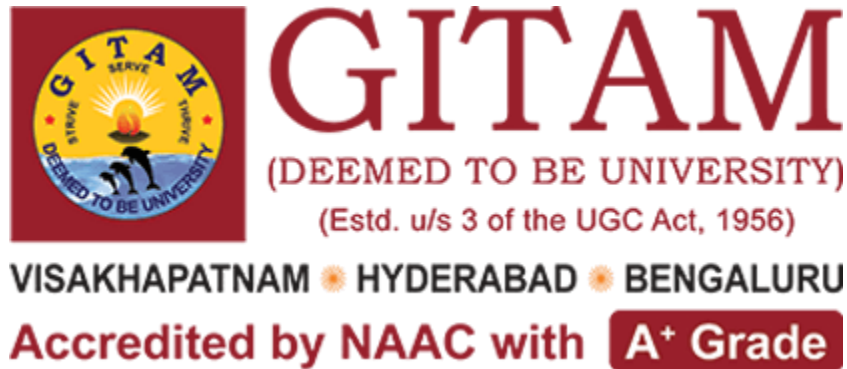


**MACHINE LEARNING
(MOVIE REVIEWS CLASSIFICATION)**

*Summer Internship Report Submitted in partial fulfillment of the
requirement for undergraduate degree of*

**Bachelor of Technology In
COMPUTER SCIENCE AND ENGINEERING
By CHAKILAM VARUN
221710309012**

Under the Guidance of



**Department of COMPUTER SCIENCE AND ENGINEERING
GITAM School of Technology
GITAM (Deemed to be University) Hyderabad-502329
July 2020**

DECLARATION

I submit this industrial training work entitled “MOVIE REVIEWS CLASSIFICATION” to GITAM (Deemed to Be University), Hyderabad in partial fulfillment of the requirements for the award of the degree of “Bachelor of Technology” in “Computer Science and Engineering”. I declare that it was carried out independently by me under the guidance of, GITAM (Deemed to Be University), Hyderabad, India.

The results embodied in this report have not been submitted to any other University or Institute for the award of any degree or diploma.

PLACE: Hyderabad

CHAKILAM VARUN

DATE:

221710309012



GITAM (DEEMED TO BE UNIVERSITY)

Hyderabad-502329, India Dated:

CERTIFICATE

This is to certify that the Industrial Training Report entitled “MOVIE REVIEWS CLASSIFICATION” is being submitted by CHAKILAM VARUN (221710309012) in partial fulfillment of the requirement for the award of Bachelor of Technology in Computer Science & Engineering at GITAM (Deemed to Be University), Hyderabad during the academic year 2019-2020.

It is faithful record work carried out by him at the Computer Science & Engineering Department, GITAM University Hyderabad Campus under my guidance and supervision.

ACKNOWLEDGEMENT

Apart from my effort, the success of this internship largely depends on the encouragement and guidance of many others. I take this opportunity to express my gratitude to the people who have helped me in the successful competition of this internship.

I would like to thank respected **Dr. N. Siva Prasad**, Pro Vice Chancellor, GITAM Hyderabad and **Prof.N.SeetaRamaiah**, Principal, GITAM Hyderabad

I would like to thank respected **Mr.S.Phani Kumar**, Head of the Department of Computer Science & Engineering for giving me such a wonderful opportunity to expand my knowledge for my own branch and giving me guidelines to present an internship report. It helped me a lot to realize of what we study for.

I would like to thank the respected faculties who helped me to make this internship a successful accomplishment.

I would also like to thank my friends who helped me to make my work more organized and well-stacked till the end.

CHAKILAM VARUN

221710309012

ABSTRACT

Machine learning algorithms are used to predict the values from the data set by splitting the data set in to train and test and building Machine learning algorithms. The primary task is to build the models of higher accuracy .My perception of understanding the given data set has been in the view of undertaking the reviews of various viewers who watch a movie and give their reviews.

As human's opinion help in enhance product's efficiency, and since the success or the failure of a movie depends on its reviews, there is an increase in the demand and need to build a good sentiment analysis model that classifies movie reviews. In this research, tokenization is performed to transfer the input string into a word vector. The model is evaluated on a real world dataset. We use three algorithms to evaluate our data set.

The results show that Logistic Regression outperforms the other classifiers. Later, hyper parameter tuning is performed on the Logistic Regression method to improve the accuracy

Table of Contents

CHAPTER 1	1
MACHINE LEARNING.....	1
1.1. INTRODUCTION:	1
1.2. IMPORTANCE OF MACHINELEARNING:	1
1.3. USES OF MACHINELEARNING:	2
1.4. TYPES OF LEARNINGALGORITHMS:.....	3
1.4.1. Supervised Learning:.....	3
1.4.2. UnsupervisedLearning:	4
1.4.3. Semi Supervised Learning:	5
1.5. RELATION BETWEEN DATA MINING, MACHINEAND DEEP LEARNING:	5
CHAPTER 2	6
PYTHON.....	6
2.1. INTRODUCTION TOPYTHON:.....	6
2.2. HISTORY OF PYTHON:	6
2.3. FEATURES OF PYTHON:	6
2.4. HOW TO SETUP PYTHON:	7
2.4.1. INSTALLATION (using pythonIDLE):	7
2.4.2. Installation (using Anaconda):	8
2.5. PYTHON VARIABLE TYPES:.....	9
2.5.1. PythonNumbers:.....	10
2.5.2. PythonStrings:	10
2.5.3. PythonLists:.....	10
2.5.4. PythonTuples:.....	10
2.5.5. PythonDictionary:	11
2.6. PYTHON FUNCTION:.....	12
2.6.1. Defining aFunction:.....	12
2.6.2. Calling aFunction:	12
2.7. PYTHON USING OOPS CONCEPTS:.....	12
2.7.1. Class:	12
2.7.2. __init__ method inClass:.....	13
CHAPTER 3	14
CASE STUDY.....	14
3.1. PROBLEMSTATEMENT:	15
3.2. DATA SET:	15
3.3. OBJECTIVE OF THE CASESTUDY:	15
CHAPTER 4	16
MODELBUILDING.....	16
4.1. PREPROCESSING OF THE DATA:	16
4.1.1. GETTING THEDATASET:	16
4.1.2. IMPORTING THELIBRARIES:	16
4.1.3. IMPORTING THEDATA-SET:	17
4.1.4. HANDLING MISSINGVALUES:	17

4.1.5.	STATISTICAL ANALYSIS:	18
4.1.6.	CATEGORICAL DATA:	19
4.2.	TRAINING THE MODEL:	20
4.2.1.	COUNT VECTORIZATION :	22
4.2.2.	TFIDF VECTORIZER:	23
4.3.	EVALUATING THE CASE STUDY:	25
4.3.1.	NAIVE BAYES ALGORITHM:	25
4.3.2.	LOGISTIC REGRESSION:	30
4.3.3.	RANDOM FOREST CLASSIFICATION:	35
4.4.	COMPARISON OF ACCURACY SCORES OF TRAIN AND TEST DATA:	39
4.5.	HYPER PARAMETER TUNING ON LOGISTIC REGRESSION	41
4.6.	TESTING WITH UNKNOWN INPUT:	44
4.7.	ROC AUC CURVE:	45
CONCLUSION:		48
REFERENCES:		49

List of Figures

Figure 1.2. 1: The Process Flow	2
Figure 1.4.2.1: Unsupervised Learning	4
Figure 1.4.3.1.Semi Supervised Learning	5
Figure 2.4.1.1. Python download	7
Figure 2.4.2.1: Anaconda download	8
Figure 2.4.2.2.Jupyter notebook.....	9
Figure 2.7.1.1: Defining a Class.....	13
Figure 4.1.2.1: Importing Libraries	16
Figure 4.1.3.1:Reading the dataset	17
Figure 4.1.4.1:Checking for missing values.....	18
Figure 4.1.5 1: Shape of the data.....	18
Figure 4.1.5.2:Features of sentiment column.....	18
Figure 4.1.5.3: Visualization of the data set.....	19
Figure 4.1.6.1: Categorical data	20
Figure 4.2.1 : importing train_test_split.....	21
Figure 4.2.2: Train and Test data	21
Figure 4.2.1.1:Importing CountVectorizer.....	22
Figure 4.2.1.2:Generating the word counts	22
Figure 4.2.2.1:TFIDF Formula.....	23
Figure 4.2.2.2:Importing tfidf Vectorizer.....	23
Figure 4.2.2.3:Vectorisation on train and test data	24
Figure 4.2.2.4:Position of words	24
Figure 4.3.1.1: Naïve Baye's Formula	25
Figure 4.3.1.2: Importing Bernoulli Naive Bayes Method	26
Figure 4.3.1.3: Fitting the model using BernoulliNb method	26
Figure 4.3.1.4: Prediction on train data (Naive Bayes)	27
Figure 4.3.1.5: Confusion matrix for train data(Naive Bayes).....	27

Figure 4.3.1.6: Classification report for train data(Naive Bayes)	27
Figure 4.3.1.7: Heat Map of train data(Naive Bayes)	28
Figure 4.3.1.8: Prediction on test data(Naive Bayes).....	28
Figure 4.3.1.9:Confusion matrix for testing data(Naive Bayes)	28
Figure 4.3.1.10: Classification report for test data(Naive Bayes)	29
Figure 4.3.1.11: Heat Map of test data(Naive Bayes)	29
Figure 4.3.2.1:Importing Logistic Regression Method and fitting the model.....	31
Figure 4.3.2 2Figure 34:Prediction on train data(Logistic Regression)	31
Figure 4.3.2.3:Confusion matrix for train data(Logistic Regression)	32
Figure 4.3.2.4:Classification report for train data(Logistic Regression).....	32
Figure 4.3.2.5:Heat Map of train data(Logistic Regression).....	33
Figure 4.3.2.6:Prediction on test data(Logistic Regression)	33
Figure 4.3.2.7:Confusion matrix for test data(Logistic Regression).....	33
Figure 4.3.2 8:Classification report for test data(Logistic Regression)	34
Figure 4.3.2.9:Heat Map of test data(Logistic Regression)	34
Figure 4.3.3.1:Importing Random Forest Classifier Method and fitting the model.....	36
Figure 4.3.3.2:Prediction on train data and classification report(Random Forest Classifier)	36
Figure 4.3.3.3:Accuracy for train data(Random Forest Classification)	37
Figure 4.3.3.4:Heat Map of train data(Random Forest Classifier).....	37
Figure 4.3.3.5:Prediction on test data and classification report(Random Forest Classifier).....	38
Figure 4.3.3.6:Accuracy for test data(Random Forest Classification)	38
Figure 4.3.3.7:Heat Map of test data(Random Forest Classifier)	38
Figure 4.4.1:defining labels and printing accuracies.....	39
Figure 4.4.2:Visulaising the train accuracies	39
Figure 4.4.3:Visulaising the test accuracies	40
Figure 4.5.1:Passing a list of values as dictionary	41
Figure 4.5.2:Importing GridSearchCV and passing the dictionary of parameters.....	41
Figure 4.5.3:Printing the best parameters.....	42
Figure 4.5.4:Building the model with best parameters	42
Figure 4.5.5:Prediction on test data.....	42
Figure 4.5.6:Confusion Matrix of Test data(Hyper Parameter Tuning)	42
Figure 4.5.7:Classification report and accuracy of test data (Hyper Parameter Tuning)	43
Figure 4.5.8:Heat Map of test data(Hyper Parameter Tuning)	43
Figure 4.6.1: Reading unknown Input.....	44
Figure 4.6.2:Predicting the unknown input	44
Figure 4.7.1:Importing packages for Roc-Auc Curve.....	45
Figure 4.7.2:Calculating the predicted probabilities and importing roc auc curve	46

Figure 4.7.3:Calculating the auc scores 46

Figure 4.7.4:Plotting ROC Curve..... 47

Figure 4.7.5: ROC Curve 47

CHAPTER 1

MACHINE LEARNING

1.1. INTRODUCTION:

Machine Learning (ML) is the scientific study of algorithms and statistical models that computer systems use in order to perform a specific task effectively without using explicit instructions, relying on patterns and inference instead. It is seen as a subset of Artificial Intelligence (AI).

1.2. IMPORTANCE OF MACHINELEARNING:

Consider some of the instances where machine learning is applied: the self-driving Google car, cyber fraud detection, online recommendation engines—like friend suggestions on Face book, Netflix showcasing the movies and shows you might like, and “more items to consider” and “get yourself a little something” on Amazon—are all examples of applied machine learning. All these examples echo the vital role machine learning has begun to take in today’s data-rich world.

Machines can aid in filtering useful pieces of information that help in major advancements, and we are already seeing how this technology is being implemented in a wide variety of industries.

With the constant evolution of the field, there has been a subsequent rise in the uses, demands, and importance of machine learning. Big data has become quite a buzzword in the last few years; that’s in part due to increased sophistication of machine learning, which helps analyze those big chunks of big data. Machine learning has also changed the way data extraction, and interpretation is done by involving automatic sets of generic methods that have replaced traditional statistical techniques.

The process flow depicted here represents how machine learning works

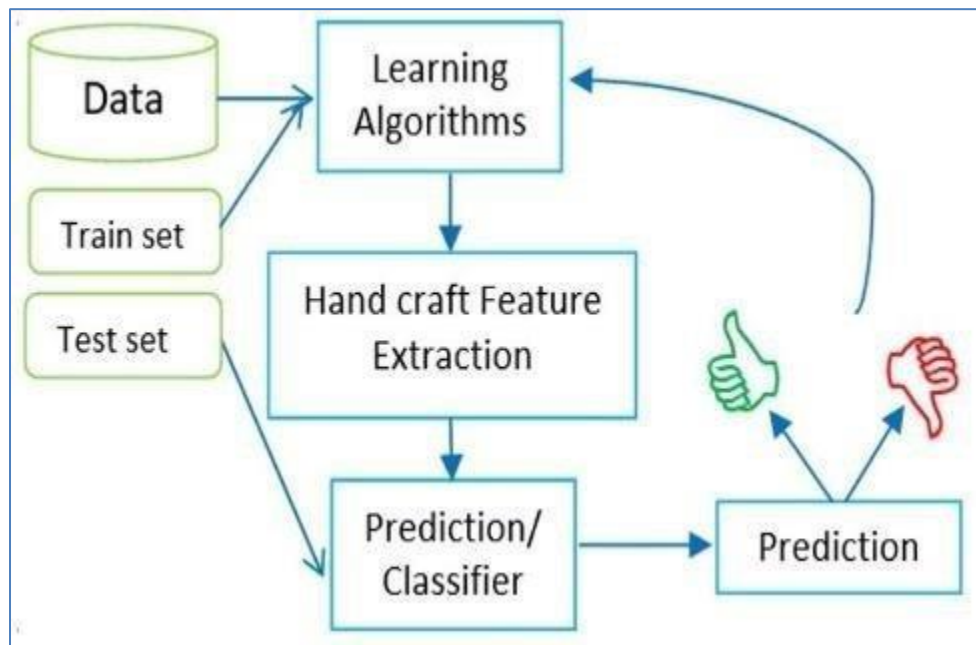


Figure 1.2.1: The Process Flow

1.3. USES OF MACHINELEARNING:

Earlier in this article, we mentioned some applications of machine learning. To understand the concept of machine learning better, let's consider some more examples: web search results, real-time ads on web pages and mobile devices, email spam filtering, network intrusion detection, and pattern and image recognition. All these are by-products of applying machine learning to analyze huge volumes of data

Traditionally, data analysis was always being characterized by trial and error, an approach that becomes impossible when data sets are large and heterogeneous. Machine Learning comes as the solution to all this chaos by proposing clever alternatives to analyzing huge volumes of data.

By developing fast and efficient algorithms and data-driven models for real-time processing of data, machine learning can produce accurate results and analysis.

1.4. TYPES OF LEARNINGALGORITHMS:

The types of machine learning algorithms differ in their approach, the type of data they input and output, and the type of task or problem that they are intended to solve.

1.4.1. Supervised Learning:

When an algorithm learns from example data and associated target responses that can consist of numeric values or string labels, such as classes or tags, in order to later predict the correct response when posed with new examples comes under the category of supervised learning.

Supervised machine learning algorithms uncover insights, patterns, and relationships from a labeled training dataset – that is, a dataset that already contains a known value for the target variable for each record. Because you provide the machine learning algorithm with the correct answers for a problem during training, it is able to “learn” how the rest of the features relate to the target, enabling you to uncover insights and make predictions about future outcomes based on historical data.

Examples of Supervised Machine Learning Techniques are Regression, in which the algorithm returns a numerical target for each example, such as how much revenue will be generated from a new marketing campaign.

Classification, in which the algorithm attempts to label each example by choosing between two or more different classes. Choosing between two classes is called binary classification, such as determining whether or not someone will default on a loan. Choosing between more than two classes is referred to as multiclass classification.

1.4.2. Unsupervised Learning:

When an algorithm learns from plain examples without any associated response, leaving to the algorithm to determine the data patterns on its own. This type of algorithm tends to restructure the data into something else, such as new features that may represent a class or a new series of uncorrelated values. They are quite useful in providing humans with insights into the meaning of data and new useful inputs to supervised machine learning algorithms.

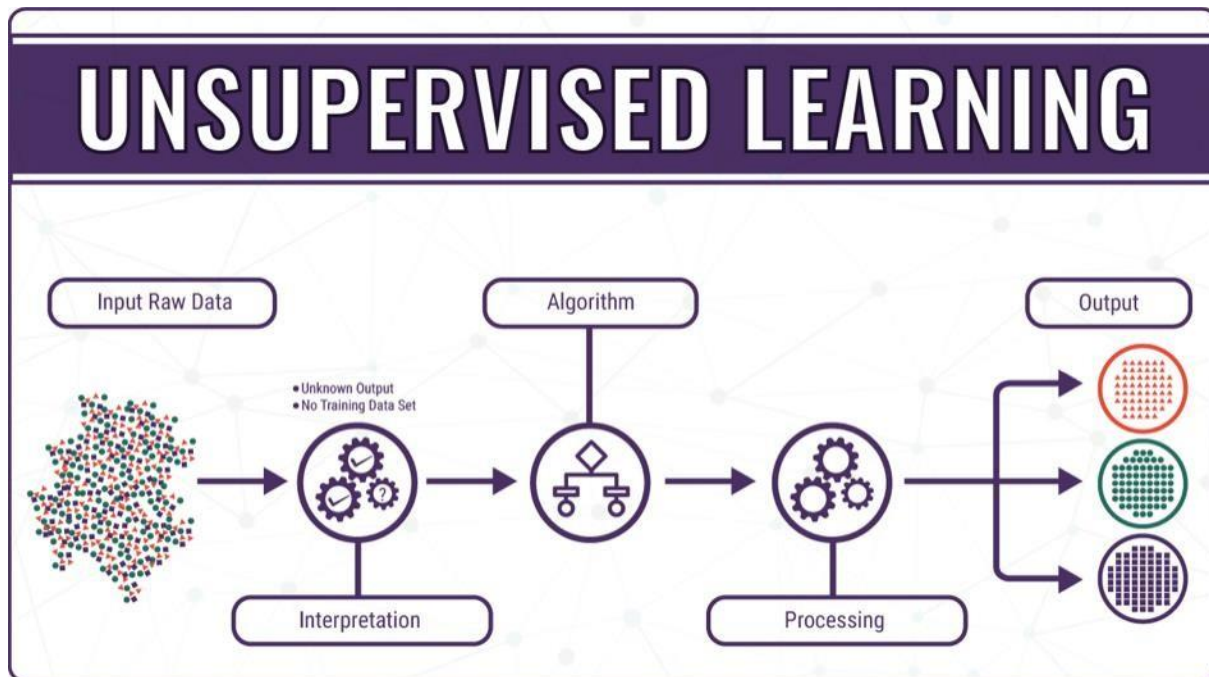


Figure 1.4.2.1: Unsupervised Learning

Popular techniques where unsupervised learning is used also include self-organizing maps, nearest neighbor mapping, singular value decomposition, and k-means clustering. Basically, online recommendations, identification of data outliers, and segment text topics are all examples of unsupervised learning.

1.4.3. Semi Supervised Learning:

As the name suggests, semi-supervised learning is a bit of both supervised and unsupervised learning and uses both labeled and unlabeled data for training. In a typical scenario, the algorithm would use a small amount of labeled data with a large amount of unlabeled data.

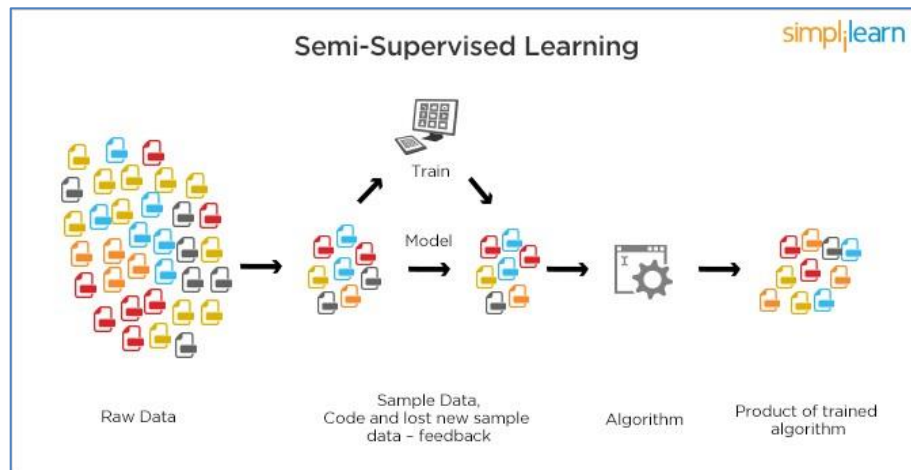


Figure 1.4.3.1.Semi Supervised Learning

1.5. RELATION BETWEEN DATA MINING, MACHINE AND DEEP LEARNING:

Machine learning and data mining use the same algorithms and techniques as data mining, except the kinds of predictions vary. While data mining discovered previously unknown patterns and knowledge, machine learning reproduces known patterns and knowledge—and further automatically applies that information to data, decision-making.

Deep learning, on the other hand, uses advanced computing power and special types of neural networks and applies them to large amounts of data to learn, understand, and identify complicated patterns. Automatic language translation and medical diagnoses are examples of deep learning.

CHAPTER 2

PYTHON

Basic programming language used for machine learning is: PYTHON

2.1. INTRODUCTION TO PYTHON:

- Python is a high-level, interpreted, interactive and object-oriented scripting language.
- Python is a general purpose programming language that is often applied in scripting roles
- Python is interpreted: Python is processed at runtime by the interpreter. You do not need to compile your program before executing it. This is like PERL and PHP.
- Python is Interactive: You can sit at a Python prompt and interact with the interpreter directly to write your programs.
- Python is Object-Oriented: Python supports the Object-Oriented style or technique of programming that encapsulates code within objects.

2.2. HISTORY OF PYTHON:

- Python was developed by GUIDO VAN ROSSUM in early 1990's
- Its latest version is 3.7, it is generally called as python3

2.3. FEATURES OF PYTHON:

- Easy-to-learn: Python has few keywords, simple structure, and a clearly defined syntax. This allows the student to pick up the language quickly.
- Easy-to-read: Python code is more clearly defined and visible to the eyes
- Easy-to-maintain: Python's source code is fairly easy-to-maintaining.
- A broad standard library: Python's bulk of the library is very portable and cross-platform compatible on UNIX, Windows, and Macintosh
- Portable: Python can run on a wide variety of hardware platforms and has the same interface on all platforms.
- Extendable: You can add low-level modules to the Python interpreter. These modules enable programmers to add to or customize their tools to be more efficient.
- Databases: Python provides interfaces to all major commercial databases.

- GUI Programming: Python supports GUI applications that can be created and ported to many system calls, libraries and windows systems, such as Windows MFC, Macintosh, and the X Window system of UNIX.

2.4. HOW TO SETUP PYTHON:

- Python is available on a wide variety of platforms including Linux and Mac OSX. Let's understand how to set up our Python environment.
- The most up-to-date and current source code, binaries, documentation, news, etc., is available on the official website of Python.

2.4.1. INSTALLATION (using python IDLE):

- Installing python is generally easy, and nowadays many Linux and Mac OS distributions include a recent python.
- Download python from www.python.org
- When the download is completed, double click the file and follow the instructions to install it.
- When python is installed, a program called IDLE is also installed along with it. It provides a graphical user interface to work with python.



Figure 2.4.1.1. Python download

2.4.2. Installation (using Anaconda):

- Python programs are also executed using Anaconda.
- Anaconda is a free open source distribution of python for large scale data processing, predictive analytics and scientific computing.
- Conda is a package manager that quickly installs and manages packages.

In WINDOWS:

- Step 1: Open Anaconda.com/downloads in web browser
- Step 2: Download python 3.4 version for (32-bits graphic installer/64 -bit graphic installer)
- Step 3: select installation type(all users)
- Step 4: Select path (i.e. add anaconda to path & register anaconda as default python3.4) next click install and next click finish.
- Step 5: Open Jupyter notebook (it opens in default browser)

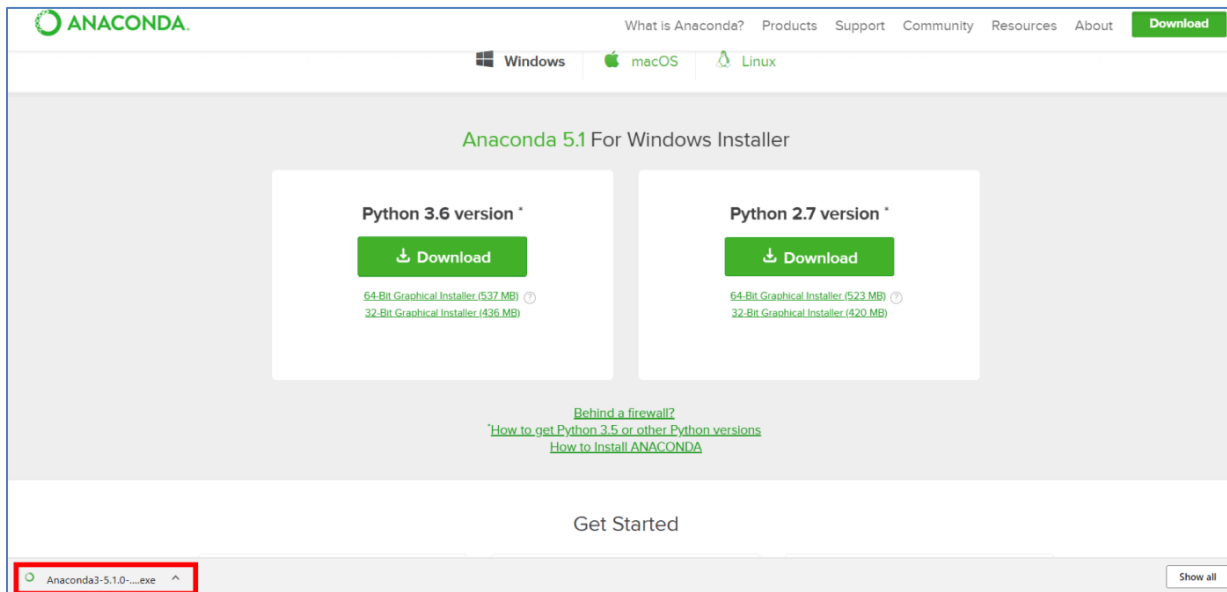


Figure 2.4.2.1: Anaconda download

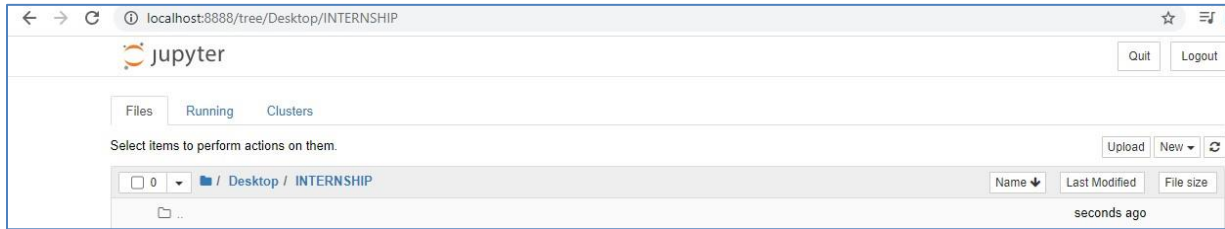


Figure 2.4.2.2.Jupyter notebook

2.5. PYTHON VARIABLE TYPES:

- Variables are nothing but reserved memory locations to store values. This means that when you create a variable you reserve some space in memory.
- Variables are nothing but reserved memory locations to store values.
- Based on the data type of a variable, the interpreter allocates memory and decides what can be stored in the reserved memory.
- Python variables do not need explicit declaration to reserve memory space. The declaration happens automatically when you assign a value to variable.
- Python has various standard data types that are used to define the operations possible on them and the storage method for each of them.
- Python has five standard data types—
 - Numbers
 - Strings
 - Lists
 - Tuples
 - Dictionary

2.5.1. Python Numbers:

- Number data types store numeric values. Number objects are created when you assign a value to them.
- Python supports four different numerical types – int (signed integers) long (long integers, they can also be represented in octal and hexadecimal) float (floating point real values) complex (complex numbers).

2.5.2. Python Strings:

- Strings in Python are identified as a contiguous set of characters represented in the quotation marks.
- Python allows for either pairs of single or double quotes.
- Subsets of strings can be taken using the slice operator ([] and [:]) with indexes starting at 0 in the beginning of the string and working their way from -1 at the end.
- The plus (+) sign is the string concatenation operator and the asterisk (*) is the repetition operator.

2.5.3. Python Lists:

- Lists are the most versatile of Python's compound data types
- A list contains items separated by commas and enclosed within square brackets ([]).
- To some extent, lists are similar to arrays in C. One difference between them is that all the items belonging to a list can be of different data types.
- The values stored in a list can be accessed using the slice operator ([] and [:]) with indexes starting at 0 in the beginning of the list and working their way to end-1.
- The plus (+) sign is the list concatenation operator, and the asterisk (*) is the repetition operator.

2.5.4. Python Tuples:

- A tuple is another sequence data type that is similar to the list.
- A tuple consists of a number of values separated by commas. Unlike lists, however, tuples are enclosed within parentheses.

- The main differences between lists and tuples are: Lists are enclosed in brackets ([]) and their elements and size can be changed, while tuples are enclosed in parentheses (()) and cannot be updated.
- Tuples can be thought of as read-only lists.
- For example – Tuples are fixed size in nature whereas lists are dynamic. In other words, a tuple is immutable whereas a list is mutable. You can't add elements to a tuple. Tuples have no append or extend method. You can't remove elements from a tuple. Tuples have no remove or pop method.

2.5.5. Python Dictionary:

- Python's dictionaries are a kind of hash table type. They work like associative arrays
- Or hashes found in Perl and consist of key-value pairs. A dictionary key can be almost any Python type, but are usually numbers or strings. Values, on the other hand, can be any arbitrary Python object.
- Dictionaries are enclosed by curly braces ({ }) and values can be assigned and accessed using square braces ([]).
- You can use numbers to "index" into a list, meaning you can use numbers to find out what's in lists. You should know this about lists by now, but make sure you understand that you can only use numbers to get items out of a list.
- What a dict does is let you use anything, not just numbers. Yes, a dict associates one thing to another, no matter what it is.

2.6. PYTHON FUNCTION:

2.6.1. Defining a Function:

You can define functions to provide the required functionality. Here are simple rules to define a function in Python. Function blocks begin with the keyword `def` followed by the function name and parentheses (i.e. `()`).

Any input parameters or arguments should be placed within these parentheses. You can also define parameters inside these parentheses

The code block within every function starts with a colon (`:`) and is indented. The statement `return [expression]` exits a function, optionally passing back an expression to the caller. A return statement with no arguments is the same as `return None`.

2.6.2. Calling a Function:

Defining a function only gives it a name, specifies the parameters that are to be included in the function and structures the blocks of code. Once the basic structure of a function is finalized, you can execute it by calling it from another function or directly from the Python prompt.

2.7. PYTHON USING OOPs CONCEPTS:

2.7.1. Class:

- **Class:** A user-defined prototype for an object that defines a set of attributes that characterize any object of the class. The attributes are data members (class variables and instance variables) and methods, accessed via dot notation.
- **Class variable:** A variable that is shared by all instances of a class. Class variables are defined within a class but outside any of the class's methods. Class variables are not used as frequently as instance variables are.
- **Data member:** A class variable or instance variable that holds data associated with a class and its objects.

- Instance variable: A variable that is defined inside a method and belongs only to the current instance of a class.
- Defining a Class:
 - We define a class in a very similar way how we define a function.
 - Just like a function, we use parentheses and a colon after the class name (i.e. () :) when we define a class. Similarly, the body of our class is indented like a functions body is.



The figure shows two code snippets side-by-side, illustrating the syntax for defining a function and a class. The function definition on the left uses 'def' and the class definition on the right uses 'class'. Both show a function/class name followed by parentheses and a colon, then an indented block of code.

```
def my_function():  
    # the details of the  
    # function go here
```

```
class MyClass():  
    # the details of the  
    # class go here
```

Figure 2.7.1.1: Defining a Class

2.7.2. `__init__` method in Class:

- The init method — also called a constructor — is a special method that runs when an instance is created so we can perform any tasks to set up the instance.
- The init method has a special name that starts and ends with two underscores: `init ()`.

CHAPTER 3

CASE STUDY

PROJECT NAME-MOVIE REVIEWS CLASSIFICATION

Packages used:

1. Pandas (1.0.1)
2. Numpy (1.18.1)
3. Matplotlib.pyplot
4. Seaborn (0.10.0)
5. Sklearn.model_selection.train_test_split
6. sklearn.feature_extraction.text.CountVectorizer
7. sklearn.feature_extraction.text.TfidfVectorizer
8. sklearn.naive_bayes.BernoulliNB
9. sklearn.metrics.classification_report,accuracy_score,confusion_matrix
10. sklearn.linear_model.LogisticRegression
11. sklearn.ensemble.RandomForestClassifier
12. sklearn.metrics.roc_curve,roc_auc_score
13. sklearn.model.selection.GridSearchCV

Algorithms used:

1. NAIVE BAYESALGORITHM
2. LOGISTIC REGRESSIONALGORITHM
3. RANDOM FOREST CLASSIFIERALGORITHM

3.1. PROBLEM STATEMENT:

To classify a given review into a positive or a negative review using sentiment analysis in Machine Learning

3.2. DATA SET:

The given data set consists of the following parameters:

A-Review

B-Sentiment

3.3. OBJECTIVE OF THE CASE STUDY:

Movie reviews are an important way to gauge the performance of a movie. While providing a numerical/stars rating to a movie tells us about the success or failure of a movie quantitatively, a collection of movie reviews is what gives us a deeper qualitative insight on different aspects of the movie. Movie Review classification helps us in classifying the review into a positive or a negative review. Sentiment Analysis is a major subject in machine learning which aims to extract subjective information from the textual reviews. Using sentiment analysis, we can find the state of mind of the reviewer while providing the review and understand if the person was “happy”, “sad”, “and angry” and soon, i.e. if they liked the movie or they hated it. We aim to utilize the relationships of the words in the review to predict the overall polarity of the review.

CHAPTER 4

MODELBUILDING

4.1. PREPROCESSING OF THE DATA:

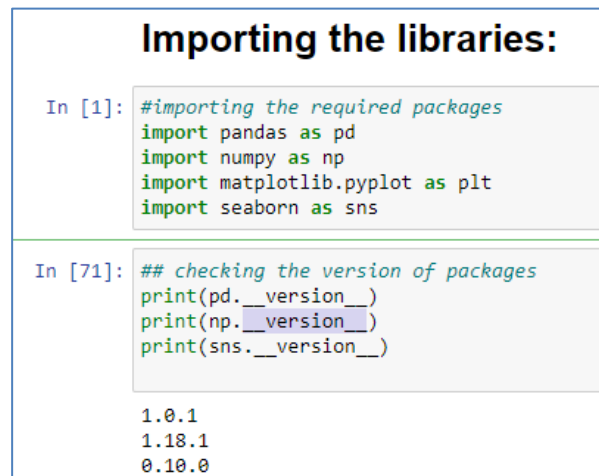
Preprocessing of the data actually involves the following steps:

4.1.1. GETTING THEDATASET:

We can get the data set from the database or we can get the data from client.

4.1.2. IMPORTING THELIBRARIES:

We have to import the libraries as per the requirement of the algorithm.



```
Importing the libraries:

In [1]: #importing the required packages
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

In [71]: ## checking the version of packages
print(pd.__version__)
print(np.__version__)
print(sns.__version__)

1.0.1
1.18.1
0.10.0
```

Figure 4.1.2.1: Importing Libraries

The versions of the packages used are:

- Pandas:1.0.1
- Numpy:1.18.1
- Seaborn:0.10.0

4.1.3. IMPORTING THEDATA-SET:

Pandas in python provide an interesting method read_csv(). The read_csv function reads the entire dataset from a comma separated values file and we can assign it to a Data Frame to which all the operations can be performed. It helps us to access each and every row as well as columns and each and every value can be access using the dataframe. Any missing value or NaNvalue have to be cleaned.



Figure 4.1.3.1:Reading the dataset

4.1.4. HANDLING MISSINGVALUES:

Missing values can be handled in many ways using some inbuilt methods:

- (a) Dropna()
- (b) fillna()
- (c) interpolate ()
- (d) Mean imputation and median imputation

Data Cleaning/ Data Pre-processing:

```
In [137]: #checking for any missing values  
data.isnull().sum()
```

```
Out[137]: review      0  
          sentiment    0  
          dtype: int64
```

Figure 4.1.4.1: Checking for missing values

Since there were no missing values, we can proceed with the further steps.

4.1.5. STATISTICAL ANALYSIS:

- The number of rows and columns of the dataset can be found using `data.shape()`

```
In [70]: #checking the shape of the data  
data.shape
```

```
Out[70]: (50000, 2)
```

Figure 4.1.5 1: Shape of the data

- Our data set contains 50000 rows and 2 columns
- The 2 columns in our dataset are review and sentiment columns
- We need to describe the sentiment column and find the counts of the features. This can be done as follows:

```
In [73]: #describing the target column  
data['sentiment'].describe()
```

```
Out[73]: count      50000  
         unique        2  
         top      negative  
         freq      25000  
         Name: sentiment, dtype: object
```

```
In [74]: #checking the count of positives and negatives  
data.sentiment.value_counts()
```

```
Out[74]: negative    25000  
         positive    25000  
         Name: sentiment, dtype: int64
```

Figure 4.1.5.2: Features of sentiment column

- Visualization of the target column i.e. the SENTIMENT column using CATPLOT.



Figure 4.1.5.3: Visualization of the data set

4.1.6. CATEGORICALDATA:

- Machine Learning models are based on equations; we need to replace the text by numbers. So that we can include the numbers in the equations.
- Categorical Variables are of two types: Nominal and Ordinal
- Nominal: The categories do not have any numeric ordering in between them. They don't have any ordered relationship between each of them. Examples: Male or Female, any color
- Ordinal: The categories have a numerical ordering in between them. Example: Graduate is less than Post Graduate; Post Graduate is less than Ph.D. customer satisfaction survey, high low medium.

```

In [138]: #splitting data into input(x)
          X=data.review
          X.head()

Out[138]: 0    One of the other reviewers has mentioned that ...
          1    A wonderful little production. <br /><br />The...
          2    I thought this was a wonderful way to spend ti...
          3    Basically there's a family where a little boy ...
          4    Petter Mattei's "Love in the Time of Money" is...
          Name: review, dtype: object

In [78]: #splitting data into output(y)
          y=data.sentiment
          y.head()

Out[78]: 0    positive
          1    positive
          2    positive
          3    negative
          4    positive
          Name: sentiment, dtype: object

```

Figure 4.1.6.1: Categorical data

4.2. TRAINING THE MODEL:

- Splitting the data : after the preprocessing is done then the data is split into train and test sets
- In Machine Learning in order to access the performance of the classifier. You train the classifier using 'training set' and then test the performance of your classifier on unseen 'test set'. An important point to note is that during training the classifier only uses the training set.
- The test set must not be used during training the classifier. The test set will only be available during testing the classifier.
- training set - a subset to train a model.(Model learns patterns between Input and Output)
- test set - a subset to test the trained model.(To test whether the model has correctly learnt)
- The amount or percentage of Splitting can be taken as specified (i.e. train data = 75% , test data =25% or train data = 80% , test data=20%)
- First we need to identify the input and output variables and we need to separate the input set and output set
- Here we have taken the test size as 0.2 indicating train data =80% and test data=20%

```
In [79]: # split the data
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(data.review, data.sentiment,
                                                    test_size= 0.2, random_state=1 )
```

Figure 4.2.1 : importing train_test_split

- After splitting train and test data the number of columns are as follows:

```
In [80]: #printing the shapes of inputs and outputs
print(X_train.shape)
print(y_train.shape)
print(X_test.shape)
print(y_test.shape)

(40000,)
(40000,)
(10000,)
(10000,)
```

Figure 4.2.2: Train and Test data

- In order to use textual data for predictive modeling, the text must be parsed to remove certain words – this process is called TOKENIZATION. It also enables the pre-processing of text data prior to generating the vector representation.

4.2.1. COUNT VECTORIZATION:

Count Vectorization involves counting the number of occurrences each words appears in a document. Python's Sci-kit learns library has a tool called CountVectorizer to accomplish this.

```
In [81]: ##Importing CountVectorizer  
from sklearn.feature_extraction.text import CountVectorizer  
  
# creating an object for CountVectorizer  
count_vect = CountVectorizer()
```

Figure 4.2.1.1:ImportingCountVectorizer

Generating the word counts for the words in the document and getting the feature names

```
In [82]: # Generate the word counts for the words in the documents  
word_count_vector = count_vect.fit(X_train)  
  
# To get the feature names  
word_count_vector.get_feature_names()  
  
'007s',  
'0080',  
'0083',  
'0093638',  
'00am',  
'00o',  
'00pm',  
'00s',  
'00schneider',  
'01',  
'0148',  
'01pm',
```

Figure 4.2.1.2:Generating the word counts

4.2.2. TFIDF VECTORIZER:

TF-IDF stands for “Term Frequency — Inverse Document Frequency”.TF- IDF are word frequency scores that try to highlight words that are more interesting, e.g. frequent in a document but not across documents. The TfidfVectorizer will tokenize documents, learn the vocabulary and inverse document frequency weightings, and allow you to encode new documents.

$$w_{i,j} = tf_{i,j} \times \log \left(\frac{N}{df_i} \right)$$

$tf_{i,j}$ = number of occurrences of i in j
 df_i = number of documents containing i
 N = total number of documents

Figure 4.2.2.1:TFIDF Formula

The data present in our document is text data.Vectorisation must be performed to convert the text data, TfidfVectoriser is used for this.

```
In [83]: #TFIDF VECTORIZER
from sklearn.feature_extraction.text import TfidfVectorizer
#initialize an object for the tfidf vectorizer
tfidf = TfidfVectorizer()
```

Figure 4.2.2.2: Importing tfidf Vectorizer

Applying tfidf vectorizer to test and train data to objects and generating the feature names. The input data X_train is stored in X_train_transformed; X_test is stored in X_test_transformed after applying Tf-idfvectorizer. After applying tf-idfvectoriser, our input is ‘X_train_transformed’ and its corresponding output is ‘y_train’.It is the same for test data also.

```

In [93]: #apply the tfidf to the data(X_train)

X_train_transformed = tfidf.fit_transform(X_train)
X_train_transformed

Out[93]: <40000x93163 sparse matrix of type '<class 'numpy.float64'>'
        with 5471200 stored elements in Compressed Sparse Row format>

In [94]: X_test_transformed = tfidf.transform(X_test)
X_test_transformed

Out[94]: <10000x93163 sparse matrix of type '<class 'numpy.float64'>'
        with 1345978 stored elements in Compressed Sparse Row format>

In [95]: #Feature names
tfidf.get_feature_names()

Out[95]: ['00',
        '000',
        '00000000000',
        '0000000000001',
        '00001',
        '00015',
        '000dm',
        '001',
        '003830',

```

Figure 4.2.2.3: Vectorisation on train and test data

The position of the words in our data set can be found using `Tfidf.vocabulary_`

```

In [18]: # position of the words
tfidf.vocabulary_

'to': 83403,
'the': 82503,
'hilarious': 38473,
'british': 11525,
'comedy': 16954,
'fish': 30690,
'called': 12851,
'wanda': 89606,
'although': 3626,
'can': 13020,
'see': 72931,
'why': 90722,
'only': 58689,
'connection': 17782,
'find': 30513,
'is': 42841,
'monty': 54568,
'python': 65803,
'one': 58660,
'eric': 27675,

```

Figure 4.2.2.4: Position of words

4.3. EVALUATING THE CASESTUDY:

In our project, we are going to use 3 different algorithms to train and predict the data. Then, we find out the optimum algorithm that best fits our data and gives maximum accuracy. The three algorithms which we will be using are:

1. NAIVE BAYESALGORITHM
2. LOGISTIC REGRESSIONALGORITHM
3. RANDOM FOREST CLASSIFICATION ALGORITHM

4.3.1. NAIVE BAYESALGORITHM:

A Naive Bayes classifier is a supervised machine learning algorithm that uses Baye's theorem which assumes that the features are statistically independent. The theorem relies on the naive assumption that input variables are independent of each other.

The equation of Baye's theorem is:

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)}$$

Figure 4.3.1.1: Naïve Baye's Formula

There are three types of Naive Bayes model under the scikit-learn library:

- **Gaussian:** It is used in classification and it assumes that features follow a normal distribution.
- **Multinomial:** It is used for discrete counts.
- **Bernoulli:** The binomial model is useful if your feature vectors are binary (i.e. zeros and ones). One application would be text classification with a 'bag of words' model where the 1s & 0s are "word occurs in the document" and "word does not occur in the document" respectively.

We use the Bernoulli Naive Bayes model in our project.

Parameters:

- **Priors: array-like of shape (n_classes,):** Prior probabilities of the classes. If specified the priors are not adjusted according to the data.
- **var_smoothing: float, default=1e-9:** Portion of the largest variance of all features that is added to variances for calculation stability. New in version 0.20.

The BernoulliNB method is imported from sklearn.naive_bayes and is assigned to a variable 'model_BernNB'

```
In [20]: # Apply the naive bayes algorithm
         from sklearn.naive_bayes import BernoulliNB
         #creating an object
         model_BernNB = BernoulliNB()
```

Figure 4.3.1.2: Importing Bernoulli Naive Bayes Method

Then, we apply the method to our data and train the model using fit function. We pass X_train_transformed and y_train to train the model.

```
In [21]: # Applying the Algorithm to the data
         # ObjectName.fit(Input, Output)

         model_BernNB.fit(X_train_transformed, y_train)

Out[21]: BernoulliNB(alpha=1.0, binarize=0.0, class_prior=None, fit_prior=True)
```

Figure 4.3.1.3: Fitting the model using BernoulliNb method

Next, we need to predict the data for both training data and testing data and compare the actual value (y_train and y_test) with the model predicted values. For prediction on train data:

```
# Prediction on Train Data
# Syntax: objectname.predict(InputValues)
y_train_pred = model_BernNB.predict(X_train_transformed)
```

Figure 4.3.1.4: Prediction on train data (Naive Bayes)

To compare the actual values of y_train with the predicted values of y_train_pred we use a confusion matrix .

The confusion matrix is imported from sklearn.metrics package.

```
In [23]: # compare the actual values(y_train) with predicted values(y_train_pred)
from sklearn.metrics import confusion_matrix , classification_report, accuracy_score
confusion_matrix(y_train,y_train_pred)

Out[23]: array([[18376, 1580],
               [ 2494, 17550]], dtype=int64)
```

Figure 4.3.1.5: Confusion matrix for train data(Naive Bayes)

The classification report and accuracy is printed as:

```
In [155]: #printing the classification report
print(classification_report(y_train,y_train_pred))
```

	precision	recall	f1-score	support
negative	0.88	0.92	0.90	19956
positive	0.92	0.88	0.90	20044
accuracy			0.90	40000
macro avg	0.90	0.90	0.90	40000
weighted avg	0.90	0.90	0.90	40000

```
In [156]: #checking the accuracy fot training data
train_accu[0]=accuracy_score(y_train,y_train_pred)
print("Naive Baye's Train Accuracy: ",train_accu[0])

Naive Baye's Train Accuracy: 0.89815
```

Figure 4.3.1.6: Classification report for train data(Naive Bayes)

To visualize the confusion matrix, we use heat map



Figure 4.3.1.7: Heat Map of train data(Naive Bayes)

We can conclude that the accuracy score of Naive Bayes on train data is **0.89**

For prediction on test data:

```
In [26]: # prediction on test data
# syntax: objectname.predict(input)
y_test_pred = model_BernNB.predict(X_test_transformed)
```

Figure 4.3.1.8: Prediction on test data(Naive Bayes)

To compare the actual values of `y_test` with the predicted values of `y_test_pred` we use a confusion matrix .

The confusion matrix is imported from `sklearn.metrics` package.

```
In [27]: # compare the actual values(y_test) with predicted values(y_test_pred)
confusion_matrix(y_test,y_test_pred)
```

Out[27]: array([[4431, 613],
 [780, 4176]], dtype=int64)

Figure 4.3.1 9:Confusion matrix for test data(Naive Bayes)

The classification report and accuracy is printed as

```
In [160]: #checking the accuracy for testing data
test_accu[0]=accuracy_score(y_test,y_test_pred)
print("Naive Baye's Test Accuracy: ",test_accu[0])

Naive Baye's Test Accuracy:  0.8607

In [161]: #printing the classification report
print(classification_report(y_test,y_test_pred))
```

	precision	recall	f1-score	support
negative	0.85	0.88	0.86	5044
positive	0.87	0.84	0.86	4956
accuracy			0.86	10000
macro avg	0.86	0.86	0.86	10000
weighted avg	0.86	0.86	0.86	10000

Figure 4.3.1.10: Classification report for test data(Naive Bayes)

To visualize the confusion matrix, we use



Figure 4.3.1.11: Heat Map of test data(Naive Bayes)

We can conclude that the accuracy score of Naive Bayes on test data is 0.86

4.3.2. LOGISTIC REGRESSION:

Logistic regression is another technique of machine learning from the field of statistics. It is the go-to method for binary classification problems (problems with two class values). Logistic regression is used to describe data and to explain the relationship between one dependent binary variable and one or more nominal, ordinal, interval or ratio-level independent variables.

There are three main types of logistic regression:

- **binomial:** target variable can have only 2 possible types: “0” or “1” which may represent “win” vs. “loss”, “pass” vs “fail”, “dead” vs. “alive”, etc.
- **Multinomial:** target variable can have 3 or more possible types which are not ordered (i.e. types have no quantitative significance) like “disease A” vs “disease B” vs “disease C”.
- **Ordinal:** it deals with target variables with ordered categories. For example, a test score can be categorized as: “very poor”, “poor”, “good”, “and very good”. Here, each category can be given a score like 0, 1, 2, and 3.

Parameters:

- **Penalty:** {‘l1’, ‘l2’, ‘elastic net’, ‘none’}, default=‘l2’
Used to specify the norm used in the penalization. The ‘Newton-cg’, ‘sag’ and ‘lbfgs’ solvers support only l2 penalties. ‘Elastic net’ is only supported by the ‘saga’ solver. If ‘none’ (not supported by the liblinear solver), no regularization is applied. New in version 0.19: l1 penalty with SAGA solver (allowing ‘multinomial’ + L1)
- **Dual:** bool, default=False
 - Dual or primal formulation. Dual formulation is only implemented for l2 penalty with liblinear solver. Prefer dual=False when n_samples>n_features.
- **Tol:** float, default=1e-4
 - Tolerance for stopping criteria.
- **C:** float, default=1.0
 - Inverse of regularization strength; must be a positive float. Like in support vector machines, smaller values specify stronger regularization.

First we import the 'LogisticRegression' method from the 'sklearn.linear_model' package.

Then we create an object called LogisticRegression () and store it in a variable called 'reg'

Next we apply the algorithm to the data and fit the data using the syntax:
objectname.fit(input, output)

```
In [31]: # build the classifier on training data
#sklearn library: import,instantiate,fit
from sklearn.linear_model import LogisticRegression
reg = LogisticRegression()
reg.fit(X_train_transformed,y_train) #input and output will be passed to fit method

C:\Users\Varun\anaconda3\lib\site-packages\sklearn\linear_model\logistic.py:432: FutureWarning: Default solver will be changed
to 'lbfgs' in 0.22. Specify a solver to silence this warning.
  FutureWarning)
```

Figure 4.3.2.1:Importing Logistic Regression Method and fitting the model

Next , we need to predict the data for both training data and testing data and compare the actual value(y_train and y_test) with the model predicted values.

For prediction on train data:

```
In [32]: #predicting on train data
# syntax : objectname.predict(Input)
y_train_pred =reg.predict(X_train_transformed)
y_train_pred

Out[32]: array(['negative', 'positive', 'positive', ..., 'negative', 'negative',
                'positive'], dtype=object)
```

Figure 4.3.2 2: Prediction on train data (Logistic Regression)

To compare the actual values of y_train with the predicted values of y_train_pred we use a confusion matrix .

The confusion matrix is imported from sklearn.metrics package.

```
In [33]: # confusion matrix for the training data
# confusion matrix(actual values,predicted values)
from sklearn.metrics import confusion_matrix , accuracy_score
conf = confusion_matrix(y_train,y_train_pred)
conf

Out[33]: array([[18422, 1534],
               [ 1239, 18805]], dtype=int64)
```

Figure 4.3.2.3: Confusion matrix for train data (Logistic Regression)

The classification report and accuracy is printed as:

```
In [166]: #checking the accuracy fot training data
train_accu[1]=accuracy_score(y_train,y_train_pred)
print("LogisticRegression Train Accuracy: ",train_accu[1])

LogisticRegression Train Accuracy:  0.930675

In [167]: #printing the classification report
from sklearn.metrics import classification_report
print(classification_report(y_train,y_train_pred))
```

	precision	recall	f1-score	support
negative	0.94	0.92	0.93	19956
positive	0.92	0.94	0.93	20044
accuracy			0.93	40000
macro avg	0.93	0.93	0.93	40000
weighted avg	0.93	0.93	0.93	40000

Figure 4.3.2.4: Classification report for train data (Logistic Regression)

To visualize the confusion matrix, we use heat map

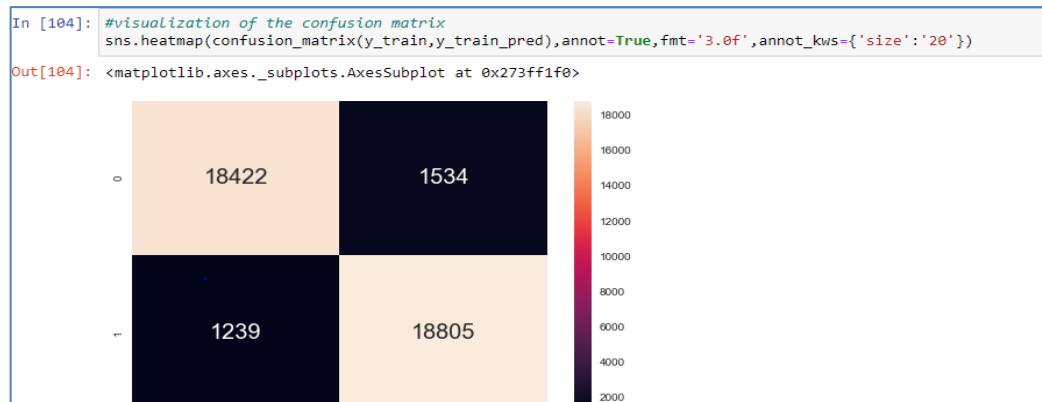


Figure 4.3.2.5: Heat Map of train data (Logistic Regression)

We can conclude that the accuracy score of Logistic Regression on train data is **0.93**

For prediction on test data:

```
In [37]: #predicting on testing data
y_test_pred = reg.predict(X_test_transformed)
y_test_pred
```

Out[37]: array(['negative', 'negative', 'negative', ..., 'negative', 'positive',
 'negative'], dtype=object)

Figure 4.3.2.6: Prediction on test data (Logistic Regression)

To compare the actual values of `y_test` with the predicted values of `y_test_pred` we use a confusion matrix.

The confusion matrix is imported from `sklearn.metrics` package.

```
In [108]: #printing the confusion matrix
conf = confusion_matrix(y_test,y_test_pred)
conf
```

Out[108]: array([[4482, 562],
 [454, 4502]], dtype=int64)

Figure 4.3.2.7: Confusion matrix for test data (Logistic Regression)

The classification report and accuracy is printed as:

```
In [171]: #checking the accuracy for testing data
test_accu[1]=accuracy_score(y_test,y_test_pred)
print("LogisticRegression Test Accuracy: ",test_accu[1])
LogisticRegression Test Accuracy: 0.8984

In [172]: #printing the classification report
print(classification_report(y_test,y_test_pred))
```

	precision	recall	f1-score	support
negative	0.91	0.89	0.90	5044
positive	0.89	0.91	0.90	4956
accuracy			0.90	10000
macro avg	0.90	0.90	0.90	10000
weighted avg	0.90	0.90	0.90	10000

Figure 4.3.2 8: Classification report for test data (Logistic Regression)

To visualize the confusion matrix, we use heat map



Figure 4.3.2.9: Heat Map of test data (Logistic Regression)

We can conclude that the accuracy score of Logistic Regression on test data is **0.89**

4.3.3. RANDOM FOREST CLASSIFICATION:

The random forest is a classification algorithm consisting of many decisions trees. It uses bagging and feature randomness when building each individual tree to try to create an uncorrelated forest of trees whose prediction by committee is more accurate than that of any individual tree. The random forest algorithm is a supervised learning model; it uses labeled data to “learn” how to classify unlabeled.

First we import the ‘RandomForestClassifier’ method from the ‘sklearn.ensemble’ package Then we create an object called RandomForestClassifier() and store it in a variable called ‘rfc’

Next we apply the algorithm to the data and fit the data using the syntax `objectname.fit(input,output)`.

Parameters:

- **n_estimators**:int, default=100: The number of trees in the forest.

Criterion: {“gini”, “entropy”}, default=“gin”

The function to measure the quality of a split. Supported criteria are “gini” for the Gini impurity and “entropy” for the information gain. Note: this parameter is tree-specific.

- **max_depth**:int, default=None

The maximum depth of the tree. If None, then nodes are expanded until all leaves are pure or until all leaves contain less than min_samples_split samples.

- **min_samples_split**:int or float, default=2

The minimum number of samples required to split an internal node:

```

In [42]: #Import, initialize and fit

#Import the RfC from sklearn
from sklearn.ensemble import RandomForestClassifier

#Initilaize the object for RFC
rfc = RandomForestClassifier(n_estimators = 40)

#fit the TFC to the dataset
rfc.fit(X_train_transformed,y_train)

Out[42]: RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
                                max_depth=None, max_features='auto', max_leaf_nodes=None,
                                min_impurity_decrease=0.0, min_impurity_split=None,
                                min_samples_leaf=1, min_samples_split=2,
                                min_weight_fraction_leaf=0.0, n_estimators=40,
                                n_jobs=None, oob_score=False, random_state=None,
                                verbose=0, warm_start=False)

```

Figure 4.3.3.1:Importing Random Forest Classifier Method and fitting the model

Next , we need to predict the data for both training data and testing data and compare the actual value(y_train and y_test) with the model predicted values.

For prediction on train data:

```

In [43]: #Prediction on training data
#Syntax: ojectname.predict(InputValues)
y_train_pred = rfc.predict(X_train_transformed)

from sklearn.metrics import confusion_matrix, classification_report
print(classification_report(y_train, y_train_pred))

```

	precision	recall	f1-score	support
negative	1.00	1.00	1.00	19956
positive	1.00	1.00	1.00	20044
accuracy			1.00	40000
macro avg	1.00	1.00	1.00	40000
weighted avg	1.00	1.00	1.00	40000

Figure 4.3.3.2: Prediction on train data and classification report(Random Forest Classifier)

The accuracy is printed as:

```
In [176]: #checking the accuracy fot training data
train_accu[2]=accuracy_score(y_train,y_train_pred)

print("Random Forest Classification Train Accuracy: ",train_accu[2])

Random Forest Classification Train Accuracy:  1.0
```

Figure 4.3.3.3: Accuracy for train data(Random Forest Classification)

To visualize the confusion matrix, we use heat map

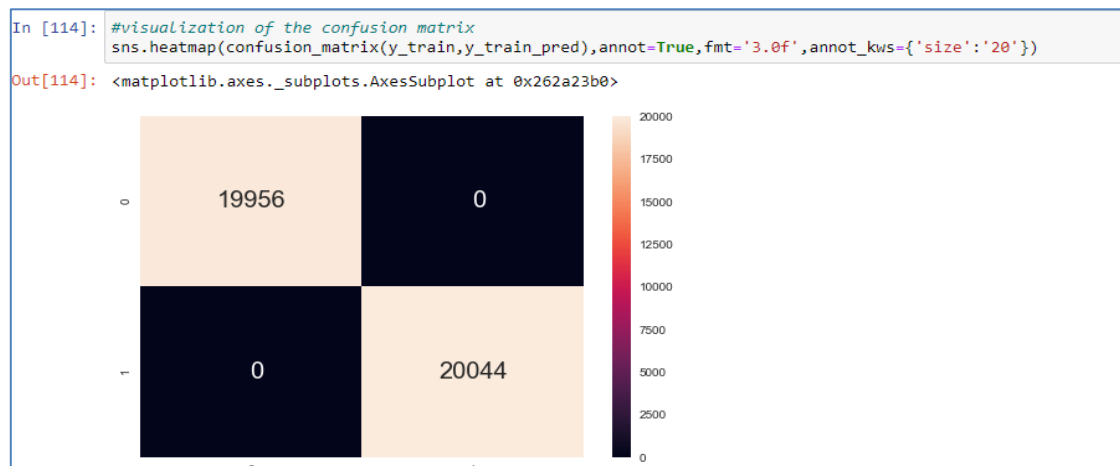


Figure 4.3.3.4: Heat Map of train data (Random Forest Classifier)

We can conclude that the accuracy score of Random Forest Classifier on train data is **1.0**

For prediction on test data:

```
In [116]: #Prediction on testing data
y_test_pred = rfc.predict(X_test_transformed)
#printing the classification report
print(classification_report(y_test, y_test_pred))
```

	precision	recall	f1-score	support
negative	0.82	0.84	0.83	5044
positive	0.83	0.81	0.82	4956
accuracy			0.82	10000
macro avg	0.82	0.82	0.82	10000
weighted avg	0.82	0.82	0.82	10000

Figure 4.3.3.5: Prediction on test data and classification report (Random Forest Classifier)

The accuracy is printed as:

```
In [179]: #checking the accuracy for testing data
test_accu[2]=accuracy_score(y_test,y_test_pred)
print("Random Forest Classification Test Accuracy: ",test_accu[2])

Random Forest Classification Test Accuracy: 0.8214
```

Figure 4.3.3.6: Accuracy for test data (Random Forest Classification)

To visualize the confusion matrix, we use heat map



Figure 4.3.3.7: Heat Map of test data (Random Forest Classifier)

We can conclude that the accuracy score of Random Forest Classifier on test data is **0.82**

4.4. Comparison of accuracy scores of train and test data:

The accuracy scores of Naive Bayes Algorithm, Logistic Regression and Random Forest Classification on train data are 0.89, 0.93, 1 and on test data are 0.86, 0.89 and 0.82

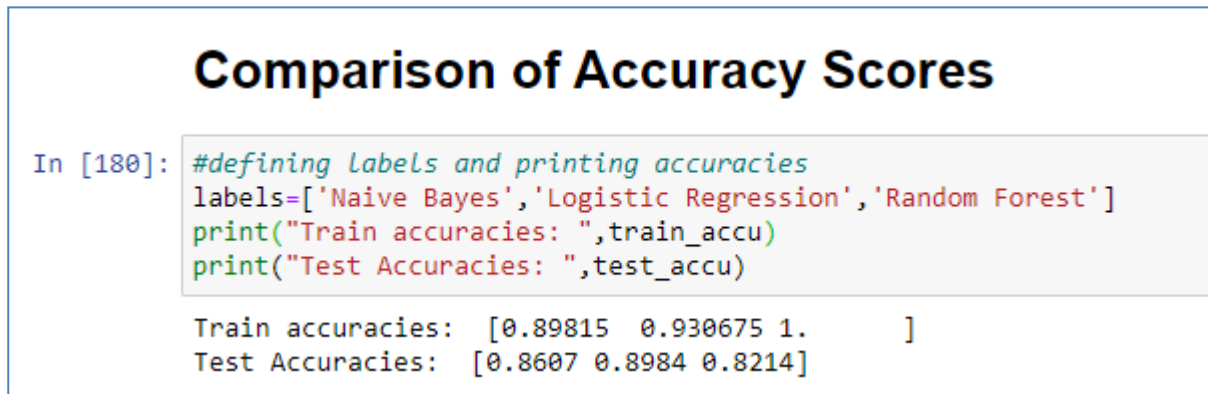


Figure 4.4.1: defining labels and printing accuracies

```
In [182]: #visualizing the train accuracies of all the models
print("Accuracy Comparision for TRAIN data for all the models:\n\n")
plt.subplots(figsize=(5,5))
plt.bar(labels,train_accu,color=['green','red','black'])
plt.xlabel('Classifiers')
plt.ylabel('Accuracy')
plt.title('Train Accuracy of various algorithms')
```

Accuracy Comparision for TRAIN data for all the models:

: Text(0.5, 1.0, 'Train Accuracy of various algorithms')



Figure 4.4.2: Visulaising the train accuracies

```
In [183]: #visualizing the test accuracies of all the models
print("Accuracy Comparision for TEST data for all the models:\n\n")
plt.subplots(figsize=(5,5))
plt.bar(labels,test_accu,color=['green','red','black'])
plt.xlabel('Classifiers')
plt.ylabel('Accuracy')
plt.title('Test Accuracy of various algorithms')
```

Accuracy Comparision for TEST data for all the models:

```
Out[183]: Text(0.5, 1.0, 'Test Accuracy of various algorithms')
```

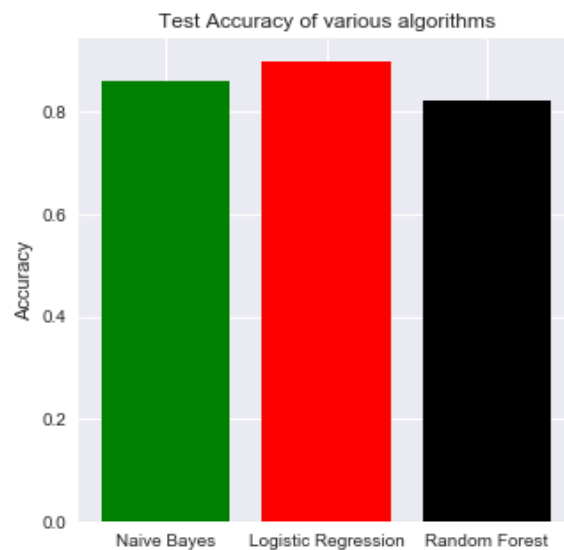


Figure 4.4.3: Visulaising the test accuracies

We can infer that Logistic Regression algorithm has the highest accuracy with 0.89 on test data among all the 3 algorithms

So we perform Hyper Parameter tuning on logistic regression to increase the accuracy score

4.5. HYPER PARAMETER TUNING ON LOGISTIC REGRESSION

Hyper Parameter Tuning is used to increase the accuracy score of the algorithm by finding out the best parameters. We now perform Hyper Parameter tuning on Logistic Regression as it has the highest accuracy score among all the 3 methods applied on the data set.

First, we pass a list of values in a dictionary to find the optimum value for each parameter and assign it to a variable 'grid_param'

```
In [49]: #Passing list of values in a dictionary to find the optimum value for each parameter
grid_param = {
    'penalty': ['l1', 'l2'],
    'C': [0.001, 0.01, 0.1, 1, 10, 100, 1000],
    'intercept_scaling': [0.001, 0.01, 0.1, 1, 10, 100, 1000],
    'fit_intercept': ['True', 'False']
}
```

Figure 4.5.1: Passing a list of values as dictionary

Now, we import GridSearchCV from sklearn.model.selection. Then the initialization of Grid Search with the parameters- Model Name and the dictionary of parameters is done. Finally we apply grid search onto dataset.

```
In [50]: #Import the GridSearchCV
from sklearn.model_selection import GridSearchCV

# initialization of GridSearch with the parameters- ModelName and the dictionary of parameters
clf = LogisticRegression()
grid_search = GridSearchCV(estimator=clf, param_grid=grid_param)

# applying gridsearch onto dataset
grid_search.fit(X_train_transformed, y_train)

C:\Users\Varun\anaconda3\lib\site-packages\sklearn\model_selection\_split.py:1978: FutureWarning: The default value of cv will
change from 3 to 5 in version 0.22. Specify it explicitly to silence this warning.
  warnings.warn(CV_WARNING, FutureWarning)
C:\Users\Varun\anaconda3\lib\site-packages\sklearn\linear_model\logistic.py:432: FutureWarning: Default solver will be changed
to 'lbfgs' in 0.22. Specify a solver to silence this warning.
  FutureWarning)

Out[50]: GridSearchCV(cv='warn', error_score='raise-deprecating',
```

Figure 4.5.2: Importing GridSearchCV and passing the dictionary of parameters.

We then print the best parameters as:

```
In [51]: # return the optimal parameters
        grid_search.best_params_

Out[51]: {'C': 10, 'fit_intercept': 'True', 'intercept_scaling': 0.01, 'penalty': 'l2'}
```

Figure 4.5.3: Printing the best parameters

Now, we build the model with the best parameters and fit the model with train data.

```
In [52]: #Build the model with best parameters
        #Initialized the DT Classifier
        clf = LogisticRegression(C=10, fit_intercept=True, intercept_scaling=0.01, penalty='l2')

        # We need to fit the model to the data
        clf.fit(X_train_transformed, y_train)

Out[52]: LogisticRegression(C=10, class_weight=None, dual=False, fit_intercept=True,
        intercept_scaling=0.01, l1_ratio=None, max_iter=100,
        multi_class='warn', n_jobs=None, penalty='l2',
        random_state=None, solver='warn', tol=0.0001, verbose=0,
        warm_start=False)
```

Figure 4.5.4: Building the model with best parameters

Prediction on test data

```
In [53]: #predicting on testing data
        y_test_pred = clf.predict(X_test_transformed)
        y_test_pred

Out[53]: array(['negative', 'negative', 'negative', ..., 'positive', 'positive',
        'negative'], dtype=object)
```

Figure 4.5.5: Prediction on test data

To compare the actual values of `y_test` with the predicted values of `y_test_pred` we use a confusion matrix

```
In [124]: #printing the confusion matrix
        conf = confusion_matrix(y_test, y_test_pred)
        conf

Out[124]: array([[4522,  522],
        [ 446, 4510]], dtype=int64)
```

Figure 4.5.6: Confusion Matrix of Test data (Hyper Parameter Tuning)

The classification report and accuracy scores are printed as:

```
In [126]: #checking the accuracy for testing data
accuracy_score(y_test,y_test_pred)

Out[126]: 0.9032

In [127]: #printing the classification report
print(classification_report(y_test,y_test_pred))
```

	precision	recall	f1-score	support
negative	0.91	0.90	0.90	5044
positive	0.90	0.91	0.90	4956
accuracy			0.90	10000
macro avg	0.90	0.90	0.90	10000
weighted avg	0.90	0.90	0.90	10000

Figure 4.5.7: Classification report and accuracy of test data (Hyper Parameter Tuning)

To visualize the confusion matrix, we use heat map



Figure 4.5.8: Heat Map of test data (Hyper Parameter Tuning)

We can clearly see that the accuracy score of Logistic Regression on test data after hyper parameter tuning has increased to 0.90 from 0.89

4.6. Testing with unknown input:

We take an unknown review and take it as a list.

```
In [184]: ## Sample Input  
new=["Beautiful Mind is a biographical movie about John Forbes Nash junior, mathematical genius with hard fate. At the beginn
```

Figure 4.6.1: Reading unknown Input

We perform vectorization on 'new' and predict the review

```
In [188]: ## Applying tfidf vectorization  
X_new_transformed = tfidf.transform(new)  
  
In [189]: ## Predicting the sample input  
y_new_pred = clf.predict(X_new_transformed)  
y_new_pred  
  
Out[189]: array(['positive'], dtype=object)
```

Figure 4.6.2: Predicting the unknown input

We see that our model has predicted the review as positive review.

4.7. ROC AUC Curve:

After building the model using the above three algorithms and hyper parameter tuning, to find the best fit among the three, we visualize the data using the roc-curve and auc_score

ROC Curve: An ROC curve (receiver operating characteristic curve) is a graph showing the performance of a classification model at all classification thresholds. This curve plots two parameters: True Positive Rate. False Positive Rate.

AUC Curve: Area under the ROC Curve :

AUC stands for “Area under the ROC Curve.” That is, AUC measures the entire two-dimensional area underneath the entire ROC curve.

Accuracy is measured by the area under the ROC curve. An area of 1 represents a perfect test; an area of 0.5 represent a worthless test.

Any classifier can be used to create an ROC curve, but it must be able to fit and predict_proba.

Here, we use:

- sklearn.naive_bayes.BernoulliNB,
- sklearn.linear_model.LogisticRegression,
- sklearn.ensemble.RandomForestClassifier

```
In [128]: #training the models
from sklearn.naive_bayes import BernoulliNB
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier

In [129]: #naive bayes,logistic regression(after hyper parameter tuning),random forest
model_BernNB = BernoulliNB()
clf = LogisticRegression(C=10, fit_intercept=True, intercept_scaling=0.01, penalty='l2')
rfc = RandomForestClassifier(n_estimators = 40)

In [130]: #fitting the models
model_BernNB.fit(X_train_transformed,y_train)
clf.fit(X_train_transformed,y_train)
rfc.fit(X_train_transformed,y_train)

C:\Users\Varun\anaconda3\lib\site-packages\sklearn\linear_model\logistic.py:432: FutureWarning: Default solver will be changed
to 'lbfgs' in 0.22. Specify a solver to silence this warning.
  FutureWarning)

Out[130]: RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
max_depth=None, max_features='auto', max_leaf_nodes=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, n_estimators=40,
n_jobs=None, oob_score=False, random_state=None,
verbose=0, warm_start=False)
```

Figure 4.7.1: Importing packages for Roc-Auc Curve

Here, we find out the predicted probabilities using 'predict_proba' and store them in 'pred_prob1','pred_prob2','pred_prob3' for the respective three models.

Next , we import the 'roc_curve' from the 'sklearn.metrics' package and pass the fpr, tpr, threshold values into them.

```
In [131]: #predicting the probabilities
pred_prob1=model_BernNB.predict_proba(X_test_transformed)
pred_prob2=clf.predict_proba(X_test_transformed)
pred_prob3=rfc.predict_proba(X_test_transformed)

In [132]: from sklearn.metrics import roc_curve

# roc curve for models
fpr1,tpr1,thresh1 = roc_curve(y_test_pred,pred_prob1[:,1],pos_label='positive')
fpr2,tpr2,thresh2 = roc_curve(y_test_pred,pred_prob2[:,1],pos_label='positive')
fpr3,tpr3,thresh3 = roc_curve(y_test_pred,pred_prob3[:,1],pos_label='positive')
```

Figure 4.7.2: Calculating the predicted probabilities and importing roc auc curve

We calculate the auc scores as:

```
In [63]: from sklearn.metrics import roc_auc_score

# auc scores
auc_score1 = roc_auc_score(y_test_pred,pred_prob1[:,1])
auc_score2 = roc_auc_score(y_test_pred,pred_prob2[:,1])
auc_score3 = roc_auc_score(y_test_pred,pred_prob3[:,1])

print(auc_score1, auc_score2, auc_score3)

0.9505193332718909 1.0 0.9271553562833933

In [64]: import matplotlib.pyplot as plt
```

Figure 4.7.3: Calculating the auc scores

Now, we plot the roc-auc curve using pyplot from matplotlib

```
In [64]: import matplotlib.pyplot as plt
plt.style.use('seaborn')

#plot roc curves
plt.plot(fpr1,tpr1,linestyle='-',color='red',label='Naive Bayes')
plt.plot(fpr2,tpr2,linestyle='-',color='green',label='Logistic Regression After Hyper Parameter Tuning')
plt.plot(fpr3,tpr3,linestyle='-',color='blue',label='Random Forest Classification')

#title
plt.title('ROC CURVE')
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.legend(loc='best')
plt.savefig('ROC',dpi=300)
plt.show()
```

Figure 4.7.4:PlottingROC Curve

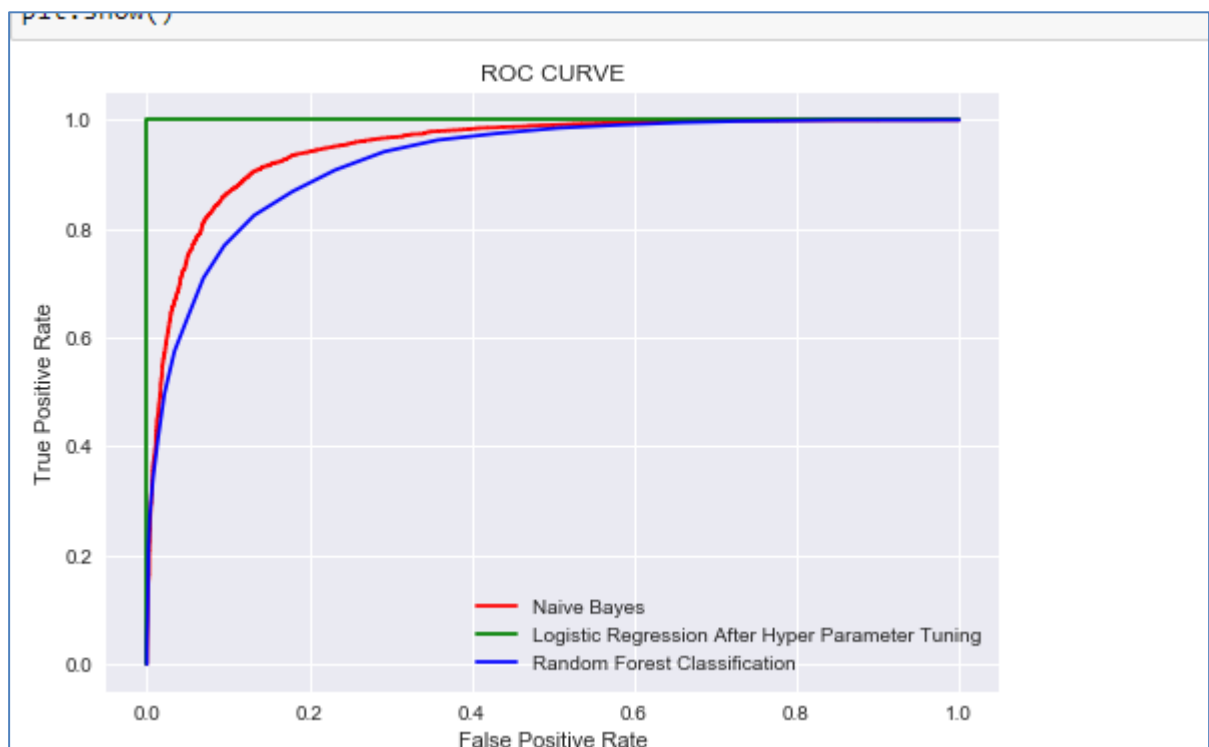


Figure 4.7.5: ROC Curve

CONCLUSION:

From the above results, we can infer that for our problem statement, the Logistic Regression Model (with hyper parameter tuning) gave the highest accuracy than the other two models Naive bayes and random forest classification. Apart from this, we can see that the next best algorithm is the Naive Bayes. Then is the random forest classifier. Heat maps are computed to get a clear understanding of the data set (which parameter has most abundant effect on the study case) .The general order of performance for the model was

LogisticRegression (Hyper Parameter Tuning) >NaïveBayes>RandomForestClassifier

REFERENCES:

1. DATASET:<https://drive.google.com/drive/folders/1vZgoh1o5xT6vehsRdfAQEX4nA79yoSVe?usp=sharing>
2. <https://www.analyticsvidhya.com/blog/2020/06/auc-roc-curve-machine-learning/>
3. https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html
4. https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.BernoulliNB.html?highlight=naive%20bayes#sklearn.naive_bayes.BernoulliNB
5. <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html?highlight=random%20forest%20classifier#sklearn.ensemble.RandomForestClassifier>
6. GIT HUB Link: <https://github.com/varun2704/Movie-Review-Classification>