

# INSTAGRAM FAKE VS GENUINE ACCOUNT DETECTION USING MACHINE LEARNING

---

## Internship Project Report

**Submitted by:** Monu Ramkesh Gupta

**Course / Internship:** Data Analyst Internship

**Tools:** Python, Pandas, NumPy, Scikit-learn, Jupyter Notebook, Power BI

---

## Abstract

This project aims to detect fake Instagram accounts using machine learning techniques. The study focuses on analyzing user account data such as followers, follows, posts, and description length to classify accounts as either fake or genuine. By applying Random Forest classification, we achieved an accuracy of approximately 94%. The project demonstrates the potential of data analytics in detecting inauthentic online activity and can serve as a foundation for social media monitoring systems.

## 1. Introduction

With the exponential rise of social media platforms, fake profiles have become a significant issue. These accounts often spread misinformation or manipulate engagement metrics. The goal of this project is to identify patterns that differentiate fake accounts from genuine ones using data analytics and machine learning. By analyzing metrics like followers, follows, posts, and profile descriptions, a prediction model is built to classify users accurately.

## 2. Objectives

- To identify fake and genuine Instagram accounts using data analytics and machine learning.
- To analyze user profile features such as followers, follows, and post count.
- To train and evaluate models for classification.
- To visualize data patterns and model results using Power BI.

## 3. Dataset Description

The dataset contains 576 Instagram account records with multiple attributes such as:

- profile pic
- username length
- fullname words
- name==username
- description length
- external URL
- private status
- #posts
- #followers
- #follows
- fake (target column)

A new feature 'ratio\_followers' was created as  $(\text{\#followers} / (\text{\#follows} + 1))$  to understand following patterns.

## 4. Methodology

The project was completed in the following steps:

1. Data collection and preprocessing.
2. Feature engineering – created 'ratio\_followers' to improve prediction accuracy.
3. Model training using Random Forest Classifier.
4. Model evaluation using accuracy score and classification report.
5. Visualization in Power BI to present results and insights.

## 5. Model Building

A Random Forest Classifier was implemented to train the model. The features used were:

- ratio\_followers
- #followers
- #follows
- #posts
- description length

After training, the model achieved 94.4% accuracy.

#### Classification Report:

Precision, Recall, and F1-Score values were 0.94 for both classes (Fake and Genuine), indicating a balanced and effective model.

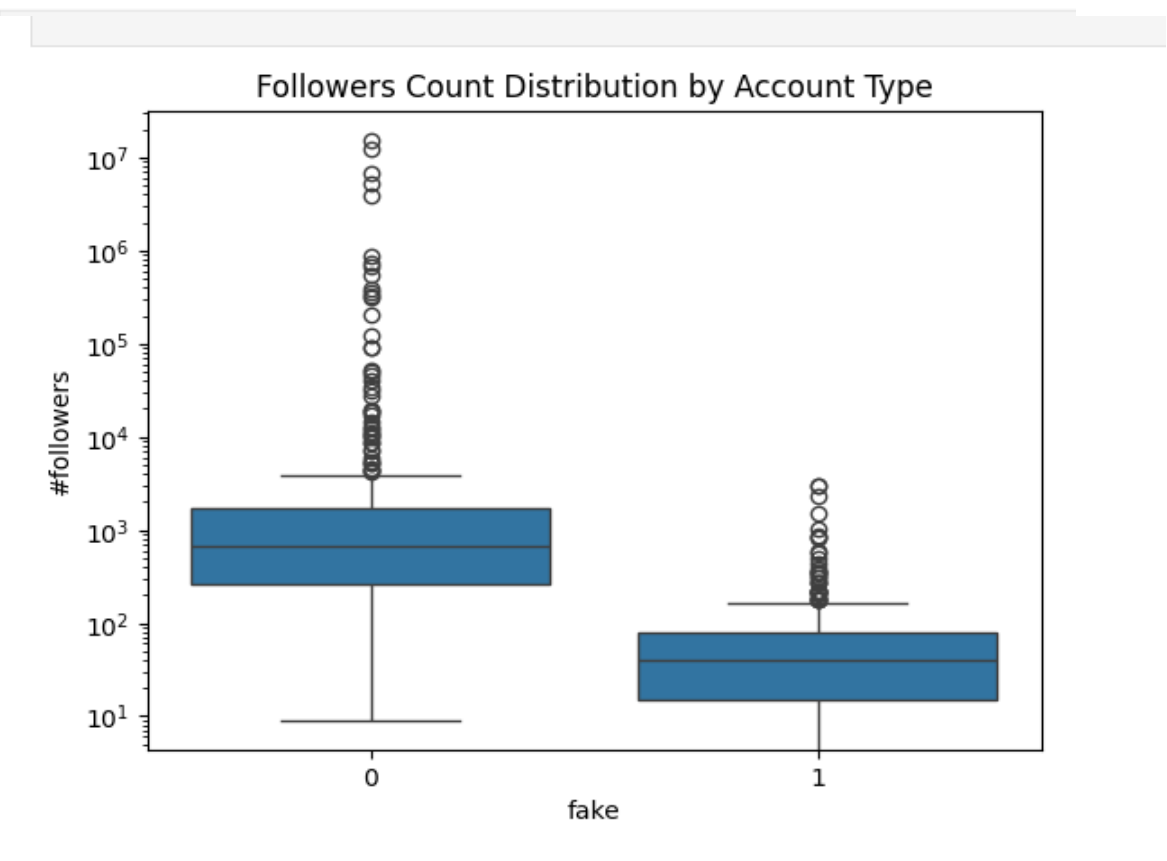
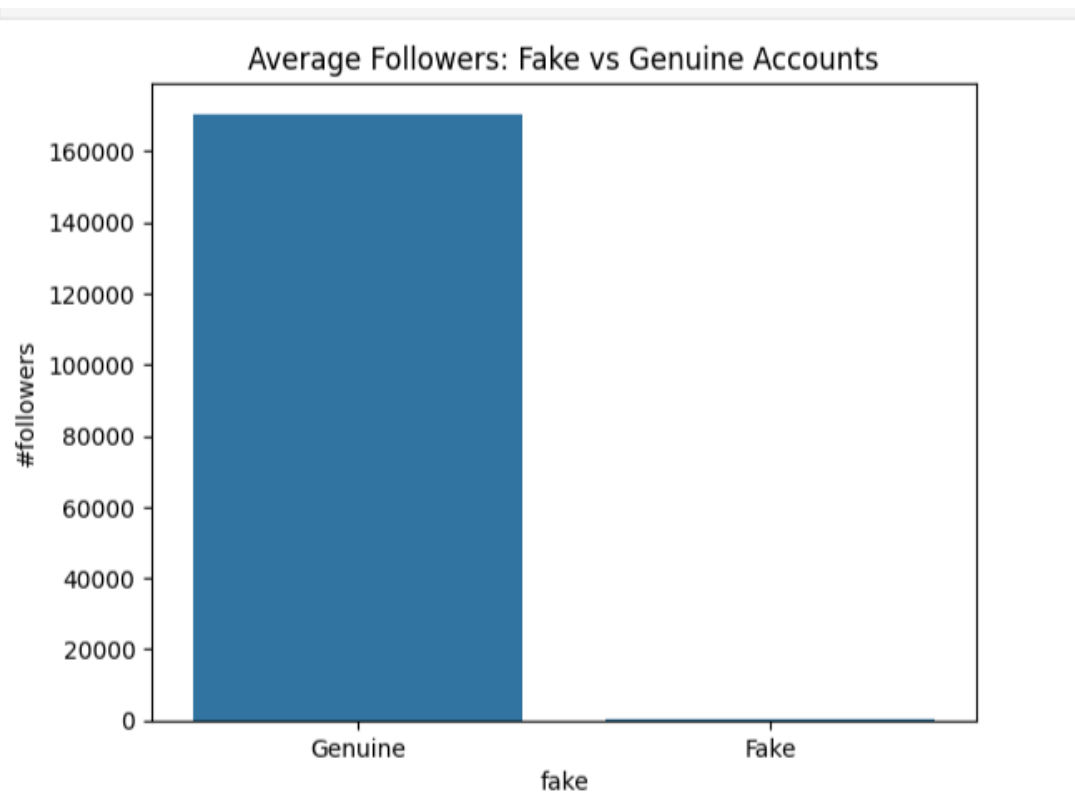
#### Feature Importance:

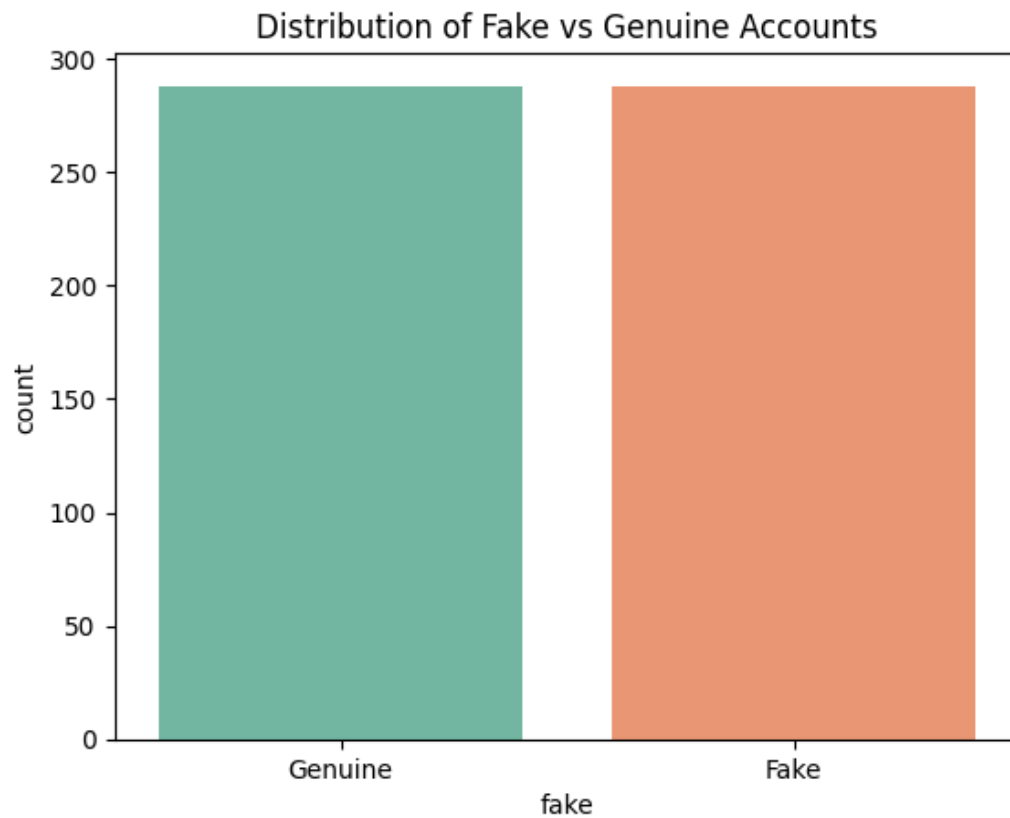
- #posts – 30.35%
- description length – 24.96%
- #followers – 17.22%
- ratio\_followers – 14.80%
- #follows – 12.66%

## 6. Data Analysis and Visualization

The analysis and visualization part was performed in Power BI to better understand the dataset and display the prediction insights. The following visuals were created:

1. Fake vs Genuine Count (Bar Chart)
2. Average Followers, Posts, and Follows (Cards)
3. Followers vs Follows Comparison (Stacked Column Chart)
4. Ratio\_Followers Distribution (Histogram or Bar Chart)
5. Feature Importance Representation (Bar Chart)





accuracy of y predict 0.9248554913294798

#### Classification report:

	precision	recall	f1-score	support
0	0.94	0.91	0.93	93
1	0.90	0.94	0.92	80
accuracy			0.92	173
macro avg	0.92	0.93	0.92	173
weighted avg	0.93	0.92	0.92	173

	Feature	Importance
1	#followers	0.315508
3	#posts	0.243270
0	ratio_followers	0.173121
5	profile pic	0.144099
2	#follows	0.108468
4	private	0.015535

## Dashboard Visualization

To provide a clear and interactive view of the findings, a Power BI dashboard was designed. The dashboard presents key insights from the analysis of Instagram fake vs genuine accounts through multiple visuals, including:

- **Total Fake vs Genuine Accounts (%)**
- **Average Followers, Follows, and Posts (by account type)**
- **Followers–Follows Relationship (Clustered Column Chart)**
- **Followers-to-Follows Ratio Comparison**
- **Correlation between Followers and Posts**
- **Account Engagement Overview**

The dashboard enables easy comparison between fake and genuine profiles, helping identify behavioral differences visually.



## 7. Results

The Random Forest Classifier performed well with an overall accuracy of 94%. The classification report shows balanced precision and recall for both fake and genuine accounts. The most important features for classification were #posts and description length.

## 8. Conclusion

This project demonstrates that fake accounts can be effectively detected by analyzing behavioral and profile-based features. The Random Forest algorithm provided high accuracy and interpretability. Visualization using Power BI enhanced the understanding of data trends and helped present insights clearly.

## 9. References

1. YouTube tutorials on Data Cleaning and Random Forest in Python.
2. Google – Documentation on scikit-learn, pandas, and Power BI.
3. UM- dataset for Instagram fake account detection.
4. Analytics Vidhya – Guides on data visualization best practices.