

A PROJECT REPORT ON
SENTIMENT ANALYSIS USING MACHINE LEARNING

Submitted to

JAWAHARLAL NEHRU TECHNOLOGICAL UNIVERSITY, KAKINADA

Partial Fulfillment of Award of the Degree of

Master of Computer Applications

Submitted By

M. RAVINDRAREDDY (19X45F0002)

Under the esteemed guidance of

Mr. M V. SUMANTH

Assistant Professor, Department of CSE



DEPARTMENT OF MASTER OF COMPUTER APPLICATIONS

S.R.K INSTITUTE OF TECHNOLOGY
(AFFILIATED TO JNTU, KAKINADA)

Enikepadu, Vijayawada – 521108.

JULY-2021

S.R.K INSTITUTE OF TECHNOLOGY
DEPARTMENT OF MASTER OF COMPUTER APPLICATIONS



CERTIFICATE

This is to certify that this project report entitled “**SENTIMENT ANALYSIS USING MACHINE LEARNING**” is the bonafide work of **Mr. M. RAVINDRAREDDY (19X45F0002)** in partial fulfillment of the requirements for the award of the post graduate degree in MASTER OF COMPUTER APPLICATIONS during the academic year 2020-2021. This Work has carried out under our supervision and guidance.

Signature of the Guide

(Mr. M V.Sumanth)

Signature of the HOD

(Dr D. Haritha)

Signature of the External Examiner

DECLARATION

I **M. Ravindrareddy** hereby declare that the project report entitled “**SENTIMENT ANALYSIS USING MACHINE LEARNING**” is an original work done in the Department of Master of Computer Applications, SRK Institute of Technology, Enikepadu, Vijayawada, during the academic year 2020-2021, in partial fulfillment for the award of the Degree of Post Graduation in Master of Computer Applications. We assure that this project is not submitted to another College or University.

Roll No

19X45F0002

Name of the Student

M. RAVINDRAREDDY

Signature

ACKNOWLEDGEMENT

Firstly I would like to convey my heart full thanks to the Almighty for the blessings on me to carry out this project work without any disruption.

I am extremely thankful to **Mr. M V. Sumanth, Project guide** who guided me throughout the project. I am thankful to her for giving me the most independence and freedom throughout various phases of the project.

I am also thankful for my projects coordinator **M.Anitha** for her valuable guidance which helped me to bring this project successfully.

I am very much grateful to **Dr. D. Haritha**, H.O.D of M.C.A Department, for her valuable guidance which helped me to bring out this project successfully. Her wise approach made me to learn the minute details of the subject. Her matured and patient guidance paved a way for completing my project with the sense of satisfaction and pleasure.

I am very much thankful to our principal **Dr. M. Ekambaram Naidu** for his kind support and facilities provided at our campus which helped me to bring out this project successfully.

Finally, I would like to convey my heart full thanks to all Technical Staff, for their guidance and support in every step of this project. I convey my sincere thanks to all the faculty and friends who directly or indirectly helped me for the successful completion of this project.

PROJECT ASSOCIATE

M. RAVINDRAREDDY (19X45F0002)

ABSTRACT

Sentiment analysis is that the process of analyzing the emotion of the users. You'll categorize their emotions as positive, negative or neutral. It's widely being employed nowadays. The explanation behind this can be every company is trying to know the sentiment of their customers, if customers are happy, they'll stay. This project could show a path to scale back customer churn.

Google CoLab device investigates text for extremity from positive to negative. Via preparing model devices with tests of feelings in text, machines consequently discover how to distinguish notion without human information. Assumption investigation models are prepared to peruse past the definitions. There are number of methods and muddled calculations will not to order and Train machines to perform slant examination. LSTM and CNN calculation is utilized to mentor and order text inside our feeling extremity model.

The Data Set IMDB is split into test and training set. This can be a binary classification dataset as positive and negative. In training set contains 25000 records, which incorporates 20000 positive records and 5000 negative records. Test dataset has a complete of 25000 records. Out of which 20000 records are classified as positive and 5000 records as negative.

TABLE OF CONTENTS

Contents	Page No
1. Introduction	1
2. System Analysis	
2.1 Existing System	5
2.2 Proposed System	5
2.3 Literature Review	5
2.4 Module Description	
2.4.1 Admin Module	9
2.4.2 System Module	9
2.4.3 User Module	9
2.5 Hardware Requirements	9
2.6 Software Requirements	9
2.7 Feasibility Study:	
2.7.1 Technical Feasibility	10
2.7.2 Operational Feasibility	10
2.7.3 Economical Feasibility	10
2.8 Functional Requirements	10
3. System Design	
3.1 System architecture	11
3.2 UML diagrams:	13
3.2.1 Class Diagram	14
3.2.2 Use Case diagram	15
3.2.3 Sequence diagram	16
3.3 Datasets	17
3.4 Technologies Description	
3.4.1 Python	18
3.4.2 Ensemble Methods	21
4. Coding and implementation	24

5. Output Screens	30
6. Evaluation Metrics	34
7. Testing	39
8 .Conclusion	41
9. Future Enhancements	42
10.Bibliography	43

LIST OF FIGURES

Figure 1.1	Sentiment Analysis	01
Figure 1.2	Process of determining polarities I	02
Figure 1.3	Process of determining polarities II	03
Figure 1.4	Determining Polarities	04
Figure 3.1	System Architecture	10
Figure 3.2	Class Diagram	14
Figure 3.3	Use Case Diagram	15
Figure 3.4	Sequence Diagram	16
Figure 3.5	Data Sets	17
Figure 5.1	Loading the data and removing stop words	30
Figure 5.2	Training the data	30
Figure 5.3	Summary of IMDB dataset	31
Figure 5.4	Conv1d model Accuracy	31
Figure 5.5	Evaluating the data	32
Figure 5.6	Final Accuracy of each review	33
Figure 6.1	Positive Rate curve	37
Figure 7.1	Test Cases	40

1.INTRODUCTION

Slant Analysis manages connection among clients and item audits used by people momentarily we are prepared to say preparing machine in understanding to the matter assertion we'll determine fine quality data by utilizing basic content entered by the client through various examples and patterns driving towards the yield we expect by the different assessments and translations, consequently we sort and group our content with in regard to the matter assertion we are taking care of. Subsequently we register our yield by utilizing the information set we've utilizing various models, calculations, numerical calculations with goes under the classification of phonetics. This field is especially applied once we ought to consistently take audits or overview from our clients on items or administrations. Wistful Analysis can be a pristine variation inside the exploration region. It's fundamentally alludes to as feelings or perspectives on the differed information that is being gathered utilizing studies, remarks and surveys over the web.



Fig 1.1: Sentimental Analysis

Huge measure of data is being produced day by day which is prepared by utilizing language handling, text examination and semantics. Assessment and notion of the clients are expanding fastly by utilizing various sorts of strategies and calculations. It decides the extremity of audits and these surveys are named two sorts positive and negative. Audits of IMDB are utilized as Data and various calculations like LSTM, CNN and Ensemble classifier on various information streams has been given.

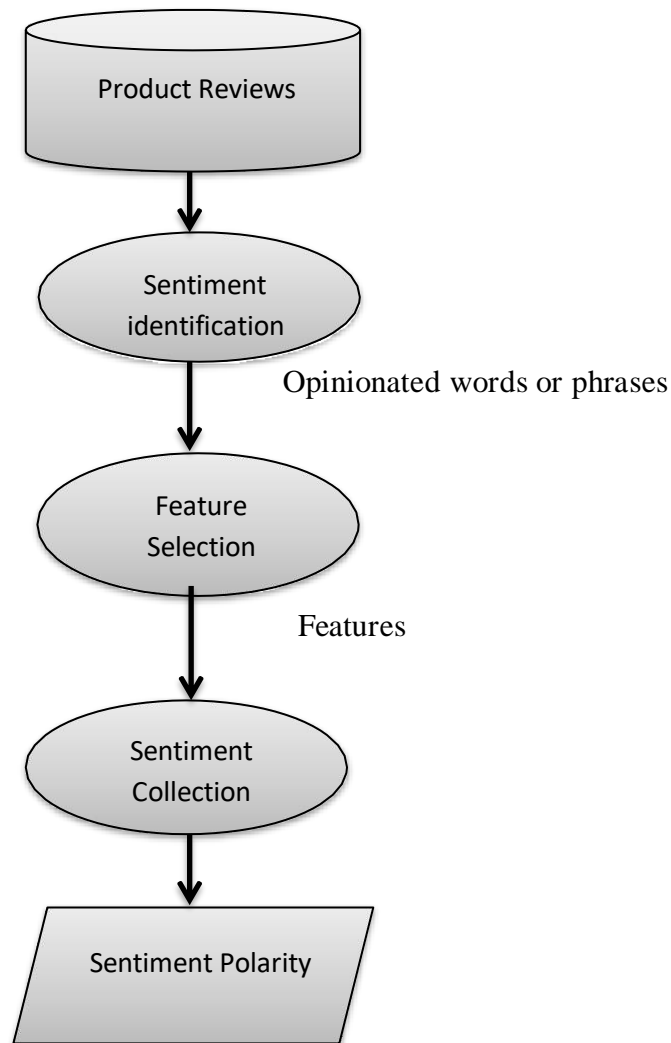


Fig 1.2: Process of determining polarities I

It likewise decides the disposition of an individual on some theme or his\her passionate response about some occurrence, records or occasions. Slant Analysis chips away at the standard of Machine learning, so we took the varying ideas and calculations of AI, taken a stab at going along with them so we are prepared to arrive at the objective of our undertaking. With the expansion in innovation, different online media stages like twitter face book, Instagram, connected in and a significant number of other.

These stages contain colossal measure of information being produced day by day inside the arranging of tweets, sites, posts, status and so on Wistful examination predicts the temperament of those writings, tweets, audits or posts which are accessible online on the stage by deciding the extremity of feelings like joy, fondness, sorrow, outrage and scorn.

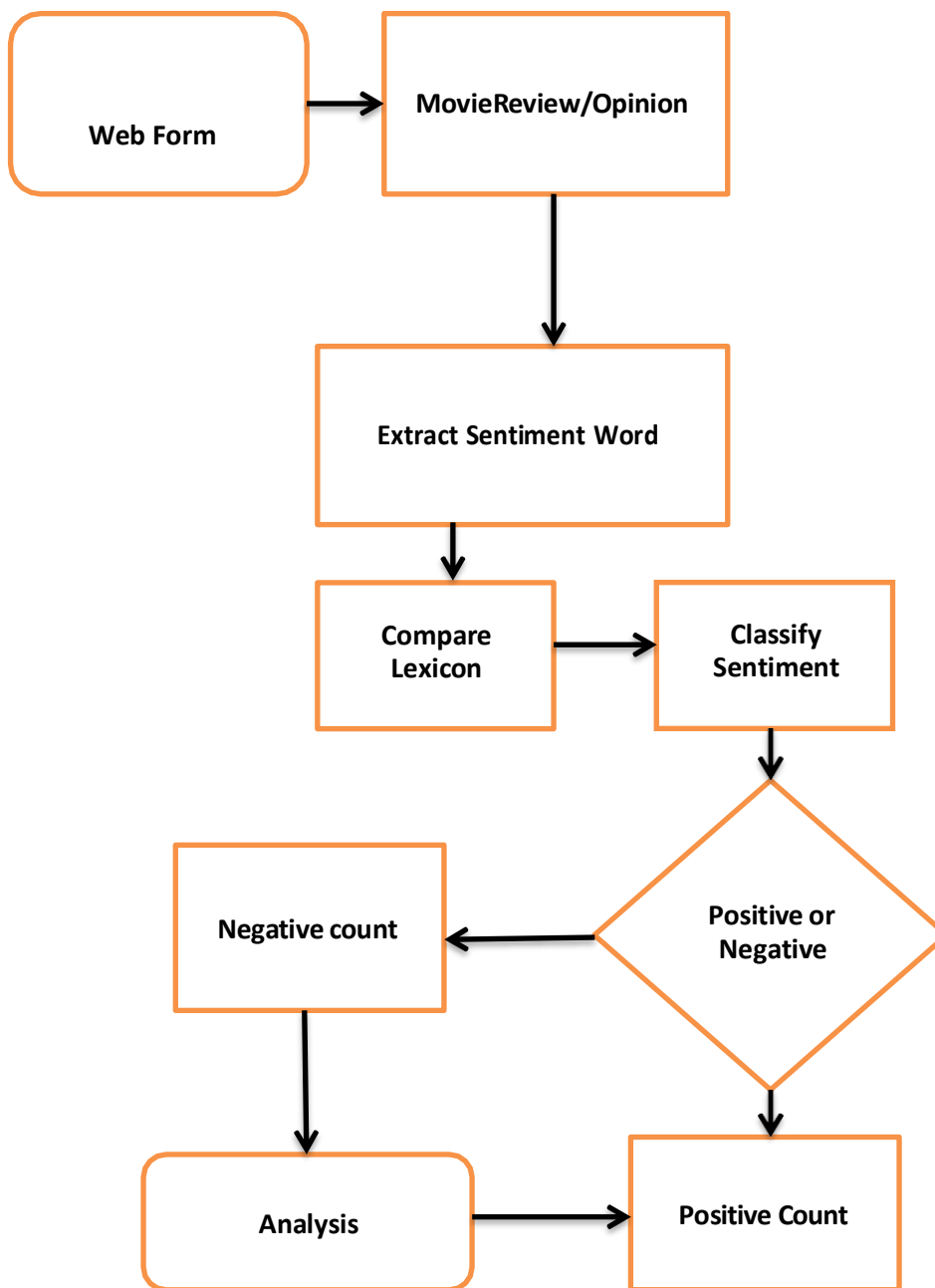


Fig 1.3: Process of determining polarities II

Be that as it may, this assignment is difficult as individuals don't in every case express inside the indistinguishable way. Remarks and surveys contrast from one individual to another as far as their territorial dialects, web slangs, and emojis. Wistful investigation is incredibly worried about ID and grouping of sentiments. It's comprehensively grouped into two kinds initially is information based methodology and other one is utilizing AI strategies. Utilizing first methodology, it requires enormous data set of predefined.

Feelings and an effective information portrayal for perceiving assessments. However, utilizing AI methods, we've train information and test information which is in an incredibly position to be utilized as an information set to foster a classifier. It's very less difficult too. Another undertaking in assessment examination is objectivity recognizable proof where it centers around arranging a given book into one in everything about two classes. Since the subjectivity of words and expressions might rely upon their unique situation and a target report might contain abstract sentences, this issue can now and then be harder than extremity grouping.

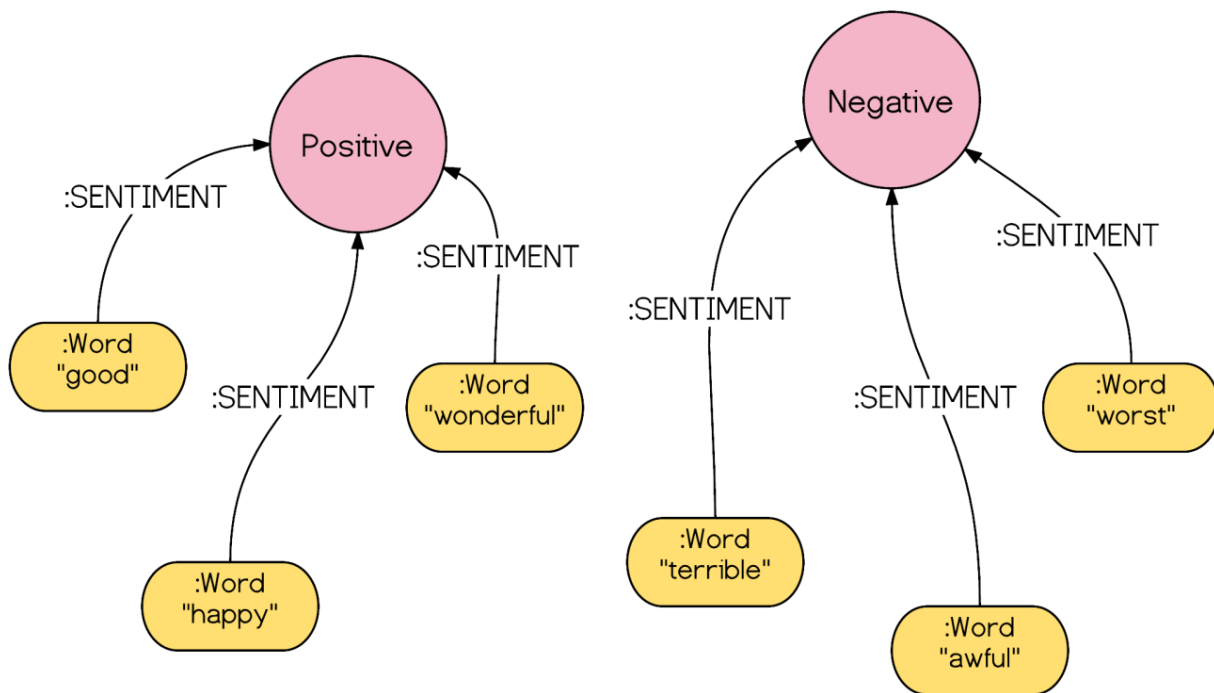


Fig 1.4: Determining polarities

2.SYSTEM ANALYSIS

2.1 EXISTING SYSTEM

Feeling investigation has for some time been an issue for business, promoting and the board regions for more worth acquired inside the choice cycle. Past work about supposition examination has been gaining practical experience in archive level, sentence-level and word-level notion extraction, with both managed and solo methodologies.

2.2 PROPOSED SYSTEM

In this framework IMDB is one in everything about chief mainstream online information bases for motion pictures and characters, a stage where different clients peruse and compose film audits. This gives a larger than usual and different dataset for notion examination. We utilized troupe techniques along with strategic relapse, SVM, choice trees which furnishes better execution with none loss of information. Client gives the predefined credits as information. The model predicts the yield. The yield is introduced to the client inside the past model. The objective variable comprises of genuine qualities so we can't compute precision rather for better execution. So during this proposed model we've made classes for scope of genuine qualities for target variable so presently we will ascertain precision utilizing different AI calculations.

2.3 LITERATURE REVIEW

In general opinion research at the starting of the 20th Century, the science of sentiment analysis and opinion mining has a strong basis. When online product reviews were required and accessible in the middle of 2000, they finally became a major research subject. Just 101 articles on this subject were published in 2005, while almost 5,699 were published in 2015. This means that over a decade sentiment analysis has increased almost 50 times, making it one of the most quickly expanding fields of study in previous years (Mantle et al., 2018) [8]. Throughout the early days of the internet, a person was able to seek feedback from his friends, neighbors and relatives before taking any decision. Opinion sampling, surveys, and general public opinion on its products or services were conducted by organizations. As the World Wide Web has come and particularly with the production and adoration of Web 2.0, where the focus on content generated by users has changed significantly the way the individual expresses his opinion or views. Now people can offer their thoughts, opinions, feelings, blogs, social platforms, forums, and reviews on their own personal web pages. Thanks to rich and

diverse data generated in Web 2.0 applications, the field of opinion mining has advanced quickly (A. Kumar & Sebastian, 2012) [9]. Research into the shifts in the subjects found that social networking such as Twitter and Face book are more focused on the most recent articles from the year 2014 to 2016. In recent year's mobile devices, stock markets, and human emotions were other topics that have become popular (Mantle et al., 2018) [8].

(Swathi & Seshadri, 2017) [10]. Brief information about different types of algorithms used for sentiment analysis is given. Sentimental analysis is defined as the analysis of opinions, thoughts, sentiments, and subjectivity of text are given. Recently introduced algorithms, sentimental analysis techniques are discussed, and also the importance of some fields such as transfer learning, feelings detection, and constructing resources are discussed. The main purpose of this survey is the categorization of recent articles, 54 of the latest published articles which are based on sentiment analysis were categorized and summarized.

(D. Kawade & Oza, 2017) [11]. Social media acts as a vital source where one can interact and may be able to fulfill their demands. This brings both satisfactions for purchasers and also companies. The traditional-based analysis is difficult to research; there are some challenges to beat this problem. Some methods for analyzing feelings, like prediction of user subjects, polarity of feelings scores, analysis and an outsized data processing application, cross-domain classification of feelings, identification of emotional differences, meaning and theme detection, classification of hash tag sentiment rates, sales forecasts, etc. are used. It also addressed briefly the complexities of sentimental analytics to try and do the duty. a number of the challenges like parallel computing for enormous data, sarcasm, grammatically incorrect words, review the author's segmentation, handling noise, and dynamism.

(Yang et al., 2010) [13]. the reviews of sentiments are classified accurately by the algorithms of machine learning like bag-of-words, n-gram, naive Bays classifier, and linguistic communication processing. Then the user's sentiments are categorized as positive, neutral, negative, the highest features of the merchandise will make the customer attract towards that exact product. This work also says that the longer term scope of reviewing the products are supported opinions in several languages, copying drawback of mapping slangs, copying with mocking opinions, so providing comparative opinion between two products for one best and copying with anaphora resolution.

(Gopu & Swarnalatha, 2017)[14]. The sentiment analysis uses natural language processing to naturally classify and derive the emotion from the text and as a result, it has a variety of applications in the consumer sector, for example. Transfer education has also emerged as a new method of machine learning that utilizes existing knowledge to solve issues and to generate forecast results. It also contains the prospect of sentimental analysis such as the appliance of cross-domain shift learning aspects which has not been fully explored, then solving negative problems of text data by using transfer learning becomes very difficult. And they conclude that in future Aspect level sentiment analysis for small texts is considered as the most promising research technique.

(R. Liu et al., 2019) [15]. Machine learning algorithms like Naive Bays, Support vector machine, and Maximum entropy classifier algorithms may be used on the sentimental analysis of giant data. Using these techniques, an enormous volume of knowledge is often utilized to induce optimized and strategically decision-making capability. sentimental analysis is additionally called opinion mining which supplies an excellent and human-like brilliance which analyzes and respond emotions, the user show in social media like Face book, Yammer, Twitter, micro blogs which offer an incredible amount of knowledge on a daily basis in textual or numerical forms and these are classified as structured, semi-structured and non-structured then later they're categorized as positive, negative and customary supported user's attitude towards a specific topic for analysis purpose.

(J. Singh et al., 2017) [17].Sentiment analysis helps to try and do a review of the flicks, product, and customer opinion on products. The role of sentiment analysis in language processing is to get rid of positive or negative polarities from social media messages. Digital social networks are growing increasingly, and culture focused on online media has affected young scientists within their research in the study of opinion. Organizations that's actually keen to see their clients or the general public opinion on their social media goods. Internet services should be ready to test social media data on blogs, web forums, articles, tweets and user feedback.

(Xavier Sumba et al., 2019). We were asked to implement different classification models to predict the sentiment of IMDB reviews, either as positive or negative, using only text each review contains. The goal is to finish the model with the best F1 score and best generalization. We trained different models using multiple combinations of text features and hyper-parameter settings for both the classifiers and also the features, which we found could potentially impact the performance significantly. Every model was evaluated by k- fold cross validation to confirm the consistency of their performance. We found that our greatest performing model was the Naive byes, Support Vector Machines classifier with bag of words.

[Ang (Carl) Li]. (2019). Introduced classification model for sentiment analysis with con- text information participated in the feature space. The dataset was captured from IMDB movie reviews, in which I sampled 1,000 instances from the huge dataset and split it by the 20%-70%-10% ratio for development, cross-validation and final test sets. Through multiple error analysis, including stretchy patterns, character N-grams and elimination of stop words, and tuning procedure on ridge parameter, the performance on final test set hit 84% of percentage correctness and 0.6806 in kappa statistics, revealing marginal improvement to the baseline Logistic Regression model. Some exploration of data and discussion about this, for error analysis and tuning, is also included through the work. Future work of this model.

Li et al. [17] studied the impact of information quality on sentiment classification performance. They considered three criteria, namely informative-mess, readability, and subjectivity, to assess the standard of online product reviews. However, when the dimensions or domain of the information varies, the reliability of the proposed method is questionable compared studies, most papers concentrate on reliability metrics, like overall accuracy or F- score, and omit interval. Additionally, the evaluations of the models are conducted on atiny low number of datasets. This research addresses that gap by means of a comprehensive comparison of sentiment analysis methods within the literature, and an experimental study to judge the performance of deep learning models and related techniques on datasets about different topics. Our research question aims to work out whether it's possible to present outperforming methods for multiple types and sizes of datasets. We repose on previous studies of improvement of SA performance by evaluating the results from the point of view of a mixture of three criteria: overall accuracy, F- score, and time interval. The aim of this comparative study is to offer an objective overview of various techniques which will guide researchers towards the achievementof higher results.

2.4 Module Description

2.4.1 ADMIN MODULE:

In Admin Module, admin collects the Movie reviews for the latest movies from various sources or from the websites like BookMyShow, Pay TM.

2.4.2 SYSTEM MODULE:

In System Module, it loads the dataset, pre-process the data given, apply the algorithms, train the model and produce the result and whether the review is positive or negative.

2.4.3 END-USER MODULE:

In User Module they see the review of the particular movie which they want to see and they can view the reviews of the movie.

2.5 HARDWARE REQUIREMENTS:

- System : Intel i3 Core ad above.
- Hard Disk : 500GB and above
- Monitor : 15' inch LED
- Ram : 8 GB and above.

2.6 SOFTWARE REQUIREMENTS:

- Operating system: Windows 7,8 & 10.
- Coding Language: PYTHON
- Tool: Google CoLab

2.7 FEASIBILITY STUDY

Primer examination of the task possibility; probability the framework are valuable to the association. The target of the achievability study is to check the Technical, Operational and Economical possibility for adding modules and investigating running System. All frameworks are doable in case they're given limitless assets and endless time. There are angles inside the achievability study bit of the fundamental examination.

- Technical Feasibility
- Operation Feasibility
- Economic Feasibility

2.7.1 TECHNICAL FEASIBILITY:

To decide if the proposed framework is in fact achievable, we ought to consistently think about the specialized issues required behind things. Specialized plausibility community on the current framework and how much it can uphold the proposed expansion. Python and its libraries are innovation programming which are acclimated foster Data Analytics. In this way, there's no requirement for additional acquisition of any product and these are open source programming's which are uninhibitedly accessible in Internet.

2.7.2 OPERATIONAL FEASIBILITY:

Proposed projects are advantageous on condition that they will be dressed into data frameworks that might meet the client's working necessities. Functional attainability parts of the venture are to be taken as a pivotal a piece of the apparatus execution. this technique is functional plausible sincethe clients are at home with the advancements and henceforth there's no need to set up the faculty to utilize the framework. Likewise the framework is amazingly well disposed and easy to utilize.

2.7.3 ECONOMIC FEASIBILITY:

To know whether a venture is monetarily practical, we need to accept the different factors as:

- Cost advantage investigation
- Long-term returns
- Maintenance costs

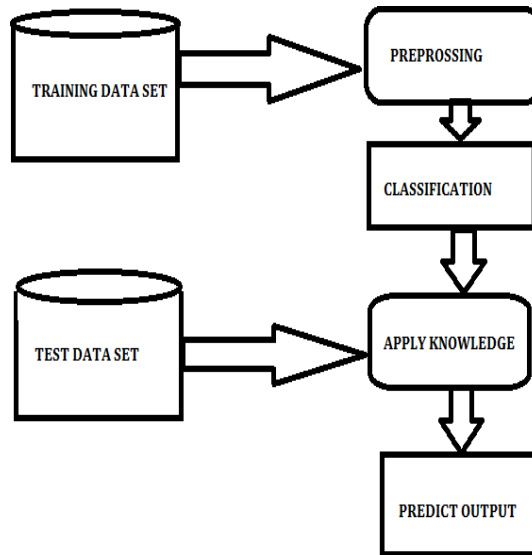
The proposed framework is PC based. It requires normal figuring capacities which is Extremely fundamental prerequisite and can be managed by an association; it doesn't bring about extra financial overheads, which delivers the framework monetarily achievable.

2.8 FUNCTIONAL REQUIREMENTS

- User gives the required movie review as input.
- The model predicts the sentiment of the review.
- The review sentiment output is presented to the user.

3.SYSTEM DESIGN

3.1 SYSTEM ARCHITECTURE



SYSTEM ARCHITECTURE

Fig 3.1: System architecture

TRAINING DATA SET:

Training Data is nothing but enriched or labeled data you wish to coach your models. you would possibly just have to collect more of it to sharpen your model accuracy. But, the probabilities of using your data is pretty low because, as you build a good model you wish great training data at scale.

TEST DATA SET:

The test set might be a bunch of perceptions wont to assess the model exhibition utilizing some preprocessing metric. It's significant that no perceptions from the preparation set are incorporated inside the test set. In the event that the test set contains models from the preparation set, it'll be hard to evaluate whether the calculation has figured out how to sum up from the preparation set or has basically remembered it.

PREPROCESSING:

Information pre-handling is a basic advance in Machine Learning on the grounds that the nature of information and accordingly the valuable data that might be gotten from it straight forwardly influences the adaptability of our model to discover.

PREPROCESSING STEPS:

The following are the preprocessing steps:

- ❖ Movie review Dataset
- ❖ Loading text data
- ❖ Clean text data
- ❖ Develop Vocabulary
- ❖ Save prepared data
- ❖ Apply classification Models

CLASSIFICATION:

Order is that the way toward anticipating the given information focuses. Classes are in some cases called as names. Arrangement prescient demonstrating is that the undertaking of approximating a planning capacity from input factors (X) to discrete yield factors (y).

CLASSIFICATION MODELS:

1. SVM
2. Conv1d

ENSEMBLE METHODS:

The objective of outfit strategies is to predict the few base assessors to work with a given calculation. In order to help generalizability strength more than one assessor.

Ensemble Models:

1. Random forest Classifier
2. Bagging Classifier
3. Gradient Boosting Classifier

3.2 UML Diagrams

In the field of Unified Modeling Language is an even visual particular language for item and examination displaying. UML is a broadly useful of displaying to plan the language that acclimated makes a theoretical model of a framework, referred to as an UML model. The model additionally contains a "Semantic backplane" documentation like composed use cases that drive the model components and outlines.

The significance of UML in Modeling:

A demonstrating language can be utilized as a planning. Whose jargon and rules planned in the theoretical and actual portrayal of a framework. A demonstrating language like UML is consequently a standard language for programming diagrams. The UML is certifiably not a visual programming language; however its models are straightforwardly associated with various programming dialects.

This implies that it's feasible to plan from a model inside the UML to a fake language Java, C++ or Visual Basic, or even to tables in an extremely very PC information base or the persevering store of an item situated data set.

The converse is also conceivable you will be prepared to recreate a model from an execution back to UML. This will be a fake language that is utilized for object-arranged programming improvement. To revise program code all the more proficiently, software engineers regularly make "protests" that are sets of organized information inside programs. UML, which has been normalized by the article Management Group and it was intended for this reason for The language has acquired to help that it's become an average language for imagining and developing programming programs.

An applied model of UML:

The three significant components of UML are

1. The UML's essential structure blocks.
2. The rules that direct how those structure blocks is additionally assembled.
3. Some normal component that applies all through the UML.

3.2.1 CLASS DIAGRAM:

In Unified Modeling Language (UML) is additionally a method of static design outline that depicts the construction of a framework by showing the framework's classes, their properties and activities between the classes.

The class chart is that the most structure block in object-arranged investigation and configuration displaying. It's utilized both for general applied displaying of the semantics of the applying and for detail demonstrating making an interpretation of the model into programming code. The classes in a very class outline address both the first items as well as associations inside the apparatus and furthermore the items to be modified. Inside the classification graph these classes are addressed by boxes which contain three sections:

1. The upper part holds the name of the class.
2. The focus part contains characteristics of a class. Properties of a classification are likewise open, private or ensured.
3. The underside part gives the techniques or activities the classification can take or attempt.

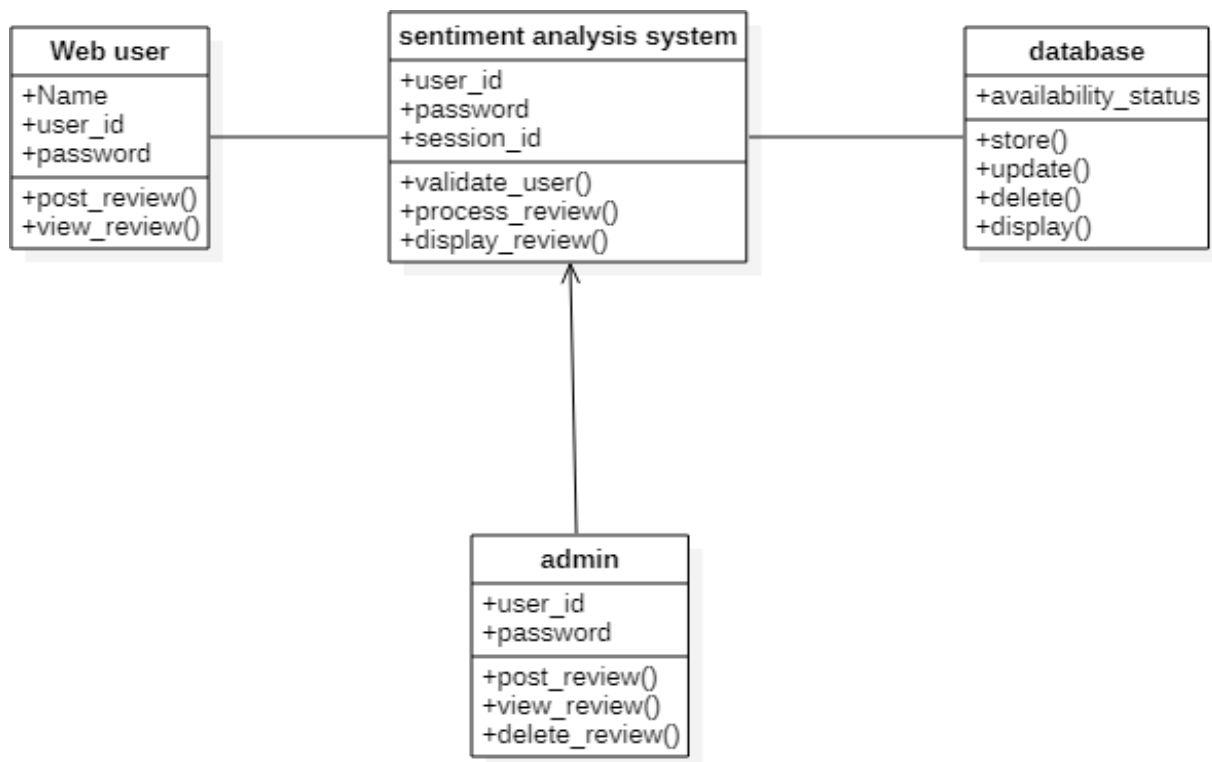


Fig 3.2: Class diagram

3.2.2 USE CASE DIAGRAM:

A utilization case chart in the UML is additionally being a method of conduct outline and it is characterized and made from a Use Case investigation. The motivation behind a utilization case outline is to point what framework capacities are performed by which entertainer. Jobs of the entertainers inside the frameworks are frequently portrayed.

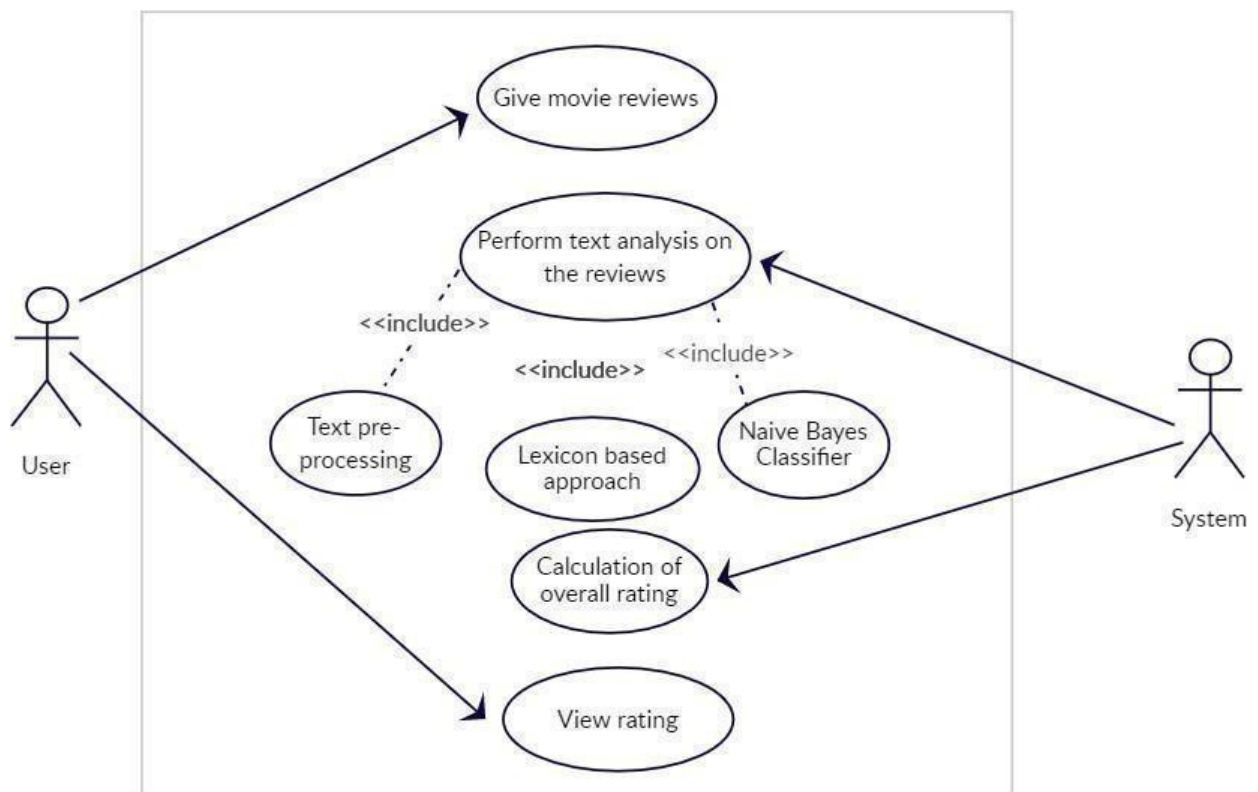


Fig 3.3: Use-case diagram

3.2.3 SEQUENCE DIAGRAM:

The arrangement outline shows the stream for a particular use case or even piece of a chose use case. There are practically plain as day. They show the calls between the different articles in their succession and might show at an exhaustive level and various calls to various items. A Sequence chart has two measurements: The upward measurement shows the arrangement of calls inside the time request that they happen. The flat measurement shows the article occasions to which the messages are sent.

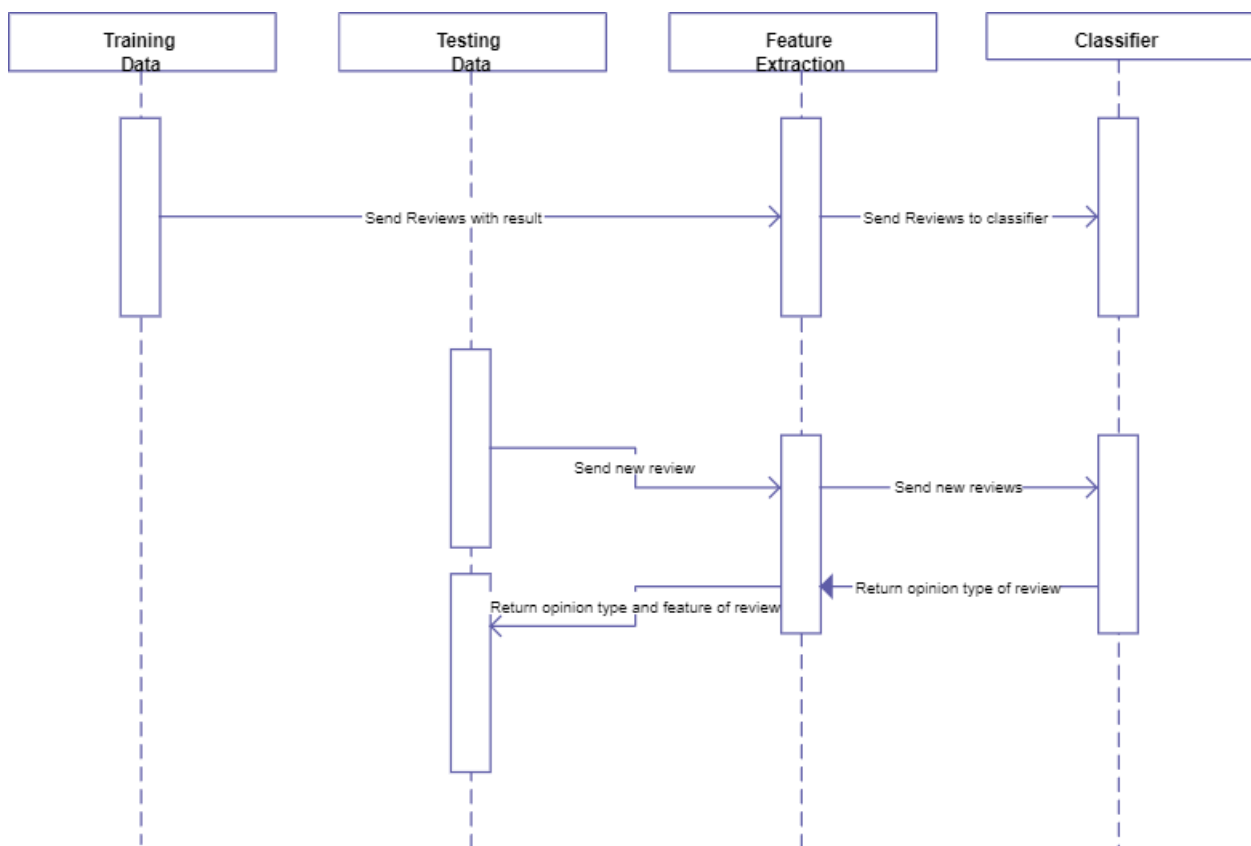


Fig 3.4: Sequence diagram

SENTIMENT ANALYSIS USING MACHINE LEARNING

3.3 Datasets:

review	sentiment
One of the other reviewers has mentioned that after watching just 1 Oz episode you'll be hooked. They are right, as this is exactly what happened with	positive
A wonderful little production. The filming technique is very unassuming- very old-time-BBC fashion and gives a comforting, and sometimes	positive
I thought this was a wonderful way to spend time on a too hot summer weekend, sitting in the air conditioned theater and watching a light-hearted comedy	positive
Basically there's a family where a little boy (Jake) thinks there's a zombie in his closet & his parents are fighting all the time. The movie is sl	negative
Petter Mattei's "Love in the Time of Money" is a visually stunning film to watch. Mr. Mattei offers us a vivid portrait about human relations. This is a mo	positive
Probably my all-time favorite movie, a story of selflessness, sacrifice and dedication to a noble cause, but it's not preachy or boring. It just never gets o	positive
I sure would like to see a resurrection of a up dated Seahunt series with the tech they have today it would bring back the kid excitement in me. I grew u	positive
This show was an amazing, fresh & innovative idea in the 70's when it first aired. The first 7 or 8 years were brilliant, but things dropped off after that. B	negative
Encouraged by the positive comments about this film on here I was looking forward to watching this film. Bad mistake. I've seen 950+ films and this is ti	negative
If you like original gut wrenching laughter you will like this movie. If you are young or old then you will love this movie, hell even my mom liked it. At fir	positive
Phil the Alien is one of those quirky films where the humour is based around the oddness of everything rather than actual punchlines. At fir	negative
I saw this movie when I was about 12 when it came out. I recall the scariest scene was the big bird eating men dangling helplessly from parachutes righ	negative
So im not a big fan of Boll's work but then again not many are. I enjoyed his movie Postal (maybe im the only one). Boll apparently bought the rights to	negative
The cast played Shakespeare. Shakespeare lost. I appreciate that this is trying to bring Shakespeare to the masses, but why ruin s	negative
This a fantastic movie of three prisoners who become famous. One of the actors is george clooney and I'm not a fan but this roll is not bad. Another goc	positive
Kind of drawn in by the erotic scenes, only to realize this was one of the most amateurish and unbelievable bits of film I've ever seen. Sort of like a hig	negative
Some films just simply should not be remade. This is one of them. In and of itself it is not a bad film. But it fails to capture the flavor and the terror of th	positive
This movie made it into one of my top 10 most awful movies. Horrible. There wasn't a continuous minute where there wasn't a fight with or	negative
I remember this film, it was the first film I had watched at the cinema the picture was dark in places I was very nervous it was back in 74/75 my Dad took	positive
An awful film! It must have been up against some real stinkers to be nominated for the Golden Globe. They've taken the story of the first famous fema	negative
After the success of Die Hard and it's sequels it's no surprise really that in the 1990s, a glut of 'Die Hard on a' movies cashed in on the wrong guy, wr	positive
I had the terrible misfortune of having to view this "b-movie" in it's entirety. All I have to say is--- save your time and money!!! This has got	negative
What an absolutely stunning movie, if you have 2.5 hrs to kill, watch it, you won't regret it, it's too much fun! Rajnikanth carries the movie on his shoul	positive
First of all, let's get a few things straight here: a) I AM an anime fan- always has been as a matter of fact (I used to watch Speed Racer all the time in Pre	negative

Fig 3.5: Datasets

ATTRIBUTES DESCRIPTION

- Review: It consists of Movies reviews or Comments on a Movie.
- Sentiment: it consists of review sentiment Whether the review is positive or negative.
- Classes in the dataset are two types. They are
 1. Positive.
 2. Negative.

3.4 Technologies Description

3.4.1 Python

Python is a significant level, deciphered, intuitive and object-situated prearranging programming language. Python is planned to be exceptionally clear. It utilizes English catchphrases every now and again where as different dialects use accentuation, and it's less linguistic developments than different dialects.

Python Identifiers:

A Python identifier might be a name acclimated distinguish a variable, work, class, module or other article. An identifier begins with a letter beginning to end or start to finish or a highlight followed by nothing or more letters, highlights and digits (0 to 9). Python doesn't permit accentuation characters like @, \$, and zippers inside identifiers. Python might be a case delicate programming language. Consequently, Manpower and labor are two distinct identifiers in Python. Here are naming shows for Python identifiers.

- Class names start with a capitalized letter. Any remaining identifiers start with a lowercase letter.
- Starting an identifier with one driving highlight shows that the identifier is private.
- Starting an identifier with two driving highlights shows an unequivocally private identifier.
- If the identifier likewise finishes with two following highlights, the identifier could be a language-characterized uncommon name.

Python Lists:

1. A posting could be an assortment of components. These components could likewise be homogeneous or heterogeneous.
2. A posting might be a worth that contains different qualities in an arranged succession. The term list esteem alludes to the actual rundown (which could be a worth which will be put away in a really factor or passed to a capacity like a few other worth)
3. Even as string esteems are composed with quote characters to stamp where the string starts and closures, a posting starts with a hole section and finishes with an end accentuation, [].
4. Values inside the rundown additionally are called things. Things are isolated with commas (that is, they're comma-delimited).
5. A posting additionally permits copy components.
6. Insertion request is safeguarded in list.
7. List components are isolated by commas and encased inside square sections ([]).

8. Every component inside the rundown has its own one of a kind file.
9. List backings both forward ordering and in reverse ordering, forward record begins from 0 and in reverse file begins from - 1.
10. We access either explicit component by utilizing ordering or set of components by utilizing cutting from the List.
11. We will make list from various perspectives. Like by utilizing list () work, by utilizing square sections "[]" and furthermore by utilizing range () work.
12. List articles are alterable.

Making List by utilizing List ()

1. This list () permits only one string esteem with set of characters.
2. If we give int type information inside the rundown () work then translator will toss 'Type Error' blunder. E.g.:

```
>>> List1=list () #creating void rundown
>>>print (List1) []
>>>type (List1)
```

Python Tuple:

1. Tuple is utilized to address an assortment of homogeneous or heterogeneous components into one element.
2. Tuple items are changeless importance once on the off chance that we make a tuple later we can't alter that Tuple.
3. All components are isolated by commas (,) and encased by brackets. Brackets are discretionary.
4. Tuple permits copy components.
5. Every component inside the Tuple has its own marker Tuple upholds both forward ordering and furthermore in reverse ordering, forward ordering begins from 0 and in reverse ordering begins from - 1.
6. If we take only one component inside the Tuple then we ought to consistently utilize comma (,) right now single component.
7. Tuples might be utilized as keys to the word reference.
8. we can make a Tuple from various perspectives, as with Tuple (), with () or without () too.
9. the most distinction among records and tuples is-Lists are encased in sections ([]) and their components and size might be changed, while tuples are encased in brackets () and can't be

Refreshed.

Making a Tuple with Tuple ():

Eg1:

```
>>>tup=Tuple ([10, 20, 30, True, 'Python'])
>>> print (tup) (10, 20, 30, True, 'Python')
>>> Type (tup)
>>> Id (tup) 52059760
```

Word reference:

Each key's isolated from its worth by a colon (:), the things are isolated by commas, and furthermore the unit is encased in wavy supports. An unfilled word reference with none things is composed with only two wavy supports, similar to this:

Keys are special inside a word reference while qualities probably won't be. The upsides of a word reference are frequently of any sort, yet the keys should be of a permanent information type like strings, numbers, or tuples. Getting to Values in Dictionary to get to word reference components, you'll utilize the comfortable square sections along with the way to get its worth. Following could be a straightforward model –

```
Dict =print "dict ['Name']: “dict ['Name'] Print "dict ['Age']: “dict ['Age']
```

At the point when the above code is executed, it creates the ensuing outcome dict ['Name']: Zara dict ['Age']: 7

1. Python is deciphered: Python is handled at runtime by the mediator. you are doing not should arrange your program prior to executing it. This can be equivalent to PERL and PHP.
2. Python is Interactive: You can really sit at a Python incite and interface with the mediator on to compose your projects. Python accompanies an intelligent mediator. After you type python in your shell or electronic correspondence, the python translator open up with a >>> brief and anticipating your guidelines.

```
>>> says that your are inside the python translator
```

```
$ Python
```

SENTIMENT ANALYSIS USING MACHINE LEARNING

Python 2.7.6 (default, Apr 24 2015, 09:38:35) [GCC 4.2.1 Compatible Apple LLVM 6.0 (crash 600.0.39)] on Darwin Type "help", "copyright", "credits" or "permit" for more data.

>>> If you might want to leave the Python console anytime, simply type exit () or utilize the easy route Ctrl + Z for Windows and Ctrl + D for Mac/Linux. Then, at that point you won't see >>> any more.

Python is Object-Oriented:

Python upholds Object-Oriented style or procedure of programming that embodies code inside objects.

Factors - In Python there are no presentations.

Progressively Typed:

Python is a dynamic-composed language. Numerous different dialects are static composed, like C/C++ and Java.

A static composed language requires the software engineer to expressly mention to the PC what sort of "thing" every information esteem is. For instance, in C on the off chance that you had a variable that was to contain the cost of something; you would need to proclaim the variable as a "glide" type. This tells the compiler that the solitary information that can be utilized for that variable should be a gliding point number,

For example a number with a decimal point. On the off chance that some other information esteem was allocated to that factor, the compiler would give a mistake when attempting to accumulate the program.

Note: Everything is an item in Python.

Python is a Beginner's Language:

Python is an incredible language for the novice software engineers and supports the advancement of a wide scope of uses from simple text preparing.

Running Python Scripts:

Open your content manager, type the accompanying content and save it as "hello.py". Print "Hi, World!"

What's more, shown this program to calling "python hello.py". Ensure you change to the directory where you saved the document prior to doing it.

```
C:\Users\USER\Desktop> python hello.py
```

```
Hello, World
```

3.4.1 ENSEMBLE-METHODS

The objective of troupe strategies is to join the expectations of a few base assessors worked with a provided learning calculation to further develop generalizability strength over a solitary assessor.

Two groups of troupe strategies are generally recognized:

- In meeting techniques, imdb_rating advertisement conv1d models are summed up utilizing the greatest length 150 cushioning and 3 boundaries. In this meeting technique three layers are announced to quantify the real boundary esteems
- By contrast, in assess technique is utilized to know the assessment measurements of the survey. Two models are referenced to assess the model and for estimating the boundaries precision.

Table 3.9: Differences and Similarities between Bagging and Boosting

S.NO	BAGGING	BOOSTING
1.	Simplest way of combining predictions that belong to the same type.	Away of combining predictions that belong to the different types.
2.	Aim to decrease variance, not bias.	Aim to decrease bias, not variance.
3.	Each model receives equal weight.	Models are weighted according to their performance.
4.	Each model is built independently.	New models are influenced by performance of previously built models.
5.	Different training data subsets are randomly drawn with replacement from the entire training dataset.	Every new subset contains the elements that were misclassified by previous models.
6.	Bagging tries to solve over-fitting problem.	Boosting tries to reduce bias.
7.	If the classifier is unstable (high variance), then apply bagging.	If the classifier is stable and simple (high bias) then apply boosting.
8.	Random forest.	Gradient boosting.

4. CODING AND IMPLEMENTATION

The objective of the coding or programming stage is to interpret the arranging of the framework created during the look stage into code during a given programming language, which might be executed by a PC which plays out the calculation determined by the arranging. The coding stage influences both testing and support. The objective of coding isn't to reduce back the execution expense yet the objective ought to be proportional back the cost of later stages. At the end of the day, the objective isn't to work on the assignment of a software engineer. Maybe the objective ought to be to improve on the obligation of the analyzer and maintainer.

Coding Approach:

There are two significant methodologies for coding any bundle. They're Top-Down approach and Bottom-up approach. Granular perspectives are regularly best appropriate for fostering the article situated frameworks. During framework configuration stage to downsize the intricacy, we disintegrate the framework into a proper number of subsystems, that articles will be displayed freely. These articles display the manner in which the subsystems play out their tasks. When the articles are demonstrated they're carried out through coding. While related with the indistinguishable framework on the grounds that the articles are autonomous of each other. The Bottom-up approach is more reasonable for coding these articles.

In this methodology, we initially do the coding of articles freely so incorporate these modules into one framework to which they have a place.

```
from google.colab import drive
drive.mount('/content/drive')

from keras.datasets import imdb
import pandas as pd
import numpy as np
from keras.layers import LSTM, Activation, Dropout, Dense, Input, Conv1D,
MaxPooling1D, GlobalMaxPooling1D
from keras.layers.embeddings import Embedding
from keras.models import Model
import string
import re
from keras.preprocessing.text import Tokenizer
from sklearn.preprocessing import LabelBinarizer
from keras.preprocessing.sequence import pad_sequences
import keras
from sklearn.model_selection import train_test_split
```


SENTIMENT ANALYSIS USING MACHINE LEARNING

```
data = pd.read_csv('/content/drive/MyDrive/fp/sa/movie_reviews.csv')
data['review'] = data['review'].str.lower()

stopwords = [ "a", "about", "above", "after", "again", "against", "all",
              "am", "an", "and", "any", "are", "as", "at", "be", "because",
              "been", "before", "being", "below", "between", "both",
              "but", "by", "could", "did", "do", "does", "doing", "down", "during",
              "each", "few", "for", "from", "further", "had", "has",
              "have", "having", "he", "he'd", "he'll", "he's", "her", "here",
              "here's", "hers", "herself", "him", "himself", "his",
              "how", "how's", "i", "i'd", "i'll", "i'm", "i've", "if", "in", "into",
              "is", "it", "it's", "its", "itself", "let's", "me", "
more", "most", "my", "myself", "nor", "of", "on", "once", "only", "or",
              "other", "ought", "our", "ours", "ourselves", "out", "over",
              "own", "same", "she", "she'd", "she'll", "she's", "should"
              ,
              "so", "some", "such", "than", "that", "that's", "the"
              , "their", "theirs", "them", "themselves", "then", "there", "there's",
              "these", "they", "they'd", "they'll", "they're", "the
y've", "this", "those", "through", "to", "too", "under", "until", "up",
              "very", "was", "we", "we'd", "we'll", "we're", "we've"
              , "were", "what", "what's", "when", "when's", "where", "where's",
              "which", "while", "who", "who's", "whom", "why", "why
's", "with", "would", "you", "you'd", "you'll", "you're", "you've"
              ,
              "your", "yours", "yourself", "yourselves" ]

def remove_stopwords(data):
    data['review without stopwords'] = data['review'].apply(lambda x : ' '.join([word for word in x.split() if word not in (stopwords)]))
    return data

def remove_tags(string):
    result = re.sub('<.*?>', '', string)
    return result

data_without_stopwords = remove_stopwords(data)
data_without_stopwords['clean_review'] = data_without_stopwords['review without stopwords'].apply(lambda cw : remove_tags(cw))
data_without_stopwords['clean_review'] = data_without_stopwords['clean_review'].str.replace('{}'.format(string.punctuation), ' ')
```

SENTIMENT ANALYSIS USING MACHINE LEARNING

```
data_without_stopwords.head()

reviews = data_without_stopwords['clean_review']
reviews

reviews_list = []
for i in range(len(reviews)):
    reviews_list.append(reviews[i])

sentiment = data_without_stopwords['sentiment']

y = np.array(list(map(lambda x: 1 if x=="positive" else 0, sentiment)))
y

X_train, X_test, Y_train, Y_test = train_test_split(reviews_list, y, test_size=0.2, random_state = 45)

len(Y_train)

tokenizer = Tokenizer(num_words=5000)
tokenizer.fit_on_texts(X_train)

words_to_index = tokenizer.word_index

len(words_to_index)

def read_glove_vector(glove_vec):
    with open(glove_vec, 'r', encoding='UTF-8') as f:
        words = set()
        word_to_vec_map = {}
        for line in f:
            w_line = line.split()
            curr_word = w_line[0]
            word_to_vec_map[curr_word] = np.array(w_line[1:], dtype=np.float64)
        return word_to_vec_map

word_to_vec_map = read_glove_vector('/content/drive/My Drive/fp/sa/we/glove.6B.50d.txt')

maxLen = 150
```

SENTIMENT ANALYSIS USING MACHINE LEARNING

```
vocab_len = len(words_to_index)
embed_vector_len = word_to_vec_map['moon'].shape[0]

emb_matrix = np.zeros((vocab_len, embed_vector_len))

for word, index in words_to_index.items():
    embedding_vector = word_to_vec_map.get(word)
    if embedding_vector is not None:
        emb_matrix[index, :] = embedding_vector

embedding_layer = Embedding(input_dim=vocab_len, output_dim=embed_vector_len,
                             input_length=maxLen, weights = [emb_matrix], trainable=False)

def imdb_rating(input_shape):
    X_indices = Input(input_shape)
    embeddings = embedding_layer(X_indices)
    X = LSTM(128, return_sequences=True)(embeddings)
    X = Dropout(0.6)(X)
    X = LSTM(128, return_sequences=True)(X)
    X = Dropout(0.6)(X)
    X = LSTM(128)(X)
    X = Dense(1, activation='sigmoid')(X)
    model = Model(inputs=X_indices, outputs=X)
    return model

def conv1d_model(input_shape):
    X_indices = Input(input_shape)
    embeddings = embedding_layer(X_indices)
    X = Conv1D(512, 3, activation='relu')(embeddings)
    X = MaxPooling1D(3)(X)
    X = Conv1D(256, 3, activation='relu')(X)
    X = MaxPooling1D(3)(X)
    X = Conv1D(256, 3, activation='relu')(X)
    X = Dropout(0.8)(X)
    X = MaxPooling1D(3)(X)
    X = GlobalMaxPooling1D()(X)
    X = Dense(256, activation='relu')(X)
    X = Dense(1, activation='sigmoid')(X)
    model = Model(inputs=X_indices, outputs=X)
    return model

model = imdb_rating((maxLen,))
model.summary()

model_1d = conv1d_model((maxLen,))
```

SENTIMENT ANALYSIS USING MACHINE LEARNING

```
model_1d.summary()

X_train_indices = tokenizer.texts_to_sequences(X_train)

X_train_indices = pad_sequences(X_train_indices, maxlen=maxLen, padding='post')
X_train_indices.shape

adam = keras.optimizers.Adam(learning_rate = 0.0001)
model_1d.compile(optimizer=adam, loss='binary_crossentropy', metrics=['accuracy'])

model_1d.fit(X_train_indices, Y_train, batch_size=64, epochs=15)

adam = keras.optimizers.Adam(learning_rate = 0.0001)
model.compile(optimizer=adam, loss='binary_crossentropy', metrics=['accuracy'])

model.fit(X_train_indices, Y_train, batch_size=64, epochs=1)

X_test_indices = tokenizer.texts_to_sequences(X_test)
X_test_indices = pad_sequences(X_test_indices, maxlen=maxLen, padding='post')

model.evaluate(X_test_indices, Y_test)

model_1d.evaluate(X_test_indices, Y_test)

preds = model_1d.predict(X_test_indices)

n = np.random.randint(0,9999)
X_test[n]
if preds[n] > 0.5:
    print('predicted sentiment : positive')
else:
    print('predicted sentiment: negative')
if (Y_test[n] == 1):
    print('correct sentiment: positive')
else:
    print('correct sentiment : negative')
```

SENTIMENT ANALYSIS USING MACHINE LEARNING

```
preds[n]

Y_test[n]

model_1d.save_weights('/content/drive/My Drive/fp/sa/weights/imdb-weights.tsv')

reviews_list_idx = tokenizer.texts_to_sequences(reviews_list)

def add_score_predictions(data, reviews_list_idx):
    data['sentiment score'] = 0
    reviews_list_idx = pad_sequences(reviews_list_idx, maxlen=maxLen, padding='post')
    review_preds = model.predict(reviews_list_idx)
    data['sentiment score'] = review_preds
    pred_sentiment = np.array(list(map(lambda x : 'positive' if x > 0.5 else 'negative', review_preds)))
    data['predicted sentiment'] = 0
    data['predicted sentiment'] = pred_sentiment
    return data

data = add_score_predictions(data,

reviews_list_idx)data
```

5. OUTPUT SCREENS

```

[9] 2 data_without_stopwords['clean_review'] = data_without_stopwords['review without stopwords'].apply(lambda cw: remove_tags(cw))
    3 data_without_stopwords['clean_review'] = data_without_stopwords['clean_review'].str.replace('{}'.format(string.punctuation), ' ')

[10] 1 data_without_stopwords.head()

```

	review	sentiment	review without stopwords	clean_review
0	one of the other reviewers has mentioned that...	positive	one reviewers mentioned watching just 1 oz epi...	one reviewers mentioned watching just 1 oz epi...
1	a wonderful little production. the...	positive	wonderful little production. the f...	wonderful little production the filming techn...
2	i thought this was a wonderful way to spend ti...	positive	thought wonderful way spend time hot summer we...	thought wonderful way spend time hot summer we...
3	basically there's a family where a little boy ...	negative	basically family little boy (jake) thinks zomb...	basically family little boy jake thinks zomb...
4	petter matter's "love in the time of money" is...	positive	petter matter's "love time money" visually stu...	petter mattei s love time money visually stu...

```

[11] 1
    2 reviews = data_without_stopwords['clean_review']
    3 reviews

0      one reviewers mentioned watching just 1 oz epi...
1  wonderful little production the filming techn...
2  thought wonderful way spend time hot summer we...
3  basically family little boy jake thinks zomb...
4  petter mattei s love time money visually stu...
...
49995 thought movie right good job wasn t creative ...
49996 bad plot bad dialogue bad acting idiotic di...
49997 catholic taught parochial elementary schools n...
49998 going disagree previous comment side malin on...
49999 no one expects star trek movies high art fans...
Name: clean_review, Length: 50000, dtype: object

```

Fig 5.1: Loading the data And Removing Stop words

```

[12] 1 reviews_list = []
    2 for i in range(len(reviews)):
    3     reviews_list.append(reviews[i])
    4
    5

[13] 1 sentiment = data_without_stopwords['sentiment']

[14] 1 y = np.array(list(map(lambda x: 1 if x=="positive" else 0, sentiment)))

[15] 1 y
    array([1, 1, 1, ..., 0, 0, 0])

[16] 1 X_train, X_test, Y_train, Y_test = train_test_split(reviews_list, y, test_size=0.2, random_state = 45)

[17] 1 len(Y_train)
    40000

[18] 1 tokenizer = Tokenizer(num_words=5000)
    2 tokenizer.fit_on_texts(X_train)

[19] 1 words_to_index = tokenizer.word_index

[20] 1 len(words_to_index)

```

Fig 5.2: Training the Data

SENTIMENT ANALYSIS USING MACHINE LEARNING

```
[27] 1 model = imdb_rating((maxlen,))
      2 model.summary()
```

Model: "model"

Layer (type)	Output Shape	Param #
input_1 (InputLayer)	[(None, 150)]	0
embedding (Embedding)	(None, 150, 50)	4770100
lstm (LSTM)	(None, 150, 128)	91648
dropout (Dropout)	(None, 150, 128)	0
lstm_1 (LSTM)	(None, 150, 128)	131584
dropout_1 (Dropout)	(None, 150, 128)	0
lstm_2 (LSTM)	(None, 128)	131584
dense (Dense)	(None, 1)	129

Total params: 5,125,045
Trainable params: 354,945
Non-trainable params: 4,770,100

```
[28] 1 model_id = conv1d_model((maxlen,))
      2 model_id.summary()
```

Model: "model_1"

Layer (type)	Output Shape	Param #
--------------	--------------	---------

Activate Windows
Go to Settings to activate Windows.

Fig 5.3: Summary of IMDB dataset

```
[28] 1 model_id.summary()
```

Model: "model_1"

Layer (type)	Output Shape	Param #
input_2 (InputLayer)	[(None, 150)]	0
embedding (Embedding)	(None, 150, 50)	4770100
conv1d (Conv1D)	(None, 148, 512)	77312
max_pooling1d (MaxPooling1D)	(None, 49, 512)	0
conv1d_1 (Conv1D)	(None, 47, 256)	393472
max_pooling1d_1 (MaxPooling1D)	(None, 15, 256)	0
conv1d_2 (Conv1D)	(None, 13, 256)	196864
dropout_2 (Dropout)	(None, 13, 256)	0
max_pooling1d_2 (MaxPooling1D)	(None, 4, 256)	0
global_max_pooling1d (GlobalMaxPooling1D)	(None, 256)	0
dense_1 (Dense)	(None, 256)	65792
dense_2 (Dense)	(None, 1)	257

Total params: 5,503,797
Trainable params: 733,697
Non-trainable params: 4,770,100

```
[29] 1 X_train_indices = tokenizer.texts_to_sequences(X_train)
```

Activate Windows
Go to Settings to activate Windows.

Fig 5.4: Conv1dModel Accuracy

SENTIMENT ANALYSIS USING MACHINE LEARNING

```
[30] 1 X_train_indices = pad_sequences(X_train_indices, maxlen=maxLen, padding='post')
     2 X_train_indices.shape
(40000, 150)

[31] 1 adam = keras.optimizers.Adam(learning_rate = 0.0001)
     2 model.compile(optimizer=adam, loss='binary_crossentropy', metrics=['accuracy'])

[32] 1 model.fit(X_train_indices, Y_train, batch_size=64, epochs=15)
Epoch 1/15
625/625 [=====] - 176s 258ms/step - loss: 0.6219 - accuracy: 0.6460
Epoch 2/15
625/625 [=====] - 161s 258ms/step - loss: 0.4705 - accuracy: 0.7782
Epoch 3/15
625/625 [=====] - 161s 258ms/step - loss: 0.4209 - accuracy: 0.8069
Epoch 4/15
625/625 [=====] - 161s 257ms/step - loss: 0.3835 - accuracy: 0.8291
Epoch 5/15
625/625 [=====] - 161s 257ms/step - loss: 0.3547 - accuracy: 0.8449
Epoch 6/15
625/625 [=====] - 161s 258ms/step - loss: 0.3395 - accuracy: 0.8504
Epoch 7/15
625/625 [=====] - 160s 257ms/step - loss: 0.3068 - accuracy: 0.8673
Epoch 8/15
625/625 [=====] - 160s 256ms/step - loss: 0.2803 - accuracy: 0.8826
Epoch 9/15
625/625 [=====] - 159s 254ms/step - loss: 0.2582 - accuracy: 0.8927
Epoch 10/15
625/625 [=====] - 159s 255ms/step - loss: 0.2408 - accuracy: 0.9001
Epoch 11/15
625/625 [=====] - 160s 256ms/step - loss: 0.2055 - accuracy: 0.9199
Epoch 12/15
625/625 [=====] - 160s 255ms/step - loss: 0.1760 - accuracy: 0.9298
Epoch 13/15
625/625 [=====] - 160s 255ms/step - loss: 0.1534 - accuracy: 0.9425
Epoch 14/15
625/625 [=====] - 159s 255ms/step - loss: 0.1105 - accuracy: 0.9614
Epoch 15/15
625/625 [=====] - 159s 255ms/step - loss: 0.0856 - accuracy: 0.9701
<keras.callbacks.History at 0x7f14139d8e50>

[33] 1 adam = keras.optimizers.Adam(learning_rate = 0.0001)
     2 model.compile(optimizer=adam, loss='binary_crossentropy', metrics=['accuracy'])

[36] 1 model.fit(X_train_indices, Y_train, batch_size=64, epochs=1)
625/625 [=====] - 460s 736ms/step - loss: 0.4276 - accuracy: 0.8050
<keras.callbacks.History at 0x7f140f980410>

[37] 1 X_test_indices = tokenizer.texts_to_sequences(X_test)
     2
     3 X_test_indices = pad_sequences(X_test_indices, maxlen=maxLen, padding='post')

[38] 1 model.evaluate(X_test_indices, Y_test)
313/313 [=====] - 32s 98ms/step - loss: 0.4091 - accuracy: 0.8149
[0.40914812684059143, 0.8148999810218811]

[39] 1 model.evaluate(X_test_indices, Y_test)
313/313 [=====] - 9s 28ms/step - loss: 0.3557 - accuracy: 0.8433
[0.35571014881134033, 0.8432999849319458]

[40] 1 preds = model.predict(X_test_indices)
```

Fig 5.5: Evaluating the data

SENTIMENT ANALYSIS USING MACHINE LEARNING

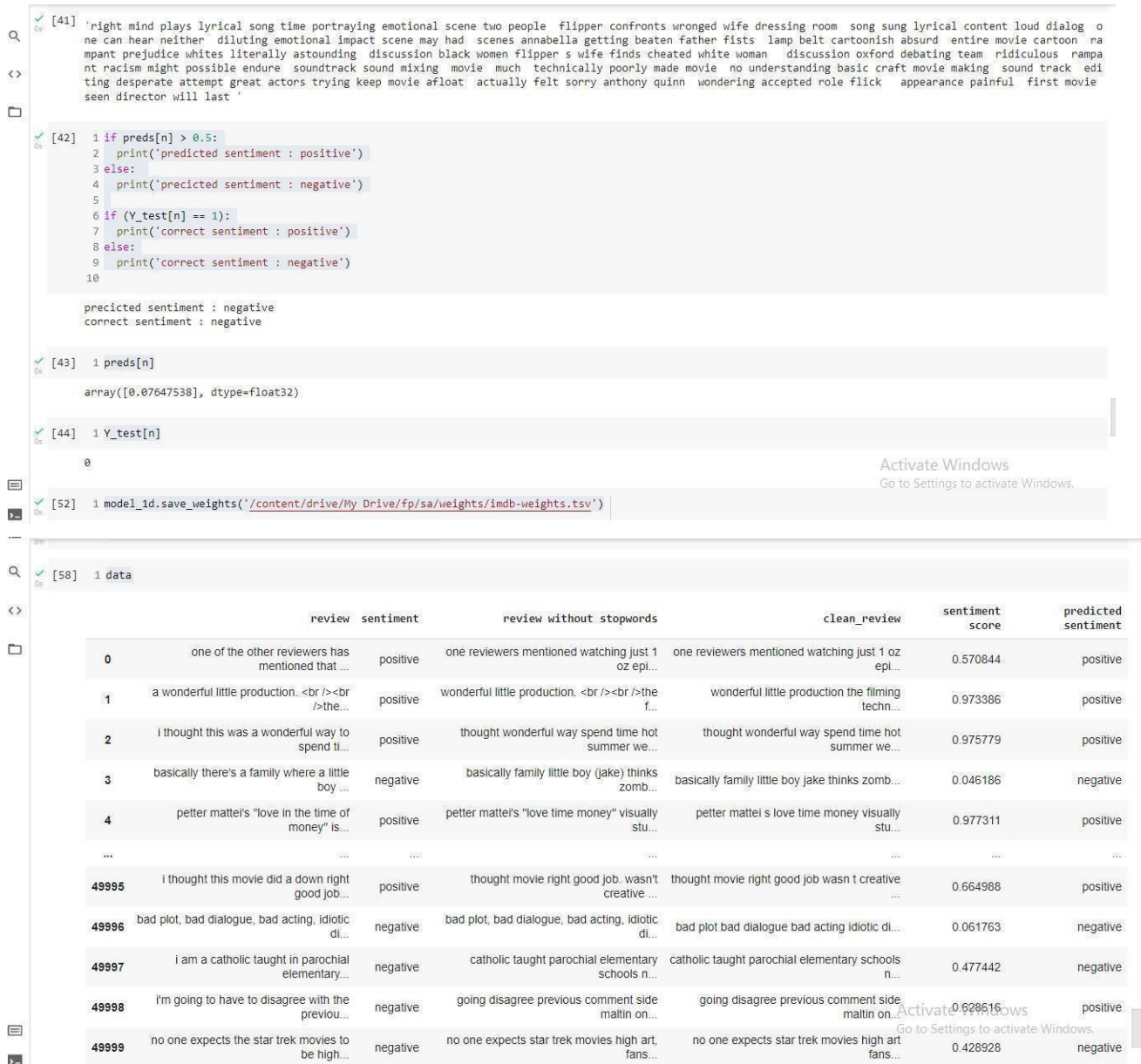


Fig 5.6: Final Accuracy of each Review

6.EVALUATION METRICS

Assessing your AI calculation is a fundamental piece of any task. Your model might give you fulfilling results when assessed utilizing a measurement say exactness score however may give helpless outcomes when considered in contrast to different measurements like logarithmic misfortune or some other such measurement. The vast majority of the occasions we use arrangement exactness to gauge the presentation of our model, anyway it isn't sufficient to genuinely pass judgment on our model. In this post, we will cover various kinds of assessment measurements accessible.

- Classification Accuracy
- Logarithmic Loss
- Confusion Matrix
- F1 Score
- Area under Curve
- Mean Squared Error

Classification Accuracy:

Grouping Accuracy is the thing that we typically mean, when we utilize the term precision. it's the proportion of number of right expectations to the full number of information tests.

$$Accuracy = \frac{\text{Number of Correct predictions}}{\text{Total number of predictions made}}$$

It functions admirably gave that there is equivalent number of tests having a place with each class.

For instance, consider that there are 98% examples of refinement an and a couple of tests of complexity B in our preparation set. Then, at that point our model can without much of a stretch get 98% preparing exactness by basically foreseeing each preparation test having a place with class A.

At the point when the indistinguishable model is tried on a test set with 60% examples of refinement An and 40% examples of complexity B, then, at that point the test exactness would drop to 60%.Classification Accuracy is incredible, however gives us the misguided feeling of accomplishing high precision.

The genuine issue emerges, when the cost of misclassification of the minor class tests are extremely high. On the off chance that we handle an uncommon however deadly illness, the cost of neglecting to analyze the infection of a victim is much more than the cost of sending a sound individual to more tests.

Logarithmic Loss:

Logarithmic Loss or Log Loss works by punishing the bogus groupings. It functions admirably for multi-class grouping. When working with Log Loss, the classifier should appoint likelihood to each class for every one of the examples. Assume, there are N tests having a place with M classes, then, at that point the Log Loss is determined as underneath:

$$\text{Logarithmic Loss} = \frac{-1}{N} \sum_{i=1}^N \sum_{j=1}^M y_{ij} * \log(p_{ij})$$

Where y_{ij} , shows whether test I has a place with class j or not p_{ij} , demonstrates the likelihood of test I having a place with class j Log Loss has no upper bound and it exists on the reach $[0, \infty)$. Log Loss closer to 0 demonstrates higher exactness, though assuming the Log Loss is away from 0, it shows lower precision. As a rule, limiting Log Loss gives more noteworthy precision for the classifier.

Disarray Matrix:

Disarray Matrix as the name proposes gives us a network as yield and depicts the total presentation of the model.

How about we accept we have a paired arrangement issue. We have a few examples having a place with two classes: YES or NO. Additionally, we have our own classifier which predicts a class for a given information test. On testing our model on 165 examples, we get the accompanying outcome.

n=165	Predicted: NO	Predicted: YES
Actual: NO	50	10
Actual: YES	5	100

Disarray Matrix

There are 4 significant terms:

- True Positives: The cases in which we anticipated YES and the genuine yield was additionally YES.
- True Negatives: The cases in which we anticipated NO and the actual output was NO.

SENTIMENT ANALYSIS USING MACHINE LEARNING

- False Positives: The cases in which we anticipated YES and the actual output was NO.
- False Negatives: The cases in which we anticipated NO and the genuine yield was YES.

Exactness for the lattice can be determined by taking normal of the qualities lying across the "principle corner to corner" for example

Confusion Matrix forms the basis for the other types of metrics.

F1 Score:

F1 Score is utilized to live a test's exactness

F1 Score is that the mean among accuracy and review. The reach for F1 Score is [0, 1]. It reveals to you the way exact your classifier is (the number of cases it characterizes accurately), still as how powerful (it doesn't miss countless cases).

High exactness however lower review, gives you a truly precise, yet it then, at that point misses a curiously large number of examples that are hard to arrange. The more noteworthy the F1 Score, the upper is that the exhibition of our model. Numerically, it are normal communicated as :

F1 Score attempts to glance out the harmony among accuracy and review.

- Precision: it is the measure of right certain outcomes partitioned by the quantity of positive outcomes anticipated by the classifier.
- Recall: it is the amount of right sure outcomes separated by the quantity of every single pertinent example (all examples that ought to are distinguished as sure).

$$Precision = \frac{TruePositives}{TruePositives + FalseNegatives}$$

Region under Curve Area under Curve (AUC) is one through and through the first broadly utilized measurements for assessment. Its utilized for parallel order issue. AUC of a classifier is up to the likelihood that the classifier will rank a haphazardly picked positive model over an arbitrarily picked negative model. Prior to characterizing AUC, permit us to get a handle on two essential terms:

- **True Positive Rate (Sensitivity):** True Positive Rate is characterized as $TP/(FN+TP)$. Genuine Positive Rate compares to the extent of positive information focuses that are accurately considered as sure, as to all sure information focuses.
- **False Positive Rate (Specificity):** False Positive Rate is characterized as $FP/(FP+TN)$. Bogus Positive Rate compares to the extent of negative information focuses that are erroneously considered as sure, with pertinence all adverse information focuses. Bogus Positive Rate and True Positive Rate both include values inside the reach $[0, 1]$. FPR and TPR both have processed at limit esteems like $(0.00, 0.02, 0.04 \dots 1.00)$ and a chart is drawn. AUC is that the domain under the bend of plot False Positive Rate versus Genuine Positive Rate at various focuses in $[0, 1]$.

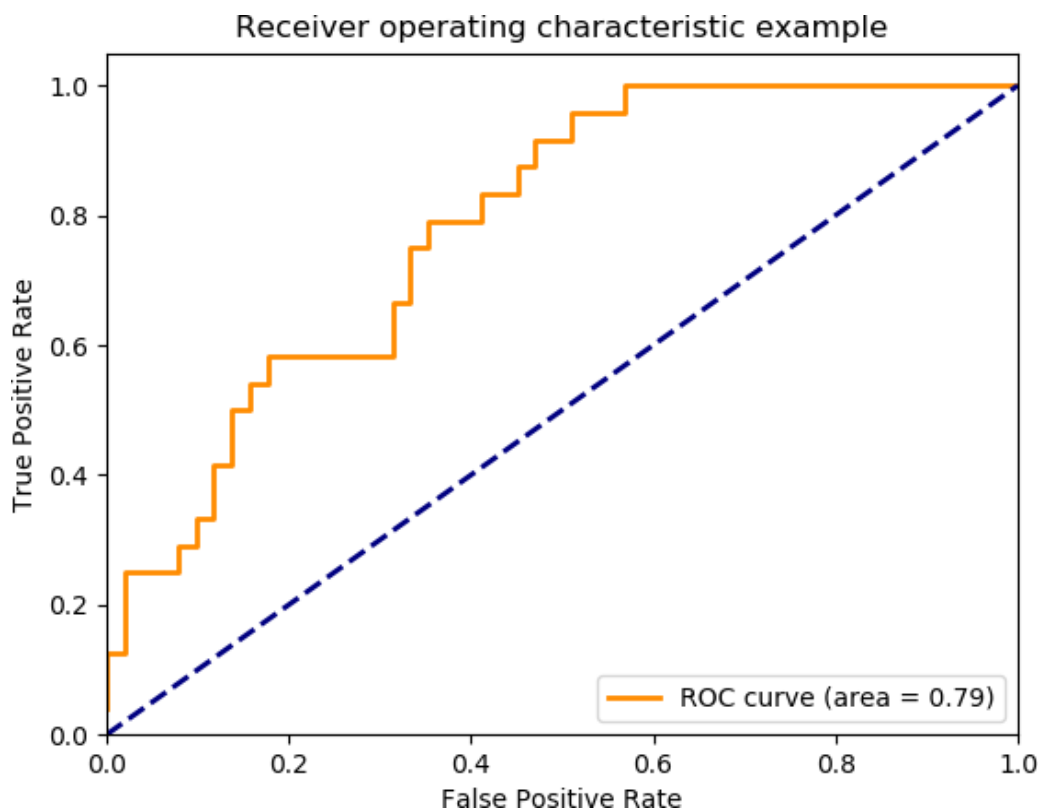


Fig 6.1: Positive rate curve

As obvious, AUC incorporates an assortment of [0, 1]. The more noteworthy the worth, the upper is that the presentation of our model.

Mean Squared Error:

Mean Squared Error (MSE) is very similar to Mean Absolute Error, the lone genuine contrast being that MSE takes the normal of the square of the distinction between the essential qualities and in this way the expected qualities. The benefit of MSE being that it's simpler to figure the slope, though Mean Absolute Error requires muddled maths apparatuses to register the angle. As, we take square of the blunder, the impact of bigger mistakes become more articulated then more modest blunder, subsequently the model would now be able to zero in additional on the bigger blunders.

$$MeanSquaredError = \frac{1}{N} \sum_{j=1}^N (y_j - \hat{y}_j)^2$$

What to undertake and do with missing values

Now and again informational collections will have missing perceptions, which makes ascertaining similitude testing. You have a few choices for filling in these missing information focuses:

- Fill inside the clear regions with zeros,
- Replace the missing qualities with the middle for the set,
- Replace the missing qualities with the mean for the set,
- Use a k-closest neighbor or EM calculation.
- Replace the missing qualities with the chief continuous happened an incentive for that property.

7. TESTING

TESTING TECHNOLOGIES

Testing is that the cycle of recognition blunders. Testing plays out a top quality job for affirmation and for guaranteeing the adaptability of programming. The aftereffects of testing are utilized anon during upkeep too.

TESTING OBJECTIVES:

The primary goal of testing is to uncover a bundle blunder, efficiently, the base exertion and time beginning officially, we can say

1. Test is that the way toward executing a program with the goal of discovering a slip.
2. A successful test is one that uncovers but then unseen blunder.
3. A good lawful activity is one that includes a high likelihood of discovering blunders, on the off chance that it exits.

WHITE BOX TESTING:

1. This can be unit trying technique where the unit are taken at a time and tested completely at a declaration level to look out the most extreme level blunders.
2. We have tried advance savvy each piece of code, dealing with every assertion inside thecode.
 - 2.1 Is executed a base once.
 - 2.2 The white box testing is also called glass box testing.

Discovery TESTING:

This testing strategy models one unit and checks the unit at interface and correspondence with different models somewhat going exhaustively levels. Here the model will be treated as an account machine that take input and creates the yield. Yield of given information blends are sent to different models.

UNIT TESTING:

Unit testing centers confirmation exertion around the tiniest unit of programming that is the model utilizing the nitty gritty plan and along these lines the cycle determinations testing is finished to reveal blunders inside the limit of the model all models should make progress inside the unit test before the start of the blend testing. In our venture unit testing includes looking at every future laid in the segment a segment performs just little a section of} the usefulness on the framework and depends on helping out other piece of the framework.

Joining TESTING:

In this venture coordinating every one of the modules shapes the principle framework when incorporating every one of the modules we have checked whether the reconciliation impacts working of any of the administrations by giving various mixes of contributions with which the administrations run impeccably before combination.

Table: 7.1 TEST CASES

1	Test Case	Upload Dataset
2	Precondition	Loading the dataset Select the path of the dataset file.
3.	Description	Loaded dataset should be valid and data should not contain null values and file should be available in the selected path.
4.	Test Steps	Load the Dataset. Select the path of the dataset
5.	Expected Output	Based o the dataset retrieved from the path data should be loaded into the system.
6.	Actual Output	Based on the Dataset retrieved from the path data should be loadedto the system
7.	Status	Success

8. CONCLUSION

By dissecting the outcomes, we presumed that ended up being best model for forecast of notion examination utilizing Machine Learning and this model creates exact outcomes with high precision.

We pick five famous classifiers thinking about their exhibition for the venture. We select one dataset which is accessible at www.kaggle.com dataset archive. To look at the grouping execution of learning, classifiers are applied on same information and results are thought about based on misclassification and right characterization rate, it tends to be presumed that each survey of the dataset gives best exactness.

9. FUTURE ENHANCEMENTS

In this paper, film surveys are arranged into positive or negative extremity. The framework can be utilized to arrange a gigantic information base of film audits. This will help film makers to check the situation with their film. Future work, this API can be prepared for different audits like PDAs, workstations or garments and so on

Also, future work can attempt more convoluted models for the examination. For instance, repetitive neural organization might have the option to give better execution since it can additionally represent the relationship of the sentences.

10. BIBLIOGRAPHY

References:

- [1] Tripathy, A. Agrawal, and S. K. Rath, "Arrangement of Sentimental Reviews Using Machine Learning Techniques", third International Conference on Recent Trends in Computing 2015 (ICRTC-2015), Procedia Computer Science, vol. 57, 2015, pp. 821 – 829.
- [2] R. Feldman, "Strategies and applications for slant examination," Communications of the ACM, vol. 56, 2013, pp. 82–89.
- [3] A.K. Sharma, S. Chaurasia, and D. K. Srivastava, "Nostalgic Short Sentences Classification by Using CNN Deep Learning Model with Fine Tuned Word2Vec", International Conference on Computational Intelligence and Data Science (ICCIDIS 2019), Procedia Computer Science vol.167, 2020, 1139–1147.
- [4] P. Vijayaragavan, R. Ponnusamy, and M. Aramudhan, "An ideal help vector machine based arrangement model for wistful examination of online item surveys", Future Generation Computer Systems, vol. 111, 2020, 234–240.
- [5] A.Kub,"Sentiment Analysis with Python (Part1)", [https://towardsdatascience.com/sentimentanalysis-with-python-section 1-5ce197074184](https://towardsdatascience.com/sentimentanalysis-with-python-section-1-5ce197074184), got to on Jun. 5, 2020. (Part2)", [https://towardsdatascience.com/sentimentanalysis-with-python-section 2-4f71e7bde59a](https://towardsdatascience.com/sentimentanalysis-with-python-section-2-4f71e7bde59a), got to on Jun. 5, 2020.
- [6] S. Basal, "A Comprehensive Guide to Understand and Implement Text Classification in Python", [https://www.analyticsvidhya.com/blog/2018/04/a-exhaustive manual for comprehend and implement-text-characterization in-python/](https://www.analyticsvidhya.com/blog/2018/04/a-exhaustive-manual-for-comprehend-and-implement-text-characterization-in-python/), got to on Jun. 5, 2020.
- [7] G. James, D. Witten, T. Hastie, and R. Tibshirani, "An Introduction to Statistical Learning: with Applications in R", Springer Publishing Company, Incorporated, 2014.
- [8] T. Hastie, R. Tibshirani, J. H. Friedman, "The components of measurable learning: information mining, surmising, and expectation. Second ed", New York: Springer, 2009.
- [9] M. Tengyu, A. Avati, K. Katanforoosh, and A. Ng, "CS 229 AI", class gift, Stanford University, 2020.
- [10] J. Leskovec, A. Rajaraman, and J. D. Ullman. "Mining of Massive Datasets (2nd.ed.)", Cambridge University Press, USA, Chapter 1, pp. 8-19, 2014.

WEB RESOURCES:

Large Movie Review Dataset – <http://ai.stanford.edu/~amaas/data/sentiment/>

Glove embeddings: <http://nlp.stanford.edu/data/glove.6B.zip>

TEXT BOOKS:

- Python Machine Learning.
- Hands-ON Machine Learning with Scikit-Learn & Tensor Flow.
- Fundamentals of Machine Learning for Predictive Data Analytics: Algorithms, Worked Examples, and Case Studies.