A

Project Report

On


# Detecting Fake News With Python And Machine Learning


Submitted in partial fulfilment of the requirement for the 3$^{rd}$ semester

BTECH

By

Yash Sharma 2017586


Under the Guidance of

Dr. Vishen Gupta Sir


Deptt. of Computer Science & Application



DEPARTMENT OF COMPUTER SCIENCE & APPLICATION

GRAPHIC ERA UNIVERSITY, DEHRADUN


CLEMENT TOWN, BELL ROAD

DISTRICT- DEHRADUN-248002

2021 – 2022

## CANDIDATE'S DECLARATION

I hereby certify that the work which is being presented in the dissertation entitled "Detecting Fake News Using Python And Machine Learning" in partial fulfilment of the requirements for the award of the Degree of Bachelor of Technology in Computer Science and Engineering, submitted in the Department of Computer Science and Engineering of the Graphic Era Deemed to be University, Dehradun is an authentic record of my own work carried out during a period from October 2021 to February 2022, under the supervision of Dr. Vishen Gupta Sir, Department of Computer Science and Engineering of the Graphic Era Deemed to be University, Dehradun (Uttarakhand).

The matter presented in this dissertation has not been submitted by me for the award of any other degree of this or any other Institute/University.

YASH SHARMA

2017586

# ACKNOWLEDGEMENT

# TABLE OF CONTENTS

| S.No. | Title | Page No. |
|-------|-------|----------|
| 1.0 | Abstract | 5-5 |
| 1.1 | Introduction | 5 |
| 1.2 | Motivation | 5 |
| 2.0 | Requirements of Project | 5-5 |
| 2.1 | Hardware Requirement | 5 |
| 2.2 | Software Requirement | 5 |
| 3.0 | Methodology Followed | 5-6 |
| 3.1 | Pandas Library | |
| 3.2 | Scikit-Learn (Sklearn) Module | |
| 3.3 | TfidfVectorizer | |
| 3.4 | Numpy | |
| 3.5 | stopwords | |
| 3.6 | PorterStemmer(NLTK Package) | |
| 4.0 | Screenshots | 7-9 |
| 5.0 | Conclusion | 9 |
| 6.0 | References | 9 |

# 1. ABSTRACT

## 1.1 Introduction

 This advanced python project of detecting fake news deals with fake and real news. Using sklearn, we build a TfidfVectorizer on our dataset. Then, Train the Model using Logistic Regression . In the end, the accuracy score and the confusion matrix tell us how well our model fares.

## 1.2 Motivation

Fake News has become one of the major problem in the existing society. Fake News has high potential to change opinions, facts and can be the most dangerous weapon in influencing society

With our world producing an ever-growing huge amount of data exponentially per second by machines, there is a concern that this data can be false (or fake). Fake news (or data) can pose many dangers to our world. Imagine what happens if due to some false information you are given the wrong medicine.

Luckily, this problem can be addressed using machine learning with python. We can develop a machine learning model in python which can detect whether the news is fake or not

# 2. REQUIREMENT OF PROJECT

### 2.1 Hardware Requirement
- A working computer
- Internet Connection

### 2.2 Software Requirement
- Internet Browser such as Chrome, Edge, etc.
- Google Colaboratory
- Some DataSets (Module)

# 3. METHODOLOGY FOLLOWED

3.1 Pandas Library   :- Pandas allows us to analyze big data and make conclusions based on statistical theories.Pandas can clean messy data sets, and make them readable and relevant.Relevant data is very important in data science.

Pandas In this Project :-
Read_csv('Path of file') by using this function we can read the csv file in this we Have to give the path of the csv file which we upload

## 3.2 Scikit-Learn (Sklearn) Library

Scikit-learn is an indispensable part of the Python machine learning toolkit at JPMorgan. It is very widely used in predictive analytics, and very many other machine learning tasks. Its straightforward API, its breadth of algorithms, and the quality of its documentation combine to make scikit-learn simultaneously very approachable and very powerful.

## 3.3 TfidfVectorizer

**TF (Term Frequency):** The number of times a word appears in a document is its Term Frequency. A higher value means a term appears more often than others, and so, the document is a good match when the term is part of the search terms.

**IDF (Inverse Document Frequency):** Words that occur many times a document, but also occur many times in many others, may be irrelevant. IDF is a measure of how significant a term is in the entire corpus.

The TfidfVectorizer converts a collection of raw documents into a matrix of TF-IDF features.

## 3.4 NUMPY

NumPy is a general-purpose array-processing package. It provides a high-performance multidimensional array object, and tools for working with these arrays. It is the fundamental package for scientific computing with Python. ... A powerful N-dimensional array object

## 3.5 stopwords

Stopwords are the English words which does not add much meaning to a sentence. They can safely be ignored without sacrificing the meaning of the sentence. For example, the words like the, he, have etc. Such words are already captured this in corpus named corpus.

## 3.6 PorterStemmer

Porters Stemmer It is a type of stemmer which is mainly known for Data Mining and Information Retrieval. As its applications are limited to the English language only. It is based on the idea that the suffixes in the English language are made up of a combination of smaller and simpler suffixes, it is also majorly known for its simplicity and speed. The advantage is, it produces the best output from other stemmers and has less error rate.

# 4. Screenshots

## Fig. 1



Fig.1 shows importing required libraries, and use stopwords lib. And then print them
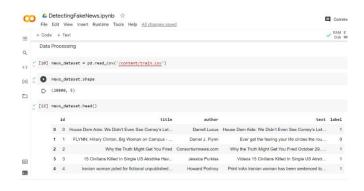
## Fig.2



Fig.2 Data Preprocessing. And print first 5 rows

# Fig. 3



Fig 3. Shows Check Nulls Value and fill them with empty string and then print data

# Fig 4.



Fig 4. Shows Performining of stemming function :- Stemming is the process of reducing a word to its Root word

Fig. 5



Fig.5 Printing Label and content (author , text). Column after seprating and then shows converting the texts into feature vectors so that we can find cosine similarity of the values. Shown below
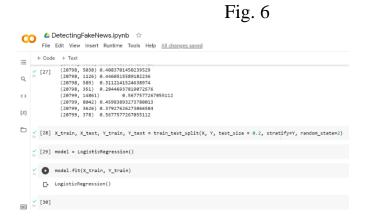
Fig. 6



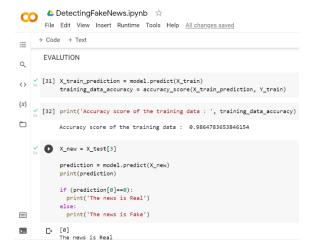Fig.6 Split the dataset into training and testing sets.

Fig. 7



Fig.7  Making the predictive System. We got an accuracy of 92.64% with this model. And get result of given news is true  or false

# 5. Conclusion

In this Project we predict the news is Fake or Not By using by Python Module Like Numpy,pandas

Etc. The completion of the project went quiet well, I learned much new things while I was building up it, and I get up to know various platforms which help us to learn all this stuff. I was able to learn the practical use of ML and Python .

Overall working on this project was great fun as I came up with great piece of knowledge and understanding of the topic

# 6. References

https://www.delftstack.com https://www.etutorialspoint.com

https://www.geeksforgeeks.org https://www.javatpoint.com

https://pandas.pydata.org/docs/reference/api/pandas.read_csv.html

https://drive.google.com/file/d/1wxbX1pmdjAinILQS21YIbDJDz6I8

8qks/view?usp=sharing

https://drive.google.com/file/d/1Ef8Pmxp6KHBx2bU6eOxnhnth53T9

mr9f/view?usp=sharing

www.youtube.com