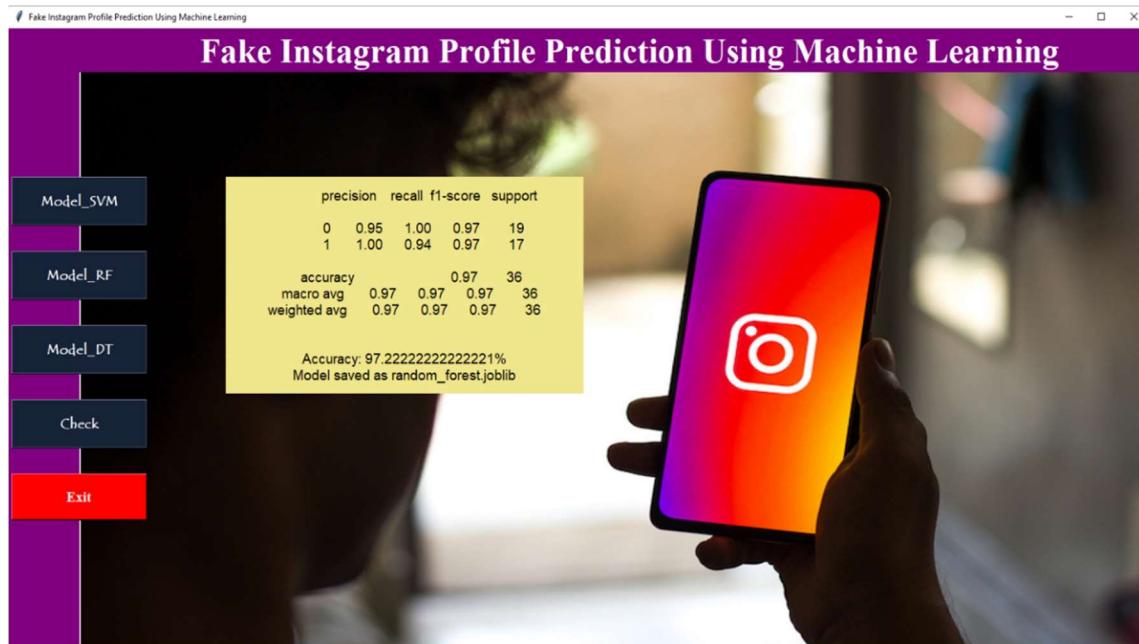


Figure 7.4.2: Output Page4

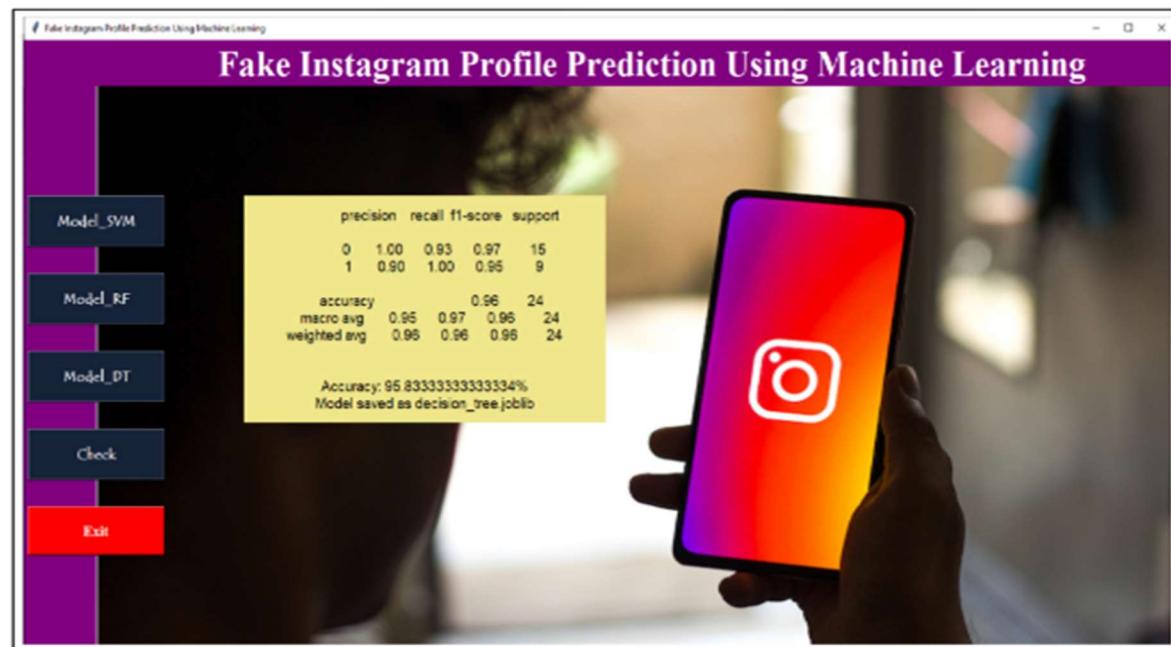


Figure 7.4.2: Output Page5

# **CHAPTER 8**

# **CONCLUSION**

The research on “Fake Instagram Profile Identification and Classification using Machine Learning” presents a comprehensive approach to tackle the persistent issue of fake profiles on social media platforms, with a specific focus on Instagram. By leveraging the power of machine learning techniques, this research contributes to creating a safer and more trustworthy online environment for users, bolstering user confidence, and upholding the integrity of social media community. The research’s outcomes extend beyond the realm of academia, impacting the lives of individuals, businesses, and society as a whole. As social media continues to shape the digital landscape, the work presented here contributes to building a foundation of trust and authenticity, reinforcing the positive potential of online interactions and collaborations.

## **CHAPTER 9**

## **REFERENCES**

1. Aleksei Romanov, Alexander Semenov, Oleksiy Mazhelis and Jari Veijalainen.2017. "Detection of Fake Profiles in Social Media". In 13th International Conference on Web Information Systems and Technologies.
2. Indira Sen,Anupama Aggarwal,Shiven Mian.2018."Worth its Weight in Likes: Towards Detecting Fake Likes on Instagram". In ACM International Conference on Information and Knowledge Management.Nazir, Atif, Saqib Raza, Chen-Nee Chuah, Burkhard Schipper, and C. A. Davis. "Ghostbusting Face- book: Detecting and Characterizing Phantom Profiles in Online Social Gam- ing Applications." In WOSN. 2010.
3. Nambouri Sravya, Chavana Sai praneetha, S. Saraswathi," Identify the Human or Bots Twitter Data using Machine Learning Algorithms", International Re- search Journal of Engineering and Technology (IRJET), Volume: 06 Issue: 03Mar 2019 www.irjet.net, e-ISSN: 2395-0056, p- ISSN: 2395-0072.
4. M. Smruthi, N. Harini," A Hybrid Scheme for Detecting Fake Accounts in Facebook", International Journal of Recent Technology and Engineering (IJRTE) ISSN: 2277-3878, Volume-7, Issue-5S3, February 2019.
5. Nazir, Atif, Saqib Raza, Chen-Nee Chuah, Burkhard Schipper, and C. A.Davis. "Ghostbusting Facebook: Detecting and Characterizing Phantom Pro- files in Online Social Gaming Applications." In WOSN. 2010.

6. Rao, P. S., J. Gyani, and G. Narsimha. "Fake profiles identification in online social networks using machine learning and NLP." *Int. J. Appl. Eng. Res* 13.6 (2018): 973-4562.
7. Raturi, Rohit. "Machine learning implementation for identifying fake ac- counts in social network." *International Journal of Pure and Applied Mathematics* 118.20 (2018): 4785-4797. J. Wang, "Fundamentals of erbium-doped fibre amplifiers arrays (Periodical style—Submitted for publication)," *IEEE J. Quantum Electron.*, submitted for publication.
8. M. Mohammadrezaei, M. E. Shiri, and A. M. Rahmani, "Identifying fake ac- counts on social networks based on graph analysis and classification algo- rithms," *Security and Communication Networks*, vol. 2018, 2018.
9. Ala'M, Al-Zoubi, Ja'far Alqatawna, and Hossam Faris. "Spam profile detec- tion in social networks based on public features." *2017 8th International Con- ference on information and Communication Systems (ICICS)*. IEEE, 2017.
10. Romanov, Aleksei, Alexander Semenov, Oleksiy Mazhelis, and Jari Veijalainen. "Detection of fake profiles in social media-Literature review." In *International Conference on Web Information Systems and Technologies*, vol. 2, pp. 363-369. SCITEPRESS, 2018.

# **APPENDIX**

## **APPENDIX A**

### **Type of Problem: NP (Non Deterministic Polynomial Time) Class**

**Definition:** Problems for which a proposed solution can be verified quickly in polynomial time by a deterministic approach are called as NP Class problems.

#### **Explanation:**

The problem of Fake Instagram Profile Identification and Classification using Machine Learning can be categorized as an NP (Nondeterministic Polynomial time) class problem for the following reasons:

#### **1. Verification in Polynomial Time:**

Once a machine learning model (e.g., SVM, Random Forest, Decision Tree) is trained, it produces a classification for each Instagram profile (either fake or genuine).

#### **2. Decision Problem Aspect:**

For any given instance (profile and model), we can quickly verify if the model's decision is correct by comparing it with the ground truth. This quick verification aligns with the NP class definition.

#### **3. Non-Deterministic Polynomial Time:**

Imagine a hypothetical non-deterministic machine that guesses the classification (fake or genuine) for each profile. We could then verify these guesses in polynomial time by comparing them against known labels.

# APPENDIX B

## 1. Paper Name: FAKE INSTAGRAM PROFILE IDENTIFICATION

**Journal Name: International Research Journal of Modernization in Engineering Technology and Science**

**Research Paper :**



### ABSTRACT

Rise of fake profiles on platforms like, Instagram possess significant challenges related to user privacy, security, and trust. This work presents a challenging task to identify and classify fake Instagram profiles using machine learning techniques. The findings of this research contribute to the ongoing efforts to combat the issue of fake profiles on Instagram and other social media platforms. By leveraging machine learning techniques and other comprehensive feature set, the proposed model demonstrates promising results in identifying and classifying fake profiles, thereby promoting a safer and more trustworthy online environment. This research opens avenues for further exploration, including the integration of real-time data streams and the adaptation of the model to other social media platforms.

**Keywords:** Fake Profile Identification, User Authentication, Data Preprocessing, Model Training, Online Security, Machine Learning.

### I. INTRODUCTION

Counterfeit Instagram profiles encompass a spectrum from automated bots inundating feeds with spam to adept impostors aiming to deceive authentic users for various ulterior motives such as financial exploitation or social engineering. Traditional methods reliant on manual scrutiny and reporting prove inadequate in handling the overwhelming influx of profiles and engagements, thus necessitating the adoption of sophisticated technological interventions. Machine learning emerges as a potent instrument in combatting the proliferation of fake profiles across social media platforms. Leveraging the computational prowess of machine learning algorithms, it becomes feasible to automatically discern and categorize fake profiles based on discernible patterns and attributes. The amalgamation of the burgeoning influence of social media, the formidable challenges posed by counterfeit profiles, and the strides made in machine learning methodologies culminate in the formulation of solutions aimed at detecting and classifying these deceptive entities.

This research endeavors to address the imperative for a more secure and credible online milieu by proposing a holistic methodology to confront the issue of fake Instagram profiles through the application of machine learning.

### II. PROPOSED SYSTEM

Developing a reliable algorithm capable of identifying counterfeit profiles on Instagram poses a formidable yet crucial endeavor. Instagram, akin to numerous other social networking platforms, grapples with the prevalence of fraudulent accounts perpetrating spam, deceit, or nefarious deeds. Below is a conceptual framework outlining a potential system designed to detect bogus profiles on Instagram.

#### 1. Data Acquisition:

Accumulate a substantial dataset of Instagram profiles, encompassing both authentic and fraudulent accounts, showcasing a diverse array of attributes.

#### 2. Feature Extraction:

Derive pertinent characteristics from user profiles. These attributes may entail:

- Examination of profile images: Scrutinize for substandard quality, incongruities, or recurrent usage.
- Analysis of activity patterns: Evaluate the frequency of posts, likes, comments, and followers.
- Evaluation of follower-to-following ratios: Identify extreme imbalances.
- Assessment of post content: Scrutinize posts for spam.



e-ISSN: 2582-5208

**International Research Journal of Modernization in Engineering Technology and Science**

( Peer-Reviewed, Open Access, Fully Refereed International Journal )

Volume:06/Issue:02/February-2024

Impact Factor- 7.868

[www.irjmets.com](http://www.irjmets.com)

### **3. Machine Learning Framework:**

Construct a machine learning framework to categorize profiles as genuine or counterfeit. Consider employing methodologies such as:

- Adoption of deep learning architectures, like neural networks, adept at capturing intricate patterns.
- Utilization of ensemble techniques such as Random Forest or Gradient Boosting to enhance accuracy.

### **4. Model Training and Validation:**

Segment the dataset into distinct subsets for training, validation, and testing purposes. Train the model using the training data, refining hyperparameters as necessary.

- Validate the model's efficacy using the validation set, iterating as required.

### **5. Continuous Surveillance:**

Implement the algorithm to operate seamlessly in real-time across Instagram profiles. Maintain ongoing vigilance over user activity, profiles, and engagements to find early on.

#### **A. Support Vector Machine :**

Support Vector Machines (SVMs) represent a prevalent form of supervised machine learning algorithm utilized for tasks involving classification and regression. SVMs exhibit versatility, capable of handling both linear and non-linear classification scenarios.

Working Steps of Support Vector Machine:

- Step 1: Load essential libraries required for implementation.
- Step 2: Acquire the dataset and segregate the independent variables (X) and dependent variable (Y).
- Step 3: Partition the dataset into distinct training and testing subsets.
- Step 4: Initialize the SVM classifier model.
- Step 5: Train the SVM classifier model by fitting it to the training data.
- Step 6: Generate predictions using the trained model.
- Step 7: Assess the performance.

#### **B. Random Forest :**

Random Forest serves as an ensemble learning approach extensively employed in machine learning for tasks encompassing classification and regression. This method amalgamates the forecasts of numerous decision trees to enhance predictive precision and resilience.

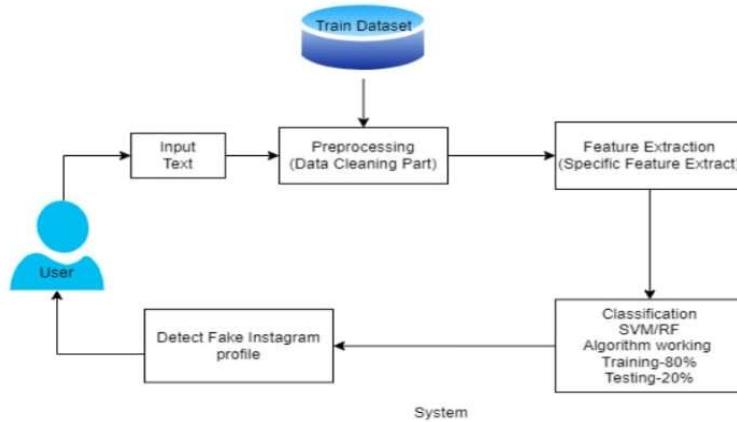
The working steps for Random Forest:

- Step 1: Import and preprocess the data.
- Step 2: Train the random forest classifier.
- Step 3: Evaluate the accuracy of predictions.
- Step 4: Visualize the outcomes derived from the classifier.

## **III. SYSTEM DESIGN**

Crafting a system architecture for an Instagram counterfeit profile detection algorithm necessitates the incorporation of various components, each serving distinct roles and functionalities. Below are the constituent elements of the architecture:

1. Data Collection Layer
2. Preprocessing Layer
3. Feature Extraction
4. Classification Algorithms
5. Implementation



**Fig. 1 System Architecture**

This system architecture provides a comprehensive view of how different components work together to detect fake profiles on Instagram. The key to success is the continuous refinement of the machine learning model, real-time monitoring, and the ability to adapt to evolving strategies used by malicious actors.

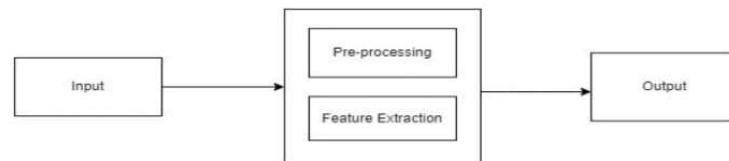
#### Data Flow Diagram :

A Data Flow Diagram (DFD) illustrates the movement of data within a system. DFD0 serves as the foundational diagram depicting input and output as rectangles, with the system represented by a circle. DFD1 elaborates on the actual inputs and outputs of the system, such as textual or image data input and rumor detection output. Conversely, DFD2 delineates the actions performed by users and administrators within the system.

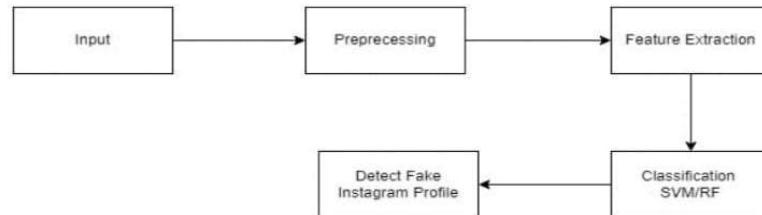
#### DFD0 :



#### DFD1:



#### DFD2:



**Fig. 2 Data Flow Diagram**

#### IV. BACK UP PLAN

If obtaining enough labeled data proves difficult, explore methods like data augmentation to generate synthetic data, aiding in enhancing model effectiveness. Establish a system enabling users to flag suspicious profiles,

which, when combined with automated detection, enhances accuracy. Employ ensemble models by merging various machine learning algorithms to boost precision and resilience. Collaborate closely with legal experts and contemplate enforcing stricter regulations for account validation and creation. Prioritize user privacy and ethical considerations, particularly when handling user data and profile details, ensuring compliance with guidelines.

## V. CONCLUSION

The research on Fake Instagram Profile Identification and Classification using Machine Learning presents a comprehensive approach to tackle the persistent issue of fake profiles on social media platforms, with a specific focus on Instagram. By leveraging the power of machine learning techniques, this research contributes to creating a safer and more trustworthy online environment for users, bolstering user confidence, and upholding the integrity of social media community. The research's outcomes extend beyond the realm of academia, impacting the lives of individuals, businesses, and society as a whole. As social media continues to shape the digital landscape, the work presented here contributes to building a foundation of trust and authenticity, reinforcing the positive potential of online interactions and collaborations.

## VI. REFERENCES

- [1] Zainab Agha, Neeraj Chatlani, Afsaneh Razi, and Pamela Wisniewski. 2020. Towards conducting responsible research with teens and parents regarding online risks. In Extended Abstracts of the 2020 CHI Conference on Human Factors.
- [2] Shiza Ali, Afsaneh Razi, Seunghyun Kim, Ashwaq Alsoubai, Joshua Gracie, Munmun De Choudhury, Pamela J Wisniewski, and Gianluca Stringhini. 2022. Understanding the Digital Lives of Youth: Analyzing Media Shared within Safe Versus Unsafe Private Conversations on Instagram. (2022), 1-14.
- [3] S. M. Din, R. Ramli and A. A. Bakar, "A Review on Trust Factors Affecting purchase Intention on Instagram", 2018 IEEE Conference on Application Information and Network Security (AINS), 2018.
- [4] Detect fake profiles on social media network.
- [5] S.C. Boerman, "The effects of the standardized Instagram disclosure for micro- and meso-influencers", Computers in Human Behavior, vol. 103, pp. 199-207, 2020.
- [6] S. Sheikhi, An Efficient Method for Detection of Fake Accounts on the Instagram Platform, 2020.
- [7] J. Kang and L. Wei, "Let me be at my funniest: Instagram users' motivations for using Finsta (a.k.a. fake Instagram)", The Social Science Journal, 2019.
- [8] M. Mondal, L. A. Silva and F. Benevenuto, "A Measurement Study of Hate Speech in Social Media", Proceedings of the 28th ACM Conference on Hypertext and Social Media - HT '17, 2017.
- [9] B. Mathew, R. Dutt, P. Goyal and A. Mukherjee, "Spread of Hate Speech in Online Social Media", Proceedings of the 10th ACM Conference on Web Science.
- [10] H. Hilal Bashir and S. A. Bhat, "Effects of Social Media on Mental Health: A Review", The International Journal of Indian Psychology.

## Certificates:





## **2. Paper Name: FAKE INSTAGRAM PROFILE IDENTIFICATION & CLASSIFICATION USING ML**

**Journal Name: International Research Journal of Modernization in Engineering Technology and Science**

**Research Paper :**



e-ISSN: 2582-5208

International Research Journal of Modernization in Engineering Technology and Science

( Peer-Reviewed, Open Access, Fully Refereed International Journal )

Volume:06/Issue:03/March-2024

Impact Factor- 7.868

[www.irjmets.com](http://www.irjmets.com)

### **FAKE INSTAGRAM PROFILE IDENTIFICATION USING ML**

**Aniruddha Lalge<sup>\*1</sup>, Apurv Badave<sup>\*2</sup>, Nikhil Elajale<sup>\*3</sup>, Pankaj Godara<sup>\*4</sup>,**

**Prof. Priyanka Kinage<sup>\*5</sup>**

<sup>\*1,2,3,4,5</sup>Department Of Computer Engineering Smt. Kashibai Navale College Of Engineering  
Pune, India.

#### **ABSTRACT**

The pervasive rise of social media has become an undeniable aspect of modern life, with its global dominance steadily increasing. However, this growth has also brought about a host of ecosystem challenges, including the proliferation of hate speech, fraudulent activities, and the spread of fake news. These issues, compounded by the staggering number of over 1.7 billion fake accounts across social media platforms, have already resulted in significant losses, and the process of removing these accounts remains daunting and time-consuming. With Instagram's user base expanding rapidly, there's a growing urgency to address the issue of identifying fake accounts on this platform. Traditionally, manual identification processes are laborious and time-intensive.

**Keywords:** Fake Profile Identification, User Authentication, Data Preprocessing, Model Training, Online Security, Machine Learning.

#### **I. INTRODUCTION**

Online Social Networks (OSNs) such as Facebook and Instagram have become increasingly pervasive in modern society, playing a vital role in communication and serving as platforms for personal and business promotion. Initially, an account's popularity is often gauged through metrics like follower count and engagement indicators such as likes, comments, and views on shared content. Consequently, users may resort to artificial means to inflate these metrics for personal gain. Common methods include the use of bots, purchasing social metrics like likes and followers, and utilizing platforms for metric trading. Notably, a significant number of Instagram accounts are automated, and bots have been known to generate more internet traffic than humans. Leveraging a comprehensive dataset comprising various features extracted from genuine and fraudulent accounts for employing advanced algorithms and feature engineering methodologies, we endeavor to construct a robust framework for automated detection, thereby enhancing the integrity and reliability of the Instagram ecosystem. Through rigorous experimentation and evaluation, our findings underscore the efficacy of the proposed system in differentiating between genuine and fake profiles, offering a promising solution to mitigate the growing menace of online deception.

#### **II. DATASET & INPUT**

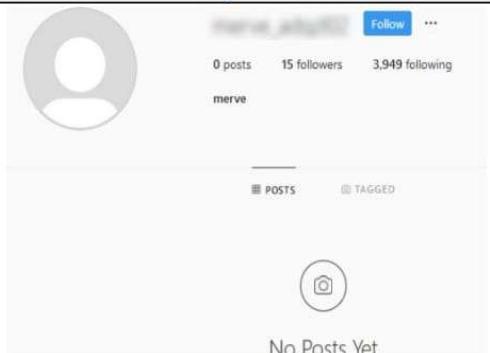
The detection of fake accounts on social media platforms is crucial due to their role in artificially inflating the popularity metrics of other users. These accounts often exhibit characteristics such as a high number of accounts followed coupled with a low number of followers, erratic liking behavior, and identifiable traits like the absence of a profile picture and unusual usernames.

##### **A. Features Of Dataset**

To facilitate the detection of fake accounts, a dataset comprising of numerous real accounts & fake accounts was meticulously curated through manual labeling. Various factors were considered during data collection, including follower and following counts, media posting frequency, comments on media, and profile characteristics such as the presence of a profile picture and the format of the username.

In the dataset, the selected base features can be listed as below:

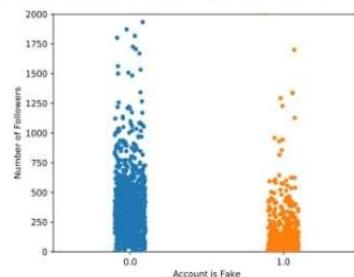
- Total media number of the account.
- Follower count of the account.
- Following count of the account.
- Number of digits present in account username.
- Whether account is private, or not (binary feature).



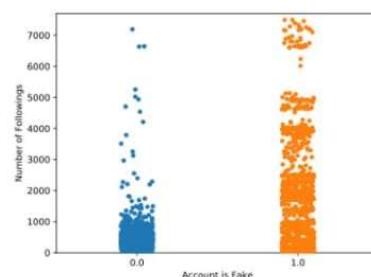
**Fig 1:** Fake account example from dataset.

#### B. Oversampling

Distribution of classes in the fake account dataset is not even. This results in poor performance for the outnumbered class. SMOTE oversampling technique is utilized to increase number of samples for fake accounts. K is chosen as 5 for this work. In the implementation of SMOTE, SMOTE-NC is applied which considers not only the quantity classes but also the categorical classes. After applying oversampling, all classifiers are trained on equal number of training samples per class.



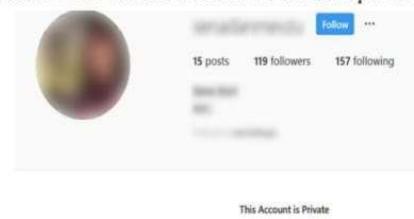
**Fig 2:** In-class data distributions for "following count" feature.



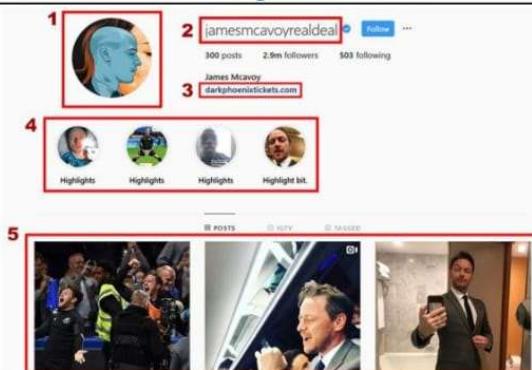
**Fig 3:** In-class data distributions for "follower count" feature.

### III. FAKE ACCOUNT DETECTION

This section delves into the identification of automated accounts, commonly referred to as bots, which engage in automated activities such as following, liking, and commenting on content. These actions are typically targeted towards specific hashtags, locations, or followers of particular accounts with the aim of boosting their own popularity metrics. Automated accounts may exhibit behavior that is entirely inorganic or a combination of organic and inorganic actions. It's worth noting that the presence of organic behavior in such accounts can occur when users continue to pursue their own interests while the bot operates in the background.



**Fig 4:** Example private account preview. Media details are not visible for private accounts.



**Fig 5:** General preview of an Instagram profile.

1: Profile picture

2: Username

3: External URL

4: Highlight reel

5: User media.

#### A. Dataset Features & Methodology

To compile the authentic accounts, we selected individuals from our personal networks, primarily from our circle of friends. In contrast, to gather the automated accounts, we analyzed the source codes of the most widely used open-source Instagram bots, which are instrumental in generating artificial engagement. We identified specific behaviors characteristic of these bots to accurately label the accounts.

Hashtags serve as a prominent avenue for identifying inorganic activity. For instance, a common method for detecting fake likes on Instagram involves scrutinizing the usage of hashtags associated with like and follow trading. In our approach, we concentrated on targeting the most popular hashtags, as it was observed that this method facilitated the detection of automated behavior more efficiently and rapidly. If a user is found to engage in activities such as following and unfollowing within predefined time intervals, as determined by parameters from online Instagram automation tools, they are categorized as automated accounts. We utilized the Instagram API alongside a Python wrapper to meticulously gather comprehensive media and user information from the accounts over a six-month period. To safeguard user privacy, any personally identifiable information such as usernames, photos, comments, and hashtag details were deliberately excluded from the dataset.

Below are the scrapped base features from these accounts:

- Total media number of the accounts.
- Follower count of the account.
- Following count of the account.
- Number of photos user is tagged by someone else.
- Average recent media hashtag number.

If the account has no media, all features scrapped from user posts are assigned as 0. Additionally helpful features are derived using the base features such as:

- Average recent media like to comment ratio (LCR).
- Follower to following ratio (FFR).
- Whether account has not any media, or not.

#### B. Bias Problem

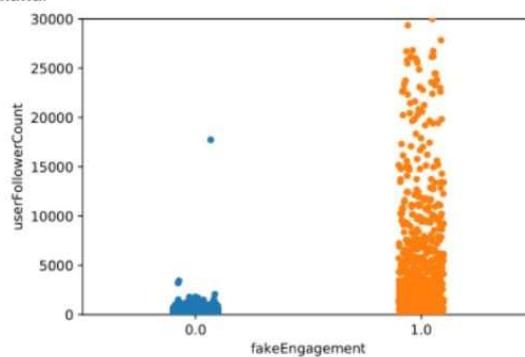
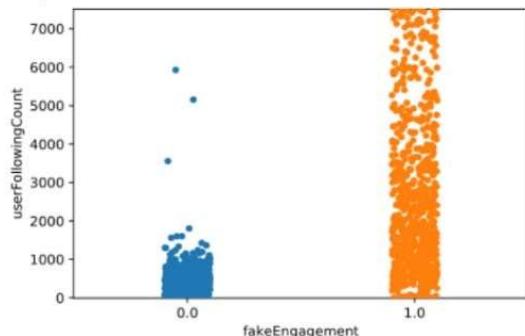
There existed some unfavorable bias within certain features of the dataset. Figures 6-7 display the distribution of the entire dataset concerning selected continuous features, where a label of "Fake Engagement" equal to 1

denotes accounts involved in fake engagement or automated behavior, while 0 corresponds to accounts with only genuine engagements or real accounts. These figures reveal bias within the dataset across the chosen features. While the bias in follower and following numbers appears unrealistic (likely due to unintentionally categorizing accounts with low follower and following counts as real accounts, which does not accurately reflect real-world scenarios), the bias in average hashtag usage per post appears more natural. Accounts engaged in automated behavior tend to employ more hashtags per post, contributing to this observed bias.

The distribution of the dataset is projected across selected binary features. These tables highlight the presence of bias within these features as well, although this time, the biases appear to be more in line with reality. For instance, the presence of highlight reels can be seen as an effective differentiator, as real engagement accounts predominantly lack this feature. Similarly, the absence of a URL in the profile can be considered a genuine bias, as it aligns with typical user behavior.

### C. Cost Sensitive Feature Selection

To address these unrealistic biases and identify the most impactful features, a novel approach employing a cost-sensitive genetic feature selection algorithm has been devised. The algorithm, outlined in Algorithm 1, begins by normalizing the continuous features while leaving the binary features unchanged. Subsequently, the normalized data is inputted into the genetic algorithm for feature selection. In this algorithm, an individual is represented as an array with a length equal to the total number of features in the dataset. Each element of the array corresponds to whether the feature is selected or not, with a value of 1 indicating selection and 0 indicating exclusion. For instance, if the second element of the individual is 1, it signifies that the second feature has been selected for inclusion in that specific individual. This representation allows for the formation of a population using randomly generated individuals. Each individual in the population represents a unique combination of selected features, with the goal of optimizing the selection process to identify the most effective features for the task at hand.


**Fig 6:** In-class data distributions for "follower count" feature.

**Fig 7:** In-class data distributions.

$$\text{Fitness} = F2 \text{ Score} - 2 \times \text{Tot.Feat.Cost} \quad (1)$$

$$\text{Precision} = \text{TP}/(\text{TP} + \text{FP}) \quad (2)$$

$$\text{Recall} = \text{TP}/(\text{TP} + \text{FN}) \quad (3)$$

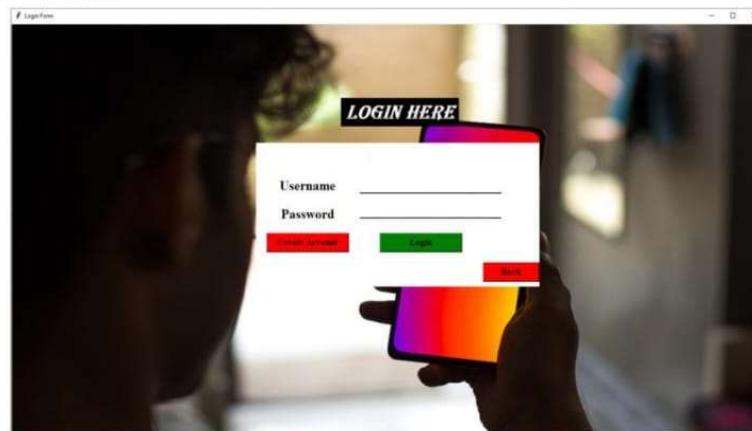
$$F1 \text{ Score} = 2 \times \text{Recall} \times \text{Precision} / (\text{Recall} + \text{Precision}) \quad (4)$$

Fitness calculation formula is given in Eq. 1 Tot.Feat.Cost is calculated by summing the individual costs of the selected features. These costs are determined based on the reliability of the data collection which is discussed in the previous Bias Problem section. Realistic biases are represented with lower features costs while the negative bias is represented with higher costs. We conducted a comprehensive evaluation of several machine learning algorithms suitable for binary classification tasks. These algorithms included logistic regression, random forest, support vector machines (SVM), gradient boosting machines (GBM), random forest (RM), decision tree (DT) etc. Each algorithm was evaluated based on its performance in terms of accuracy, precision, recall, and F1 score using cross-validation techniques. Based on the results of the evaluation, we selected the algorithm that demonstrated the highest performance across multiple metrics while considering computational efficiency. Our decision was guided by the need to balance model accuracy with practical considerations such as deployment feasibility and resource requirements. The trained model was evaluated on the validation set to estimate its performance metrics, including accuracy, precision, recall & F1 score. The training process involved optimizing the model's parameters to achieve the best possible performance on the training data. We employed techniques such as grid search or random search for hyperparameter tuning, ensuring that the model generalized well to unseen data while avoiding overfitting.

#### IV. RESULT

By employing a variety of algorithms, the objective is to leverage different facets of the dataset, including independence, separability, and complex relationships, which have not been extensively explored in existing literature. The goal is to develop an effective method for detecting fake and automated Instagram accounts. To detect automated accounts, the selected features identified through cost-sensitive feature selection are utilized. In contrast, for identifying fake accounts, the base features from the fake-real dataset are directly employed.

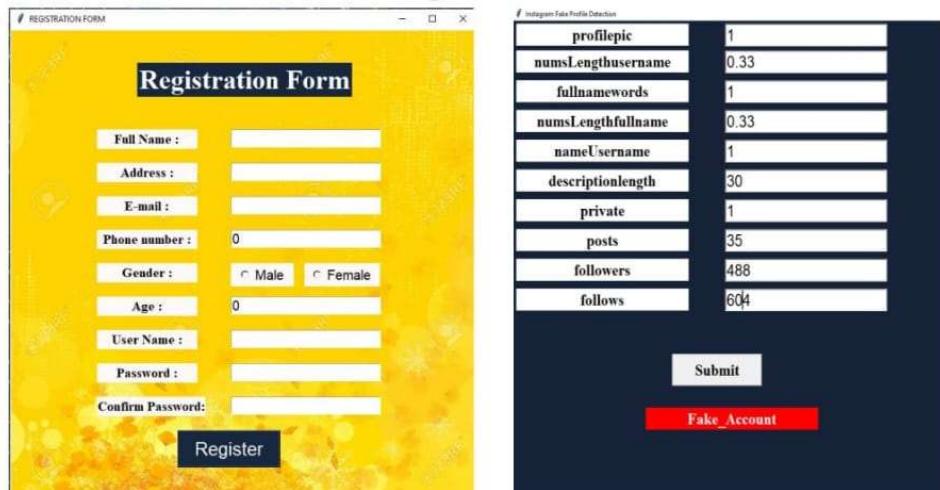
In evaluating the effectiveness of the implemented techniques, Precision, Recall, and F1 Score are utilized as evaluation metrics, as outlined in equations 3-5, respectively. These metrics provide insights into the performance of the detection algorithms. Specifically, Precision measures the proportion of true automated accounts among all accounts identified as automated, while Recall calculates the proportion of true automated accounts correctly identified from all actual automated accounts. The F1 Score, which is the harmonic mean of Precision and Recall, offers a comprehensive assessment of the algorithms' performance by considering both false positives and false negatives. More meaningful metric for evaluating overall performance compared to Precision or Recall alone.



**Fig 8:** Login Page



Fig 9: GUI Main



	Instagram_Fake_Profile_Detection
profilepic	1
numsLengthusername	0.33
fullnamewords	1
numsLengthfullname	0.33
nameUsername	1
descriptionlength	30
private	1
posts	35
followers	488
follows	604

Fig 10: Registration Page

Fig 11: Output Page 1



Fig 12: Output Page 2



e-ISSN: 2582-5208

**International Research Journal of Modernization in Engineering Technology and Science**  
( Peer-Reviewed, Open Access, Fully Refereed International Journal )

Volume:06/Issue:03/March-2024

Impact Factor - 7.868

[www.irjmets.com](http://www.irjmets.com)

## V. CONCLUSION

The study titled "Fake Instagram Profiles Identification Using Machine Learning" offers a comprehensive strategy to address the persistent challenge of fraudulent profiles on social media platforms, particularly Instagram. Through the application of advanced machine learning techniques, this research endeavors to enhance the safety and credibility of the online environment, thereby fostering user trust and maintaining the integrity of social media communities.

The outcomes of this research transcend academic boundaries, exerting a tangible impact on individuals, businesses, and society at large. By leveraging the power of machine learning techniques, this research contributes to creating a safer and more trustworthy online environment for users, bolstering user confidence, and upholding the integrity of social media community. The research's outcomes extend beyond the realm of academia, impacting the lives of individuals, businesses, and society as a whole. As social media continues to shape the digital landscape, the work presented here contributes to building a foundation of trust and authenticity, reinforcing the positive potential of online interactions and collaborations.

## VI. REFERENCES

- [1] Shiza Ali, Afsaneh Razi, Seunghyun Kim, Ashwaq Alsoubai, Joshua Gracie, Munmun De Choudhury, Pamela J Wisniewski, and Gianluca Stringhini 2022. Understanding the Digital Lives of Youth: Analyzing Media Shared within Safe Versus Unsafe Private Conversations on Instagram. (2022), 1-14.
- [2] S. M. Din, R. Ramli and A. A. Bakar, "A Review on Trust Factors Affecting purchase Intention on Instagram", 2018 IEEE Conference on Application Information and Network Security (AINS), 2018.
- [3] Detect fake profiles on social media network.
- [4] S.C. Boerman, "The effects of the standardized Instagram disclosure for micro- and meso-influencers", Computers in Human Behavior, vol. 103, pp. 199-207, 2020.
- [5] S. Sheikhi, An Efficient Method for Detection of Fake Accounts on the Instagram Platform, 2020.
- [6] J. Kang and L. Wei, "Let me be at my funniest: Instagram users' motivations for using Fake insta", The Social Science Journal, 2019.
- [7] B. Mathew, R. Dutt, P. Goyal and A. Mukherjee, "Spread of Hate Speech in Online Social Media", Proceedings of the 10th ACM Conference on Web Science.
- [8] Y. Li, O. Martinez, X. Chen, J. E. Hopcroft, "In a world that counts: Clustering and detecting fake social engagement at scale,"
- [9] Proceedings of the 25th International Conference on World Wide Web. International World Wide Web Conferences Steering Committee, 2016.
- [10] T. Clarke, "22+ Instagram Stats That Marketers Can't Ignore This Year".

## Certificates:





## APPENDIX C

### Plagiarism Report: Plagiarism Report of Published Paper.

 May 06, 2024

#### Plagiarism Scan Report

7% Plagiarized | 93% Unique

Characters: 2408 | Words: 348  
Sentences: 14 | Speak Time: 3 Min

Excluded URL: None

#### Content Checked for Plagiarism

It gives us great pleasure in presenting the preliminary project report on 'Fake Instagram Profile Identification and Classification using Machine Learning'. We would like to convey our gratitude to our Principal and all the teaching and non-teaching staff members of the Computer Engineering Department who gave us the freedom to explore and guided us the right way, also our friends and families for their valuable suggestions and support.

Fake Instagram profiles can range from automated bots posting spam to sophisticated imposters attempting to deceive genuine users for financial gain, social manipulation, or other illicit activities. Traditional methods of manual inspection and reporting are insufficient to handle the sheer volume of profiles and interactions, necessitating the use of advanced technological solutions. Machine learning has emerged as a powerful tool in addressing the issue of fake profiles on social media platforms. By harnessing the computational power of machine learning algorithms, it is possible to automatically identify and classify fake profiles based on distinctive patterns and characteristics. The combination of the growing influence of social media, the challenges posed by fake profiles, and the advancements in machine learning techniques has led to the development of solutions aimed at identifying and classifying these profiles. This research addresses the need for a safer and more trustworthy online environment by proposing a comprehensive approach to tackle the issue of fake Instagram profiles using machine learning. Research addresses the need for a safer and more trustworthy online environment by proposing a comprehensive approach to tackle the issue of fake Instagram profiles using machine learning. It gives us great pleasure in presenting the preliminary project report on 'Fake Instagram Profile Identification and Classification using Machine Learning'. It gives us great pleasure in presenting the preliminary project report on 'Fake Instagram Profile Identification and Classification using Machine Learning'. It gives us great pleasure in presenting the preliminary project report on 'Fake Instagram Profile Identification and Classification using Machine Learning'.

#### Sources

7% Plagiarized

Page 1 of 2

Fake Instagram profiles can range from automated bots posting spam to sophisticated imposters attempting to deceive genuine users for financial gain, social manipulation, or other...

[https://www.researchgate.net/profile/Apurv-Badave/publication/375112869\\_Fake\\_Instagram\\_Profile\\_Identification\\_and\\_Classification\\_using\\_Machine\\_Learning/links/654128fa0426ef6369ede3df/Fake-Instagram-Profile-Identification-and-Classification-using-Machine-Learning.pdf?origin=publication\\_detail/](https://www.researchgate.net/profile/Apurv-Badave/publication/375112869_Fake_Instagram_Profile_Identification_and_Classification_using_Machine_Learning/links/654128fa0426ef6369ede3df/Fake-Instagram-Profile-Identification-and-Classification-using-Machine-Learning.pdf?origin=publication_detail/)



[Home](#)   [Blog](#)   [Testimonials](#)   [About Us](#)   [Privacy Policy](#)

Copyright © 2024 [Plagiarism Detector](#). All rights reserved

Page 2 of 2