

Fake Profile Identification using Machine Learning Algorithms

M.Mamatha¹, M.Srinivasa Datta², Umme Hani Ansari³ Dr. Subhani Shaik⁴

Students, Dept. of IT, Sreenidhi Institute of Science and Technology (A), Hyderabad-501301

Associate Professor, Dept. of IT, Sreenidhi Institute of Science and Technology (A), Hyderabad-501301

ABSTRACT

The main theme of our paper is identifying whether the instagram profile is genuine or fake. Algorithms will be trained with all previous users fake and genuine account data and then whenever we give new test data then trained model will be applied on new test data to identify whether given new account details are from genuine or fake users. The machine learning-based methods were used to perceive false accounts that could give the wrong impression about people. The dataset is pre-processed using a variety of python libraries and a comparison form is obtained to get a realistic algorithm appropriate for the specified dataset. An effort to notice forged accounts on the social media platforms is strong-minded by a variety of machine learning algorithms. The performances of the classification algorithms Random Forest, Network and support vector machines are used for the detection of fake accounts.

Keywords - Fake Profile Identification, Machine Learning, Random Forest, Neural Network, SVM

Date of Submission: 05-07-2021

Date of Acceptance: 18-07-2021

I. INTRODUCTION

Nowadays, Instagram is dominating many ways social networking sites. Every day the amount of user's victimization of social media drastically change. The most advantage of Instagram is that we are able to hook up with individuals simply, share photos, videos and communicate with them in an exceedingly higher means. Several on-line searching stores and event planners use instagram as a media to urge quality and to sell their product and services. This provided a replacement means of a possible attack, like faux uniqueness, fake data, etc. A recent survey counsel that the amount of accounts gift within the social media is far bigger than the users victimization it. This counsel that faux accounts are raised within the recent years. On-line social media suppliers face issue in characteristic these faux accounts. The requirement for characteristic these fake accounts is that social media is busy with bogus statistics, advertisements, etc.

This paper is developed as follows. In next section, we discussed about the Methodology of proposed system. Section 3 deals with accuracy comparison and prediction Section 4 discuss the result and analysis and final section deal with conclusion of the paper.

II. METHODOLOGY

The following modules are presented in this paper. They are

1. Information assortment

2. Uploading dataset and information preprocessing
3. Ensemble Learning strategies
4. Accuracy comparison and Prediction

2.1 Information assortment

Internet hand tool is employed to gather information of instagram profiles. But gift version of instagram doesn't permit information to be scrapped.so we have a tendency to created a manual dataset.

2.2 Uploading dataset and information preprocessing

Data preprocessing is that the initial step to boost the model quality.

Steps:

information assortment

Handling missing information

Handling categorical information

We have used python for preprocessing

2.2 Ensemble Learning strategies

Ensemble learning is employed to average the predictions of various modules permanently prediction.

This technique is helpful once accuracy is very important however there's no constraint on time If we tend to run N models, then it'll be slower by N times

It uses many models with completely different algorithms to enhance the accuracy

Popular techniques square measure

Bagging
 Boosting
 Here we've got used
 Logistic regression
 Ada boost
 Random forest
 XG boost
 Gradient boosting algorithms

III. ACCURACY COMPARISON AND PREDICTION

Few signs that tells that the account is pretend square measure
 The variety of followings is much exceeds the amount of followers
 If photos within the gallery being announce on identical day then it's in all probability announce by a larva
 How several posts square measure there
 If the account has no post however the followers following count is a lot of then it should be a pretend account
 If the account is verified then it belongs to prime professionals
 Type of comments they're posting
 If it's inactive for a old however followers account exceeds
 No personal info in bio
 No real realistic posts
 If it's a default profile pic or no pic in the least.
 Addition of additional area characters in bio, user name. Then there square measure high possibilities for or not its pretend profile.

IV. ALGORITHMS

Ensemble learning says that doesn't depend upon one model. Consider all and notice the ultimate sturdy classifier.

- Boosting and textile area unit example of ensemble learning.
- In boosting models area unit in-built series. At every sequent model the weights area unit adjusted supported the educational of previous model.

Bagging technique

Various models area unit in-built parallel on varied samples and so those models vote to grant the ultimate model and thence prediction.

Adaptive Boosting

Adaptive boosting be a machine learning algorithmic program developed and conjunction with quite a few alternative forms of learning algorithms to boost throughput. The output of the conflicting learning algorithms is mutual weighted sum that represent the final output of the boosted classifier. AdaBoost is edition within the intelligence that future weak learners area unit tweak in good turn of these instance misclassified by earlier classifiers. AdaBoost is sensitive to abuzz in order and outliers.

Steps:

- Lets say if we've 'n' records then assign 1/n as sample weight to every record. This is often the sample dataset.

```

Initialization:
1. Given training data from the instance space
 $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$  where  $x_i \in \mathcal{X}$  and  $y_i \in \mathcal{Y} = \{-1, +1\}$ .
2. Initialize the distribution  $D_1(i) = \frac{1}{m}$ .
Algorithm:
for  $t = 1, \dots, T$ : do
    Train a weak learner  $h_t : \mathcal{X} \rightarrow \mathbb{R}$  using distribution  $D_t$ .
    Determine weight  $\alpha_t$  of  $h_t$ .
    Update the distribution over the training set:
        
$$D_{t+1}(i) = \frac{D_t(i)e^{-\alpha_t y_i h_t(x_i)}}{Z_t}$$

    where  $Z_t$  is a normalization factor chosen so that  $D_{t+1}$  will be a distribution.
end for
Final score:
 $f(x) = \sum_{t=0}^T \alpha_t h_t(x)$  and  $H(x) = \text{sign}(f(x))$ 
    
```

SVM algorithm

Machine learning involves predicting and classifying the data tends to use diverse these algorithms in line with the dataset. It will solve linear and non-linear issues and work well for several sensible issues. The thought of SVM is simple: The algorithmic program creates a hyper plane that separates the info into category. In machine learning, the radial basis operate kernel, could be a widespread kernel operate employed in varied kernelized learning algorithms. Especially, it's normally employed in support vector machine classification.

Decision tree

A decision tree is a graphical illustration that creates use of branch methodology to exemplify all within reach outcome of a call and supported sure conditions. The inner node represent a look at on the attribute, each branch of the tree represents the outcome of the look at and the leaf

node represent a detailed category label i.e. the choice created once computing all of the attributes.

Random Forest Tree

Random Forest is the machine learning algorithmic program that uses a textile approach to make a bunch of call trees with random set. A model is trained a lot random sample of the dataset to realize sane prediction performance many times. The output of all the choice trees within the tree, combined to create the ultimate prediction. For instance, within the higher than example - if five friends decide that you simply can like building R however solely a pair of friends decide that you simply won't just like the building then the ultimate prediction is that, you may like building R as majority continually wins.

Xgboost Algorithm

The following algorithm for XGboost.

Algorithm 1: XGboost algorithm

Data: Dataset and hyperparameters

Initialize $f_0(x)$;

for $k = 1, 2, \dots, M$ **do**

Calculate $g_k = \frac{\partial L(y, f)}{\partial f}$;

Calculate $h_k = \frac{\partial^2 L(y, f)}{\partial f^2}$;

Determine the structure by choosing splits with maximized gain

$A = \frac{1}{2} \left[\frac{G_L^2}{H_L} + \frac{G_R^2}{H_R} - \frac{G^2}{H} \right]$;

Determine the leaf weights $w^* = -\frac{G}{H}$;

Determine the base learner $\hat{b}(x) = \sum_{j=1}^T w I_j$;

Add trees $f_k(x) = f_{k-1}(x) + \hat{b}(x)$;

end

Result: $f(x) = \sum_{k=0}^M f_k(x)$

Logistic Regression

- To use supply Regression we've to use some a part of the info to estimate the form of this curve.
- Then we'd like to seek out if that curve is doing categorizing new information.
- Re-using constant information for each coaching and testing could be a dangerous plan as a result of we'd like to grasp however the model can work on information it wasn't trained on.

Use 75% data for training.
 25% data for testing.

V. RESULT AND ANALYSIS

5.1 Dataset Description

The instagram data taken from kaggle website for my research.

profile pic	nums/len	fullname	w nums/len	name==use	description	external	UF	private	#posts	#followers	#follows	fake
1	0.33	1	0.33	1	30	0	1	1	35	488	604	0
1	0	5	0	0	64	0	1	1	3	35	6	0
1	0	2	0	0	82	0	1	1	319	328	668	0
1	0	1	0	0	143	0	1	1	273	14890	7369	0
1	0.5	1	0	0	76	0	1	1	6	225	356	0
1	0	1	0	0	0	0	1	1	6	362	424	0
1	0	1	0	0	132	0	1	1	9	213	254	0
1	0	2	0	0	0	0	1	1	19	552	521	0
1	0	2	0	0	96	0	1	1	17	122	143	0
1	0	1	0	0	78	0	1	1	9	834	358	0
1	0	1	0	0	0	0	1	1	53	229	492	0
1	0.14	1	0	0	78	1	1	1	97	1913	436	0
1	0.14	2	0	0	61	0	1	1	17	200	437	0
1	0.33	2	0	0	45	0	1	1	8	130	622	0
1	0.1	2	0	0	43	0	0	0	60	192	141	0
1	0	2	0	0	56	0	1	1	51	498	337	0
1	0.33	2	0	0	86	0	1	1	25	96	499	0
1	0	1	0	0	97	0	1	1	188	202	605	0
1	0	3	0	0	46	0	1	1	590	175	199	0
1	0	2	0	0	39	0	1	1	251	223	694	0
1	0.5	1	0	0	0	0	1	1	0	189	276	0

Figure 1: Instagram dataset

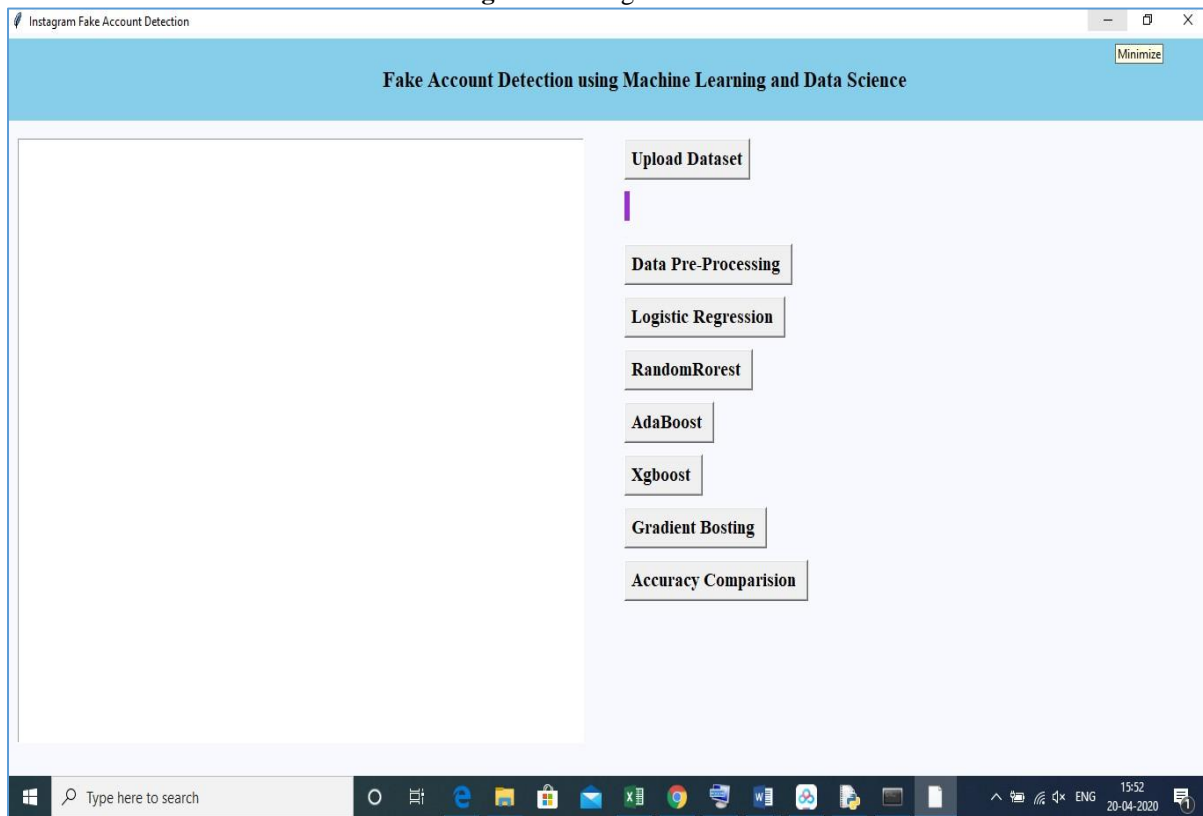


Figure 2: Fake account Detection using Machine learning techniques

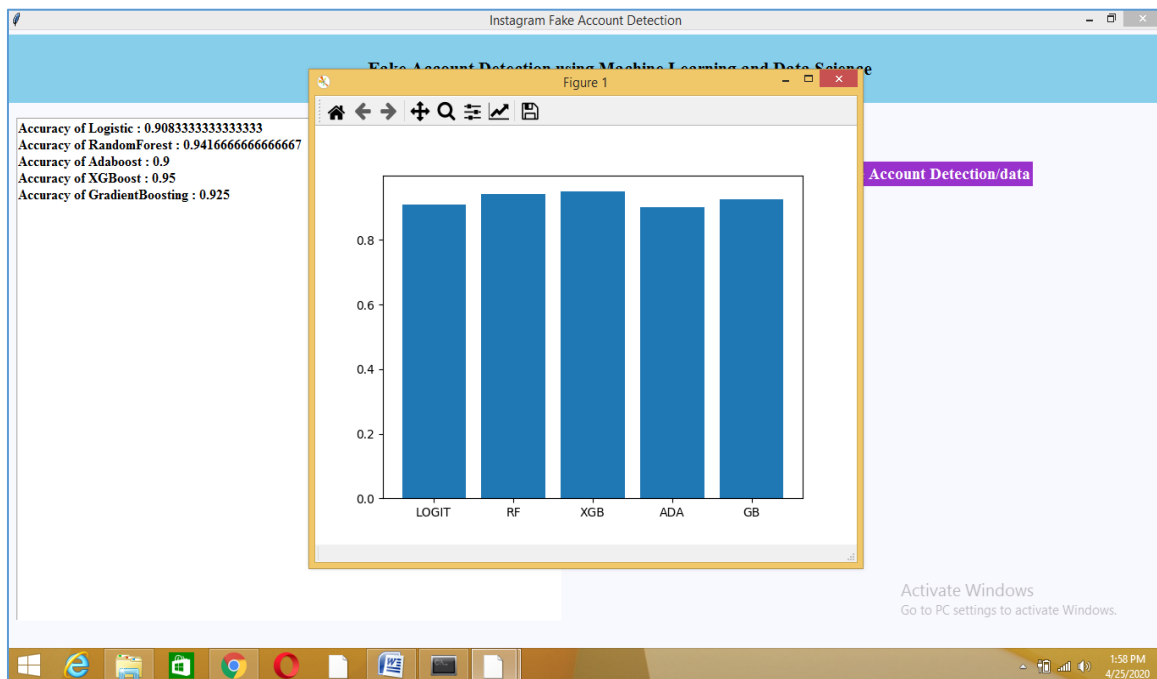


Figure 3: Comparison of machine learning techniques against Fake account detection

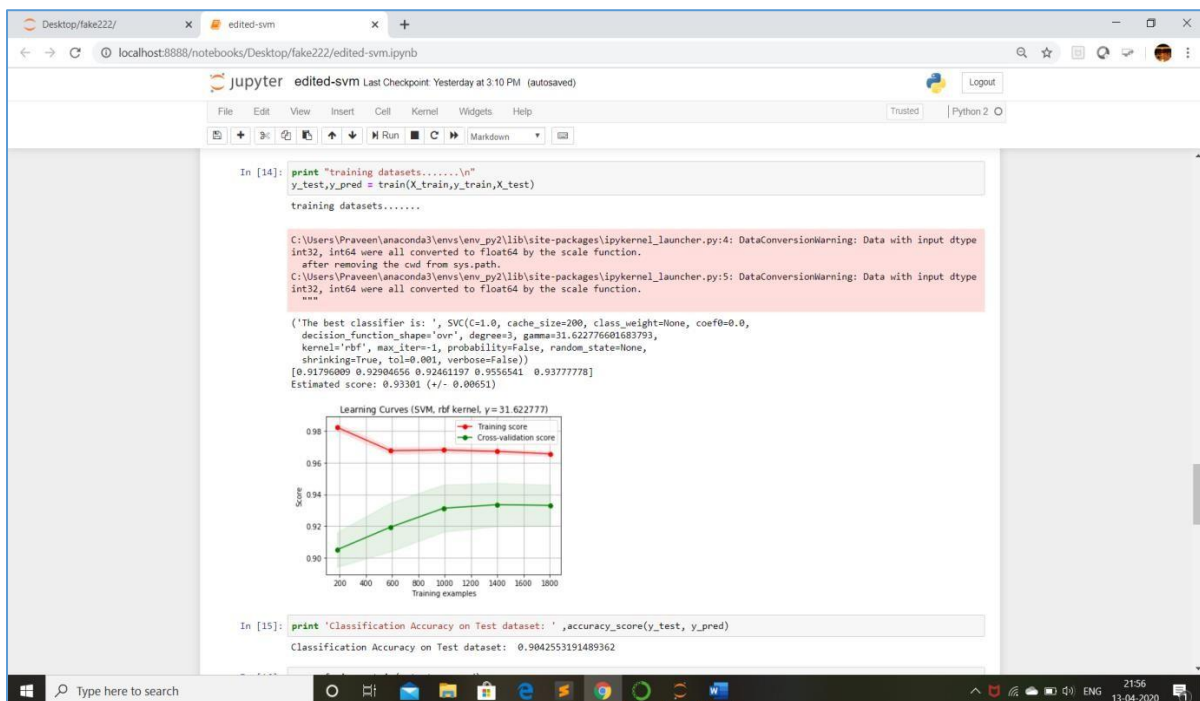


Figure 4: Among all algorithmic performance, XGboost performed well .

In this paper multiple machine learning algorithms are using against detecting fake accounts in Instagram, among these algorithms XGboost performed well and it give 96% accuracy.

VI. CONCLUSION

Algorithms will be trained with all previous users fake and genuine account data and

then whenever we give new test data then trained model will be applied on new test data to identify whether given new account details are from genuine or fake users. The machine learning-based methods were used to perceive false accounts that could give the wrong impression about people. The dataset is pre-processed using a variety of python libraries and a comparison form is obtained to get a

realistic algorithm appropriate for the specified dataset .An effort to notice forged accounts on the social media platforms is strong-minded by a variety of machine learning algorithms. The performances of the classification algorithms Random Forest, Neural Network and support vector machines are used for the detection of fake accounts. Although developing a project that has one hundred pc accuracy that predicts all accounts properly is much not possible. However this project works well even some options are missing .and accuracy is 96 %.

REFERENCES

- [1]. Ilhan aydin,Mehmet sevi, Mehmet umut salur”Detection of pretend Twitter accounts with Machine Learning Algorithms”., 2018.
- [2]. Naman singh, Tushar sharma, Abha Thakral, Tanupriya Choudhury”Detection of pretend profile in on-line social networks victimization Machine Learning”.2018.
- [3]. T.M. Mitchell, Machine Learning. McGraw-Hill, 1997.
- [4]. Yael Ben-Haim, “A Streaming Parallel Decision Tree Algorithm” , Elad Tom-Tov , 2010.
- [5]. Breiman, L., Random Forests "Bagging predictors," Machine Learning, vol. 24, no. 2, pp. 123-140, 1996.

M.Mamatha, et. al. “Fake Profile Identification using Machine Learning Algorithms.” *International Journal of Engineering Research and Applications (IJERA)*, vol.11 (7), 2021, pp 60-65.