

**A PROJECT REPORT**  
**ON**  
**“NEWS VALIDATION SYSTEM”**

**Submitted to**  
**UNIVERSITY OF MUMBAI**

**In Partial Fulfilment of the Requirement for the Award of**

**BACHELOR’S DEGREE IN**  
**COMPUTER ENGINEERING**

**BY**

<b>SHAIKH MOHD NOORALAM MOHD KHAIRULALAM BILKEESH</b>	<b>14CO48</b>
<b>PAWASKAR SUFIYAN SIRAJ SULTANA</b>	<b>14CO37</b>
<b>RANE PRADNYESH NANDKISHOR NAMRATA</b>	<b>14CO41</b>

**UNDER THE GUIDANCE OF**  
**PROF. TABREZ KHAN**



**DEPARTMENT OF COMPUTER ENGINEERING**  
**Anjuman-I-Islam's Kalsekar Technical Campus**  
**SCHOOL OF ENGINEERING & TECHNOLOGY**

**Plot No. 2 3, Sector - 16, Near Thana Naka,**  
**Khandagaon, New Panvel - 410206**

**2017-2018**

**AFFILIATED TO**  
**UNIVERSITY OF MUMBAI**

**A PROJECT II REPORT  
ON**

**“NEWS VALIDATION SYSTEM”**

**Submitted to  
UNIVERSITY OF MUMBAI**

**In Partial Fulfilment of the Requirement for the Award of**

**BACHELOR’S DEGREE IN  
COMPUTER ENGINEERING**

**BY**

<b>SHAIKH MOHD NOORALAM MOHD KHAIRULALAM BILKEESH</b>	<b>14CO48</b>
<b>PAWASKAR SUFIYAN SIRAJ SULTANA</b>	<b>14CO37</b>
<b>RANE PRADNYESH NANDKISHORE NAMRATA</b>	<b>14CO41</b>

**UNDER THE GUIDANCE OF  
PROF. TABREZ KHAN**



**DEPARTMENT OF COMPUTER ENGINEERING  
Anjuman-I-Islam's Kalsekar Technical Campus  
SCHOOL OF ENGINEERING & TECHNOLOGY  
Plot No. 2 3, Sector - 16, Near Thana Naka,  
Khandagaon, New Panvel - 410206**

**2017-2018  
AFFILIATED TO**



**UNIVERSITY OF MUMBAI**

# Anjuman-i-Islam's Kalsekar Technical Campus

Department of Computer Engineering  
SCHOOL OF ENGINEERING & TECHNOLOGY  
Plot No. 2 3, Sector - 16, Near Thana Naka,  
Khandagaon, New Panvel - 410206



## CERTIFICATE

This is certify that the project entitled

**“News Validation System“**

submitted by

<b>SHAIKH MOHD NOORALAM MOHD KHAIRULALAM BILKEESH</b>	<b>14CO48</b>
<b>PAWASKAR SUFIYAN SIRAJ SULTANA</b>	<b>14CO37</b>
<b>RANE PRADNYESH NANDKISHOR NAMRATA</b>	<b>14CO41</b>

is a record of bonafide work carried out by them, in the partial fulfilment of the requirement for the award of Degree of Bachelor of Engineering (Computer Engineering) at *Anjuman-I-Islam's Kalsekar Technical Campus, Navi Mumbai* under the University of MUMBAI. This work is done during year 2017-2018, under our guidance.

**Date:**     /     /

**Prof. TABREZ KHAN**  
Project Supervisor

**Prof. KALPANA BODKE**  
Project Coordinator

**Prof. TABREZ KHAN**  
HOD, Computer Department

**DR. ABDUL RAZAK HONNUTAGI**  
Director

**External Examiner**

## Acknowledgements

We would like to take the opportunity to express my sincere thanks to my guide **TABREZ KHAN**, Professor, Department of Computer Engineering, AIKTC , Panvel for his invaluable support and guidance throughout my project research work. Without his kind guidance & support this was not possible.

We are grateful to him for his timely feedback which helped me track and schedule the process effectively. His time, ideas and encouragement that he gave is help me to complete my project efficiently.

We would like to express deepest appreciation towards **DR. ABDUL RAZAK HONNUTAGI**, Director, AIKTC, Navi Mumbai, **Prof. TABREZ KHAN**, Head of Department of Computer Engineering and **Prof. KALPANA BODKE**, Project Coordinator whose invaluable guidance supported us in completing this project.

At last we must express our sincere heartfelt gratitude to all the staff members of Computer Engineering Department who helped me directly or indirectly during this course of work.

SHAIKH MOHD NOORALAM MOHD KHAIRULALAM BILKEESH

PAWASKAR SUFIYAN SIRAJ SULTANA

RANE PRADNYESH NANDKISHOR NAMRATA



## Project I Approval for Bachelor of Engineering

This project entitled *News Validation System* by **SHAIKH MOHD NOORALAM MOHD KHAIRULALAM BILKEESH, PAWASKAR SUFIYAN SIRAJ SULTANA, RANE PRADNYESH NANDKISHOR NAMRATA** is approved for the degree of *Bachelor of Engineering in Department of Computer Engineering*.

Examiners

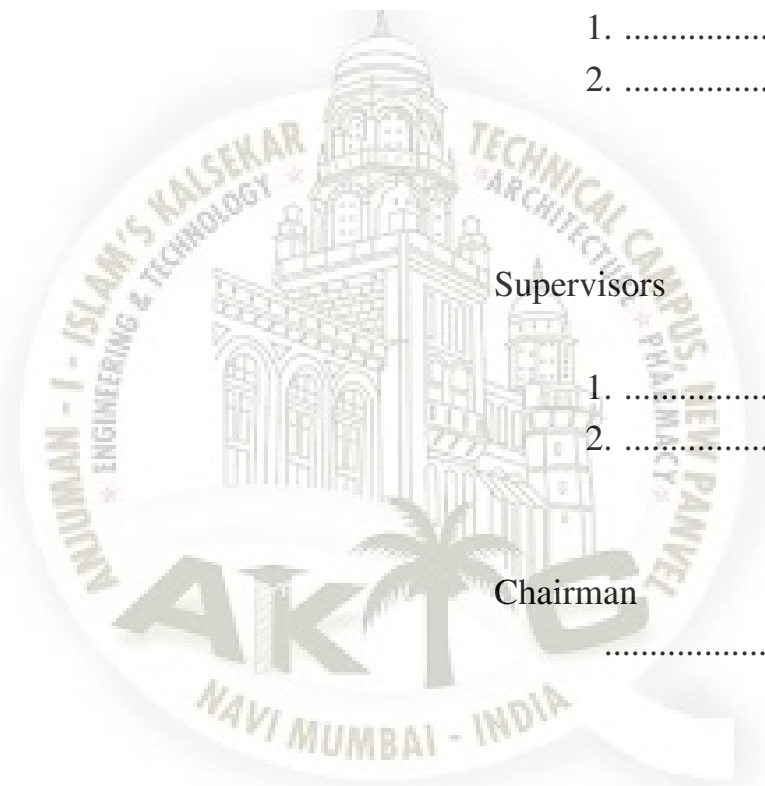
1. ....
2. ....

Supervisors

1. ....
2. ....

Chairman

.....



## Declaration

We declare that this written submission represents my ideas in my own words and where others ideas or words have been included, We have adequately cited and referenced the original sources. We also declare that we have adhered to all principles of academic honesty and integrity and have not misrepresented or fabricated or falsified any idea/data/fact/source in my submission. We understand that any violation of the above will be cause for disciplinary action by the Institute and can also evoke penal action from the sources which have thus not been properly cited or from whom proper permission has not been taken when needed.

SHAIKH MOHD NOORALAM MOHD KHAIRULALAM BILKEESH  
14CO48

PAWASKAR SUFIYAN SIRAJ SULTANA  
14CO37

RANE PRADNYESH NANDKISHOR NAMRATA  
14CO41

# ABSTRACT

Title: News Validation System

The news is an information of particular interest that was not known in the past, the information implied by news should be accurate and reliable. There are millions of news shared over social media each day without validation. Currently, there exists no resource that could validate whether any news article is giving valid information or not. Here comes the existence of News Validation System. There are various news analysis methods that would validate news. Our news validation system implements methods that would be more feasible, accurate and reliable out of all methods present for news validation. Our system will go with a hybrid approach. A hybrid approach is the combination of both linguistic and network approach, this approach is developed by integrating advantages of linguistic and network approach and by eliminating their drawbacks. It has components like authentic source identification, date validation, data collection, comparison and support analysis. All these parameters will conclude in well formed results. The system will feed upon a normal text of information in the context of news articles, the output of the system will be a probability percentage, that would tell at which extent the information is reliable and valid. Our work in news validation system can help cut down potential misleading information in the news and we can make sure that everyone gets reliable, dependable, authentic and accurate information.

**Keywords:** Validation, Analysis, Natural Language Processing,

# Contents

Acknowledgement . . . . .	iii
Project I Approval for Bachelor of Engineering . . . . .	iv
Declaration . . . . .	v
Abstract . . . . .	vi
Table of Contents . . . . .	ix
<b>1 Introduction</b>	<b>2</b>
1.1 Statement of the Project . . . . .	2
1.2 Purpose . . . . .	3
1.3 Project Scope . . . . .	3
1.4 Project Goals and Objectives . . . . .	3
1.4.1 Goals . . . . .	3
1.4.2 Objectives . . . . .	3
1.5 Organization of Report . . . . .	3
<b>2 Literature Survey</b>	<b>5</b>
2.1 Fake News or Truth? Using Satirical Cues to Detect Potentially Mis- leading News. . . . .	5
2.1.1 Description . . . . .	5
2.1.2 Advantages of Paper . . . . .	6
2.1.3 Disadvantages of Paper . . . . .	6
2.1.4 How to overcome the problems mentioned in Paper . . . . .	7
2.2 Automatic Deception Detection - Methods for Finding Fake News . . . . .	7
2.2.1 Description . . . . .	7
2.2.2 Advantages of Paper . . . . .	8
2.2.3 Disadvantages of Paper . . . . .	8
2.2.4 How to overcome the problems mentioned in Paper . . . . .	8
2.3 Developing a News Aggregation and Validation System. . . . .	8
2.3.1 Description . . . . .	8
2.3.2 Advantages of Paper . . . . .	9
2.3.3 Disadvantages of Paper . . . . .	9
2.3.4 How to overcome the problems mentioned in Paper . . . . .	9
2.4 Technical Review . . . . .	9
2.4.1 Textblob . . . . .	9

2.4.2	Advantages of Technology . . . . .	9
2.4.3	Reasons to use this Technology . . . . .	10
2.4.4	News API . . . . .	10
2.4.5	Reasons to use this technology . . . . .	10
2.4.6	Rake NLTK . . . . .	10
2.4.7	Reasons to use this technology . . . . .	10
<b>3</b>	<b>Project Planning</b>	<b>11</b>
3.1	Members and Capabilities . . . . .	11
3.2	Roles and Responsibilities . . . . .	11
3.3	Assumptions and Constraints . . . . .	11
3.4	Project Management Approach . . . . .	12
3.5	Ground Rules for the Project . . . . .	12
3.6	Project Budget . . . . .	13
3.7	Project Timeline . . . . .	13
<b>4</b>	<b>Software Requirements Specification</b>	<b>14</b>
4.1	Overall Description . . . . .	14
4.1.1	Product Perspective . . . . .	14
4.1.2	Product Features . . . . .	14
4.1.3	User Classes and Characteristics . . . . .	14
4.1.4	Operating Environment . . . . .	14
4.1.5	Design and Implementation Constraints . . . . .	14
4.2	System Features . . . . .	15
4.2.1	Scraper . . . . .	15
4.2.2	Keyword Extractor . . . . .	15
4.2.3	Properties Generator . . . . .	15
4.3	External Interface Requirements . . . . .	16
4.3.1	User Interfaces . . . . .	16
4.3.2	Hardware Interfaces . . . . .	16
4.3.3	Software Interfaces . . . . .	16
4.3.4	Communications Interfaces . . . . .	16
4.4	Nonfunctional Requirements . . . . .	17
4.4.1	Performance Requirements . . . . .	17
4.4.2	Safety Requirements . . . . .	17
4.4.3	Security Requirements . . . . .	17
<b>5</b>	<b>System Design</b>	<b>18</b>
5.1	System Requirements Definition . . . . .	18
5.1.1	Functional requirements . . . . .	18
5.1.2	System requirements (non-functional requirements) . . . . .	21
5.2	System Architecture Design . . . . .	22

5.3	Sub-system Development . . . . .	23
5.3.1	Authentic Source Identification . . . . .	23
5.3.2	Date Validation . . . . .	24
5.3.3	Data collection and comparison . . . . .	25
5.3.4	Support Analysis . . . . .	26
5.4	Systems Integration . . . . .	27
5.4.1	Class Diagram . . . . .	27
5.4.2	Sequence Diagram . . . . .	28
5.4.3	Component Diagram . . . . .	29
5.4.4	Deployment Diagram . . . . .	30
<b>6</b>	<b>Implementation</b>	<b>31</b>
6.1	Text Analysis Module . . . . .	31
6.2	Source, Author , Date Check Module . . . . .	33
6.3	Support Check Module. . . . .	36
6.4	Content Check Module . . . . .	42
6.5	Classifier . . . . .	43
<b>7</b>	<b>System Testing</b>	<b>45</b>
7.1	Test Cases and Test Results . . . . .	45
7.2	Sample of a Test Case . . . . .	45
7.2.1	Software Quality Attributes . . . . .	46
<b>8</b>	<b>Screenshots of Project</b>	<b>48</b>
8.1	Registration of client . . . . .	48
8.2	Login of client . . . . .	49
8.3	Api Key Generation for client . . . . .	50
8.4	Output for the search data . . . . .	51
<b>9</b>	<b>Conclusion and Future Scope</b>	<b>52</b>
9.1	Conclusion . . . . .	52
9.2	Future Scope . . . . .	52
	<b>References</b>	<b>52</b>
	<b>Achievements</b>	<b>53</b>

## List of Figures

3.1	<b>Spiral model</b> . . . . .	12
3.2	<b>Gantt chart</b> . . . . .	13
5.1	<b>Use Case</b> . . . . .	19
5.2	<b>Level 0 dfd for News Validation System</b> . . . . .	20
5.3	<b>Level 1 dfd for News Validation System</b> . . . . .	20
5.4	<b>Level 2 dfd for News Validation System</b> . . . . .	21
5.5	<b>Entity Relationship for News Validation System</b> . . . . .	22
5.6	<b>System Architecture for News Validation System</b> . . . . .	23
5.7	<b>Modular diagram for authentic source identification</b> . . . . .	24
5.8	<b>Modular diagram for date validation</b> . . . . .	25
5.9	<b>Modular diagram for data collection and comparison</b> . . . . .	26
5.10	<b>Modular diagram for support analysis</b> . . . . .	27
5.11	<b>Class diagram for News Validation System</b> . . . . .	28
5.12	<b>Sequence Diagram for News Validation System</b> . . . . .	29
5.13	<b>Component diagram for News Validation System</b> . . . . .	30
5.14	<b>Deployment diagram for News Validation System</b> . . . . .	30
7.1	<b>Key generation after login</b> . . . . .	46
8.1	<b>Register</b> . . . . .	48
8.2	<b>Login</b> . . . . .	49
8.3	<b>Key generation after login</b> . . . . .	50
8.4	<b>Output</b> . . . . .	51



# List of Tables

3.1	Table of Capabilities . . . . .	11
3.2	Table of Responsibilities . . . . .	11



# Chapter 1

## Introduction

### 1.1 Statement of the Project

News is nothing but information about current events. The arrival of the web and the social web brings with it a tremendous number of new sources. The accessibility of these news sources generates a large amount of information which can often times be contradicting and confusing. Facebook, for example, can be seen as a social platform that allows individuals and groups of individuals to freely exchange thoughts and opinions.

When this information travels the social web, it is difficult to distinguish between valid and unsupported news. News verification aims to implement a technology that can identify fake news. Fake news detection is defined as a news which is intentionally altered. The problem of fake news detection is more challenging and complicated task than detecting deceptive news, since the political language on TV interviews, post on Facebook and Twitter are mostly short statements

There are various approaches that can be used to develop news validation or fake news detection system. Majorly there are two types of approaches, linguistic approach and network approach. In Linguistic approach, some liar uses their language skill to avoid being caught guilty. There is some leakage of words from which we can identify that whether they are saying truth or not.

The goal in the linguistic approach is to look for such words or leakages. Network approach is innovative and varied, using network properties and behavior are ways to complement content- based approaches that rely on deceptive language and leakage cues to predict deception. Hybrid approach is the combination of both network approach as well as linguistic approach. In our system we are going to use hybrid approach because individually network or linguistic approach is not too accurate to increase efficiency and for better results we shall be using hybrid approach.

## 1.2 Purpose

In social media articles often decontextualised from source, fact can mix freely with fiction. The sharing of hoax news often results in defamation of certain entity, it plays with emotions of readers. Unverified news often causes loss of lives, initiates riots, loss in business. Often people share unverified news or text, without verifying the source and causes above problems. The content in news report are framed in such way that it provokes emotionally and gain reaction. It is high time to stop this uncontrolled flow of such news.

## 1.3 Project Scope

Our system will validate news at global level where it will only support text of the news. The system will help in filtering the fake news available on the internet and will serve quality news to the readers. This will increase the efficiency and reliability of the digital news and can be as trustworthy as the printed newspapers.

## 1.4 Project Goals and Objectives

### 1.4.1 Goals

The goal of the News Validation System is to explore how artificial intelligence technologies, particularly machine learning and natural language processing, might be leveraged to combat the fake news problem.

### 1.4.2 Objectives

The Objective of our system is to validate the news so that the user will get correct data.

We believe that these AI technologies hold promise for significantly automating parts of the procedure human fact checkers use today to determine if a story is real or a hoax.

## 1.5 Organization of Report

The remaining part of the project is organized as follows.

Chapter 2 presents a review of related work.

Chapter 3 describes the time management and time utilization during the project

implementation.

Chapter 4 introduces the Software requirement Specification of our project.

Chapter 5 proposes the project design of the project. It represents the architectural design, front end design and database design of the project.

Chapter 6 presents implementation details of our project.

Chapter 7 presents various test cases that are considered.

Chapter 8 consists of various screenshots of the project.

Chapter 9 provides some concluding remarks and direction of our future work.



## Chapter 2

# Literature Survey

### 2.1 Fake News or Truth? Using Satirical Cues to Detect Potentially Misleading News.

#### 2.1.1 Description

This paper's authors are Niall J. Conroy, Victoria L. Rubin, Yimin Chen was published in the Proceedings of the Workshop on Computational Approaches to Deception Detection at the 15th Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-CADD2016) on June 17, 2016 in San Diego, California.

This research paper enables the identification of deliberately deceptive misinformation in the institutional mainstream or non-institutional text-based online news. The resulting deception detection methodology will allow making predictions about each previously unseen news piece: is it likely to belong to the truthful or deceptive category? A system, based on this methodology, will alert users to most likely deceptive news in the incoming stream of news, and prompt the users to fact-check further.

Digital deception is a deliberate effort to create false beliefs or conclusions in technology-mediated environments. This research project focuses on deliberate misinformation in text-based online news, provided via mainstream media and citizen journalist websites, news archives and aggregators. Various deception types and degrees will be examined, categorized, and modeled: fake or fabricated news, exaggerated claims, material fact omissions, indirect responses, question-dodging, and subject-changing.

Mistaking deceptive news for authentic reports can create costly negative consequences such as sudden stock fluctuations or reputation loss. Everyday life decision-making, behavior, and mood are influenced by news we receive. When professional

analysts sift through the news, their future forecasts, fact and pattern discovery depend on veracity of the news in “big data” knowledge management and curation areas (specifically, in business intelligence, financial and stock market analysis, or national security and law enforcement). In both lay and professional contexts of news consumption, it is critical to distinguish truthful reports from deceptive ones. However, few news verification mechanisms currently exist, and the sheer volume of the information requires novel automated approaches.

Mistaking deceptive news for authentic reports can create costly negative consequences such as sudden stock fluctuations or reputation loss. Everyday life decision-making, behavior, and mood are influenced by news we receive. When professional analysts sift through the news, their future forecasts, fact and pattern discovery depend on veracity of the news in “big data” knowledge management and curation areas (specifically, in business intelligence, financial and stock market analysis, or national security and law enforcement). In both lay and professional contexts of news consumption, it is critical to distinguish truthful reports from deceptive ones. However, few news verification mechanisms currently exist, and the sheer volume of the information requires novel automated approaches.

News verification methods and tools are timely and beneficial to both lay and professional text-based news consumers.

### **2.1.2 Advantages of Paper**

The research significance is four-fold:

- a. Automatic analytical methods complement and enhance the notoriously poor human ability to discern information from misinformation.
- b. Credibility assessment of digital news sources is improved.
- c. The mere awareness of potential digital deception constitutes part of new media literacy and can prevent undesirable consequences.
- d. The proposed veracity/deception criterion is also seen as a metric for information quality assessment.

### **2.1.3 Disadvantages of Paper**

- a. This method works in the context of textual news information only.

### 2.1.4 How to overcome the problems mentioned in Paper

- a. To target other than textual means enabling support for images as well, the images may more influence than text, as it may involve visual scenes which may be obscene in nature, or not usual for normal user, indirectly it can play with reader's intellect.
- b. The images involved in the form of images can be scraped by using Google's Image Vision API, and then from the results we can extract origin, date of the image, as in fake news the image must have been used from past incident or any other incident must be framed as recent. Also, the description of searched results can be compared with the input news.

## 2.2 Automatic Deception Detection - Methods for Finding Fake News

### 2.2.1 Description

This paper's authors are Niall J. Conroy, Victoria L. Rubin, Yimin Chen was published by ASIST (Association for Information Science and Technology) on November 6-10, 2015 in St. Louis, MO, USA

News verification aims to employ technology to identify intentionally deceptive news content online, and is an important issue within certain streams of library and information science.

This paper provides researchers with a map of the current landscape of veracity deception assessment methods, their major classes and goals, all with the aim of proposing a hybrid approach to system design. These methods have emerged from separate development streams, utilizing disparate techniques. In this survey, two major categories of methods emerge: 1. Linguistic Approaches in which the content of deceptive messages is extracted and analyzed to associate language patterns with deception; and 2. Network Approaches in which network information, such as message metadata or structured knowledge network queries can be harnessed to provide aggregate deception measures.

Both forms typically incorporate machine learning techniques for training classifiers to suit the analysis. It is incumbent upon researchers to understand these different areas, yet no known typology of methods exists in the current literature. The goal is to provide a survey of the existing research while proposing a hybrid approach, which utilizes the most effective deception detection methods for the implementation of a fake news detection tool.



Structured datasets are easier to verify than non-structured (or semi-structured) data such as texts. When we know the language domain (e.g., insurance claims or health-related news) we can make better guesses about the nature and use of deception. Semi-structured non-domain specific web data come in many formats and demand flexible methods for veracity verification.

### **2.2.2 Advantages of Paper**

- a. Success was measured based on whether the machine was able to assign higher true values to true statements than to false ones.
- b. Linguistic and network-based approaches have shown high accuracy results in classification tasks within limited domains.

### **2.2.3 Disadvantages of Paper**

- a. A problem with this method is that statements must reside in a preexisting knowledge base.

### **2.2.4 How to overcome the problems mentioned in Paper**

By using Hybrid approach that is combination of linguistic and network approach we can overcome the above mentioned problem.

## **2.3 Developing a News Aggregation and Validation System.**

### **2.3.1 Description**

Fake news detection is defined as the prediction of the chances of a particular news article (news report, editorial, expose, etc.) being intentionally deceptive. The Linguistic Approach can be described as a method where the content of an item gets extracted and analyzed regarding language patterns. Sometimes a "leakage" occurs, meaning that a break in the pattern is observable. Thus the following aspects need to be discussed: stop words, stemming, lemmatization. Some words do not add value to the context of a text, thus those words do not need to be processed by the system.

Those words are described as common terms or stop words. Nonetheless, it is also stated that in some circumstances the use of stop words are essential in order to keep the context. Articles and other forms of text use different forms of a word (e.g. conjunctions) due to grammatical reasons. Examples, therefore, are 'am, are, is,' which are different forms of 'be' and 'cars, cars', 'car's' which results in car. The process of Stemming connotes that the affixes are removed from a word, meaning the ends of words are cut off, potentially resulting in the base form of a word. Lemmatization

on the other hand uses 'vocabulary and morphological analysis of words in order to reach the goal of finding the base of a word.

### 2.3.2 Advantages of Paper

- a. It is very important subset of Hybrid Approach and can be collectively used with other methods to detect fake news very accurately.

### 2.3.3 Disadvantages of Paper

- a. Linguistic Approach do not validate on the basis of real time news, it just focus on linguistic leakages in the textual news.
- b. Linguistic Approach do not provide accuracy if used independently with other method.

### 2.3.4 How to overcome the problems mentioned in Paper

Major problem with linguistic method was it only focus on linguistic leakages and not on realtime information related to that particular news and network method only focus on information about that news not about linguistic leakages so if both this method is combined as hybrid approach then major disadvantages for linguistic method can be overcome and accuracy can also be obtained for news validation.

## 2.4 Technical Review

The technology that we are using here in Python such as textblob, rake nltk, goose, newsapi.

### 2.4.1 Textblob

TextBlob is a Python (2 and 3) library for processing textual data. It provides a simple API for diving into common natural language processing (NLP) tasks such as part-of-speech tagging, noun phrase extraction, sentiment analysis, classification, translation, and more. Its is most widely used in our project like in support check module for analysis sentiments of tweets extracted.

### 2.4.2 Advantages of Technology

- a. Its very efficient
- b. Since, it is built on the shoulders of NLTK and Pattern, therefore making it simple for beginners by providing an intuitive interface to NLTK.

- c. It provides language translation and detection which is powered by Google Translate ( not provided with Spacy).

### **2.4.3 Reasons to use this Technology**

- a. It is very well designed, fast and scalable.
- b. Other library in python for do the same thing are not reliable
- c. It gives result in string format that can be easily handle.

### **2.4.4 News API**

News API is a simple HTTP REST API for searching and retrieving live articles from all over the web. Its gives the data from over 30,000 news sources.Its is freely available we just have to register on it once.After registration we get access to 1000 request per day.

### **2.4.5 Reasons to use this technology**

- a. It returns JSON metadata we can be easily use.
- b. It response time is good.
- c. Huge number of data sources that is over 30,000 news sources.

### **2.4.6 Rake NLTK**

RAKE short for Rapid Automatic Keyword Extraction algorithm, is a domain independent keyword extraction algorithm which tries to determine key phrases in a body of text by analyzing the frequency of word appearance and its co-occurrence with other words in the text.

### **2.4.7 Reasons to use this technology**

- a. As its name suggest its extract the keyword from the sentence given very rapidly.
- b. It is efficient way to extract keyword till the date.

# Chapter 3

## Project Planning

### 3.1 Members and Capabilities

Table 3.1: Table of Capabilities

SR. No	Name of Member	Capabilities
1	Nooralam Shaikh	UI Design
2	Sufiyan Pawaskar	UI Design , Database
3	Pradnyesh Rane	Database

Work Breakdown Structure

### 3.2 Roles and Responsibilities

Table 3.2: Table of Responsibilities

SR. No	Name of Member	Role	Responsibilities
1	Nooralam Shaikh	Team Leader	UI Design and core modules
2	Sufiyan Pawaskar	Team Member	UI Design ,Integration ,API design
3	Pradnyesh Rane	Team Member	Support check Module, Documentation

### 3.3 Assumptions and Constraints

#### Assumption

Assumption is that the data that is coming from the scrapper is in correct from and from genuine source.

#### Constraint

If we provide wrong training data to the system then prediction will be wrong. If news API from where we are scrapping the data goes down then the system will

fail.

### 3.4 Project Management Approach

Spiral Model is a combination of a waterfall model and iterative model. Each phase in spiral model begins with a design goal and ends with the client reviewing the progress. We have use this model for version control.

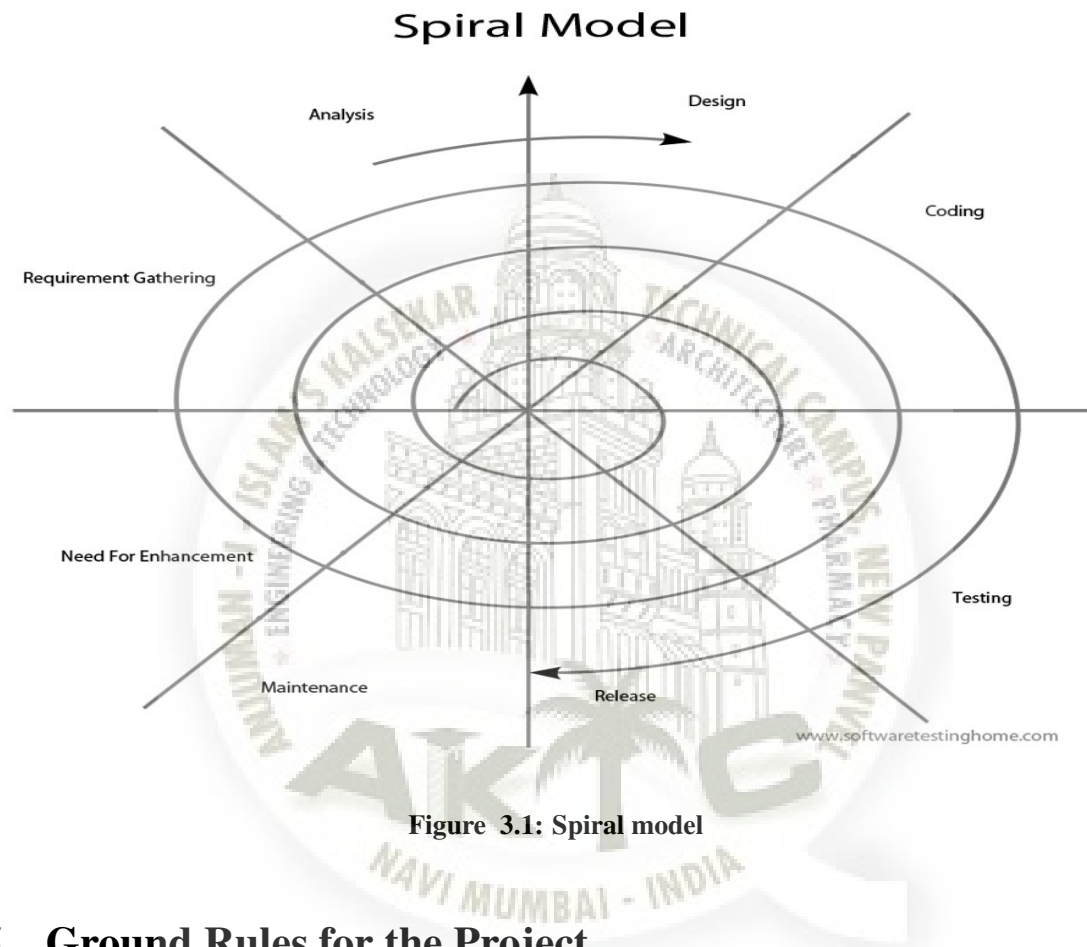


Figure 3.1: Spiral model

### 3.5 Ground Rules for the Project

1. We treat each other with respect.
2. We intend to develop personal relationships to enhance trust and open communication.
3. We value constructive feedback. We will avoid being defensive and give feedback in a constructive manner.
4. As team members, we will pitch in to help where necessary to help solve problems and catch-up on behind schedule work.
5. Additional meetings can be scheduled to discuss critical issues or tabled items upon discussion and agreement with the team leader.

6. One person talks at a time, there are no side discussions.
7. When we pose an issue or a problem, we will also try to present a solution.

### 3.6 Project Budget

The budget for this project is very low as most of the tools we have use are open source.Following are the budget for the project

1. Operating System: linux mint (Open Source).
2. Python Programming language (Open Source)
3. IDE:Andriod Studio (Open Source).
4. News API(Open Source)

### 3.7 Project Timeline

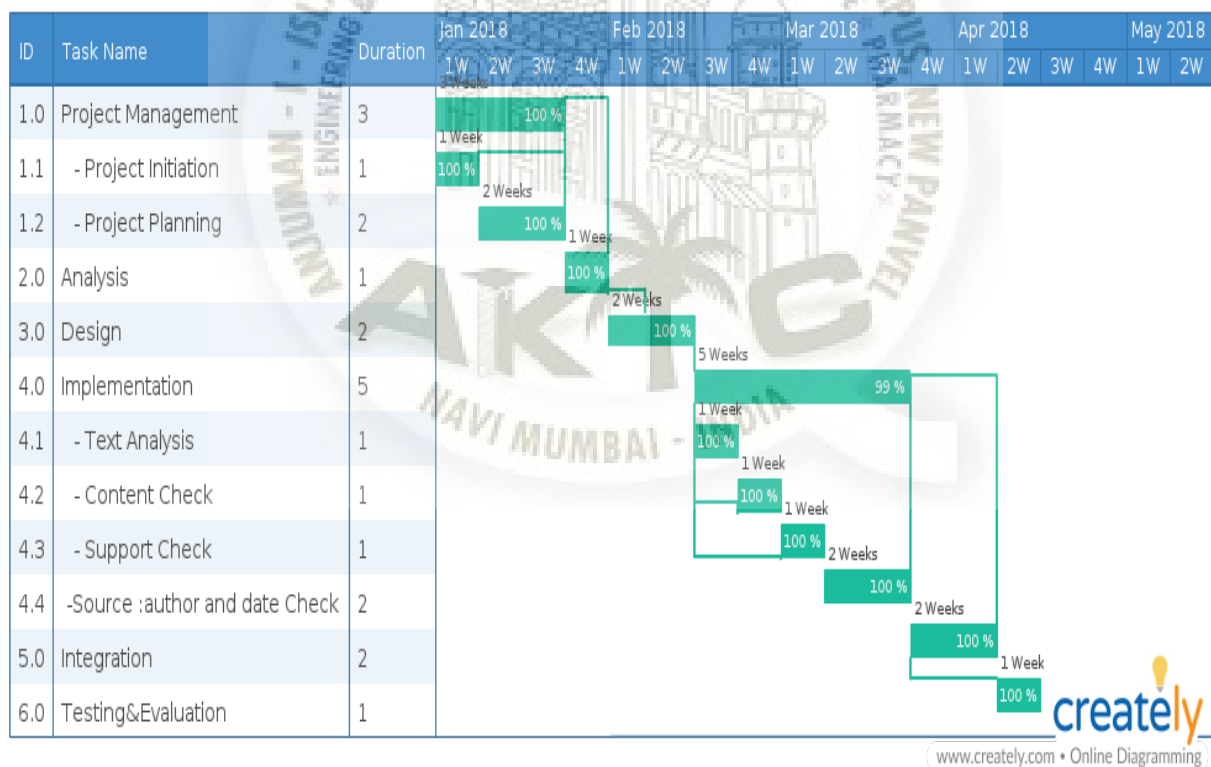


Figure 3.2: Gantt chart

## Chapter 4

# Software Requirements Specification

### 4.1 Overall Description

#### 4.1.1 Product Perspective

The origin of our product or our system is the need for validation of the news. It is a web based system implementing client-server model. The News Validation System provides simple mechanism for users to validate the news data.

#### 4.1.2 Product Features

News Validation System will provide api key to the client. Its will also validate the news given input to it showing the estimated value of the give data is true or false.

#### 4.1.3 User Classes and Characteristics

There are various users that can use the system such as normal user i.e any day to day life user who wants to check the content of the news is correct or not, any third party which use our project as a sub module of theirs.

#### 4.1.4 Operating Environment

Our System will provide both web based and app based support. One can use the website from theirs mobiles or computer they just required a normal web browser and also one can use mobiles application to do the same.

#### 4.1.5 Design and Implementation Constraints

The main decision was to select that which language should we use for developing our system. Since most of the programming language are very time consuming when it comes to data processing . So the complete system is made by using python language which is very efficient for processing the data as compare to the other language which are available in market since its provide various library to do our work simple.



## 4.2 System Features

Our system will provide an application programming interface to our client using which they can have a idea whether the input news is correct or not. There are various features as mention below.

### 4.2.1 Scrapper

It is the one of the main feature of our system.It will scrap the given data from various news sources or website in java script object notation format.

#### Description and Priority

It's is key player in our system it will extract news from over 4000 news website and get its data in java script object notation format to us.It is high priority module. Its main benefits is that it scrap the data from large number of reliable websites.

### 4.2.2 Keyword Extractor

It will only extract the keyword from the data given to it from scrapper.

#### Description and Priority

It will extract the keyword from the given the data given from the scrapper and forward those key word to properties generator. It is medium priority module.Its extract the keyword using Rapid Automatic Keyword Extraction algorithm.

### 4.2.3 Properties Generator

It will generate the properties using the data or the keyword from the keyword extractor.

#### Description and Priority

It will generate the values based on properties and then forward this values to the classifier.It is high priority module.

#### Stimulus/Response Sequences

1. The User must register in our system.
2. After Registration user must login in our system using its login credentials.In response to login the system will provide a api key.
3. By using api key provided by the system user can input for the data to validate.

4. In response to enter data the system will predict whether the enter data is correct or false. And it will return a value.
5. In case if the data is correct it will display that which is the source of the data.

### **Functional Requirements**

1. The user should register on the system User should login in the system.
2. The data input by the user must be in textual form.
3. The servers should response quickly.

## **4.3 External Interface Requirements**

### **4.3.1 User Interfaces**

1. User must first register itself on our system.
2. After registration user will get a API key through will he/she will be authenticated.
3. By using that key user can search for the data he/she needs to validate.

### **4.3.2 Hardware Interfaces**

There is no special kind of hardware required to run this system, A regular desktop or any smart phone can be enough to access our system.

### **4.3.3 Software Interfaces**

One must need a web browser or our mobile application to access our system. For database we used sequel database that is mysql. For developing core modules we used python programming language. We have used news API for extracting the required data. On the client side only a browser or our mobile application can be used.

### **4.3.4 Communications Interfaces**

The major communication between the client or user with our system is done using web-browser or an mobile application provided by us.

The communication between the system and database is done by using mysql.

## 4.4 Nonfunctional Requirements

### 4.4.1 Performance Requirements

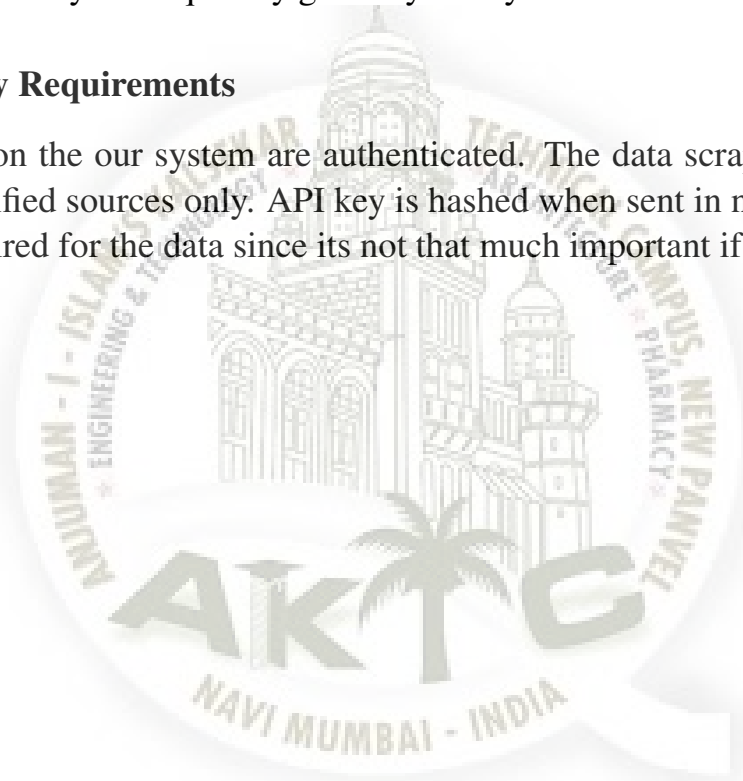
The main performance of this system depend on how much time it spends on extracting and processing the data. extracting the data from various news sources is very time consuming task but python made it simple by providing library. And also for processing the used language that is python is very efficient.

### 4.4.2 Safety Requirements

Once the client login in our system we will provide it a API key for authentication purpose.This API key is unique key given by our system to client.

### 4.4.3 Security Requirements

All the clients on the our system are authenticated. The data scrapped is also collected from verified sources only. API key is hashed when sent in network.No more security is required for the data since its not that much important if leaked.



# Chapter 5

## System Design

### 5.1 System Requirements Definition

System requirement definitions specify what the system should do, its functionality and its essential and desirable system properties. The techniques applied to elicit and collect information in order to create system specifications and requirement definitions involve consultations, interviews, requirements workshop with customers and end users. The objective of the requirements definition phase is to derive the two types of requirement:

#### 5.1.1 Functional requirements

They define the basic functions that the system must provide and focus on the needs and goals of the end users.

#### Use-case Diagram

Use case diagrams are usually referred to as behavior diagrams used to describe a set of actions (use cases) that some system or systems (subject) should or can perform in collaboration with one or more external users of the system (actors). Each use case should provide some observable and valuable result to the actors or other stakeholders of the system.

Use case diagrams are in fact twofold they are both behavior diagrams, because they describe behavior of the system, and they are also structure diagrams as a special case of class diagrams where classifiers are restricted to be either actors or use cases related to each other with associations.

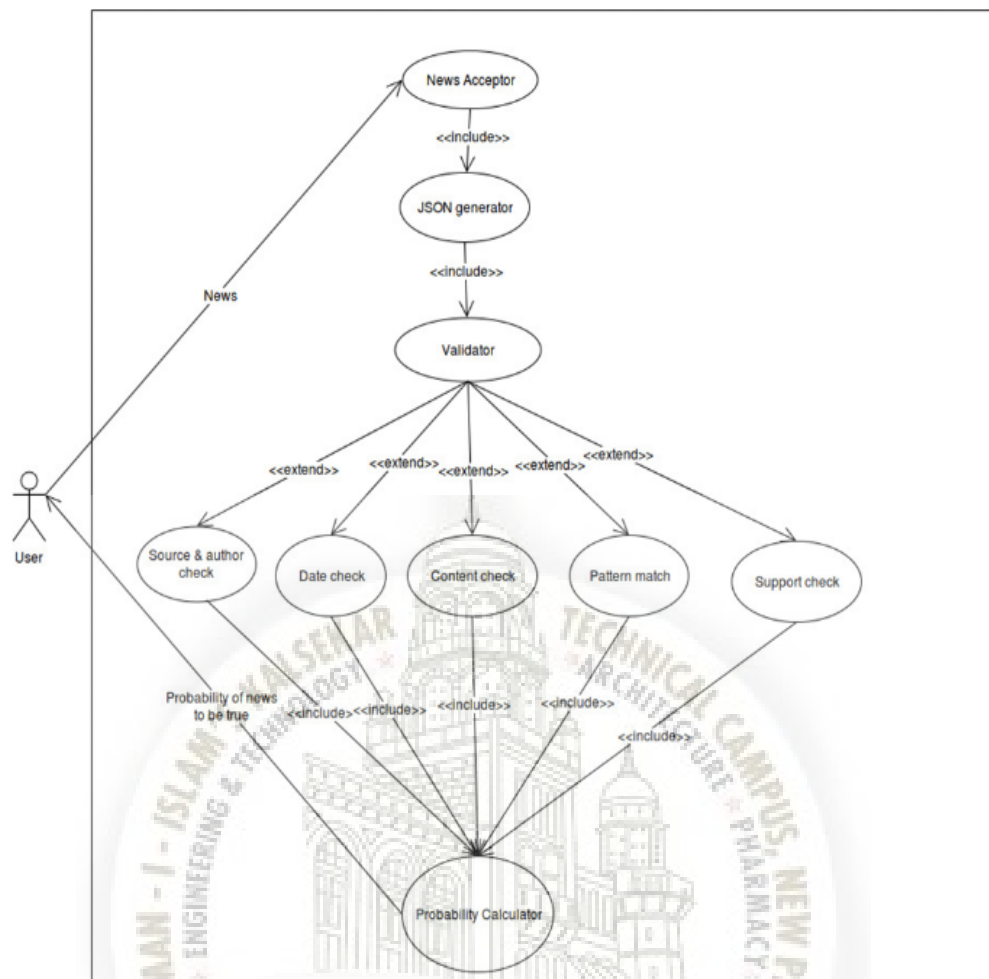


Figure 5.1: Use Case

### Data-flow Diagram

A data flow diagram (DFD) maps out the flow of information for any process or system. It uses defined symbols like rectangles, circles and arrows, plus short text labels, to show data inputs, outputs, storage points and the routes between each destination. Data flowcharts can range from simple, even hand-drawn process overviews, to in-depth, multi-level DFDs that dig progressively deeper into how the data is handled. They can be used to analyze an existing system or model a new one. Like all the best diagrams and charts, a DFD can often visually “say” things that would be hard to explain in words, and they work for both technical and nontechnical audiences, from developer to CEO. That’s why DFDs remain so popular after all these years. While they work well for data flow software and systems, they are less applicable nowadays to visualizing interactive, real-time or database-oriented software or systems.

Level 0

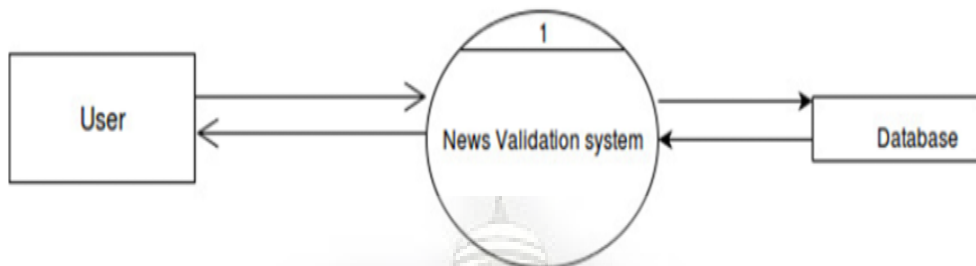


Figure 5.2: Level 0 dfd for News Validation System

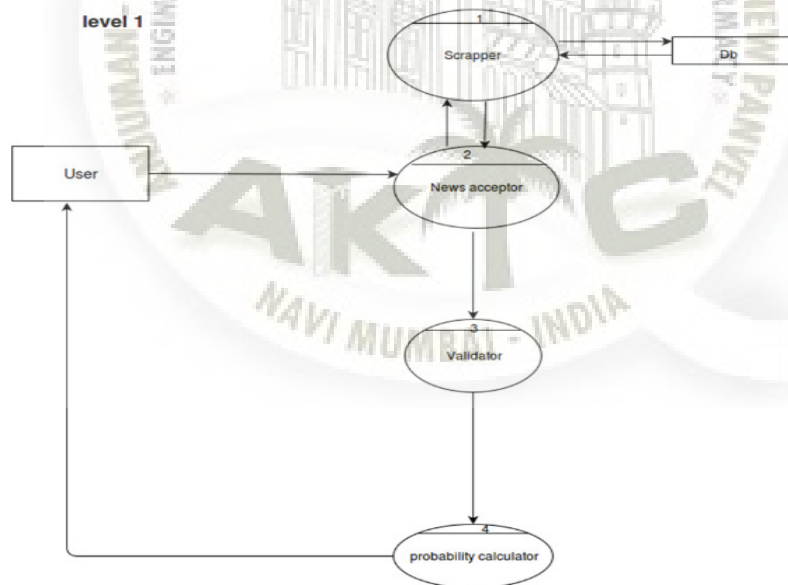


Figure 5.3: Level 1 dfd for News Validation System

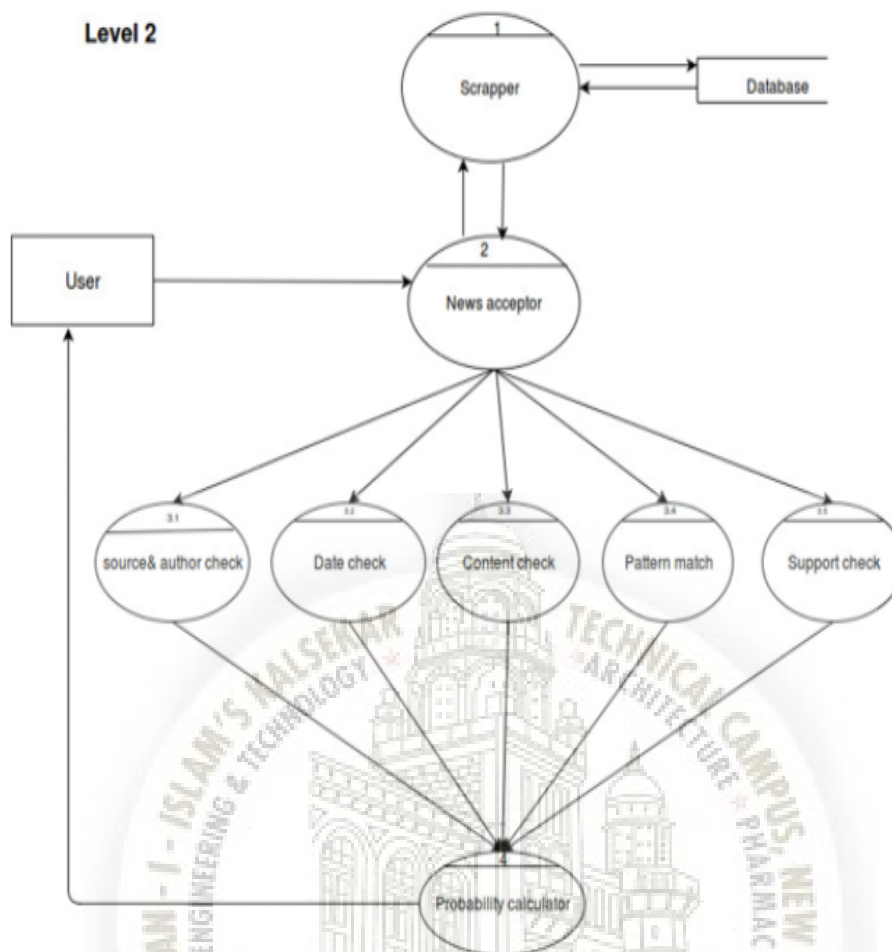


Figure 5.4: Level 2 dfd for News Validation System

### 5.1.2 System requirements (non-functional requirements)

#### Performance Requirements

The main performance of this system depend on how much time it spends on extracting and processing the data. Extracting the data from various news sources is very time consuming task but python made it simple by providing library. And also for processing the used language that is python is very efficient.

#### Safety Requirements

Once the news API server goes down then the whole system will go down. So to prevent from this we has to periodically check whether the server is up or not.

#### Security Requirements



The major security requirements for the system will be the safeguarding of the user data from any kind of exploit. In order to protect the user data the data is not stored in local databases we will be storing in the cloud for better security. And also the API key which is given to user once register for authentication is also encrypted.

### Database Schema/ E-R Diagram

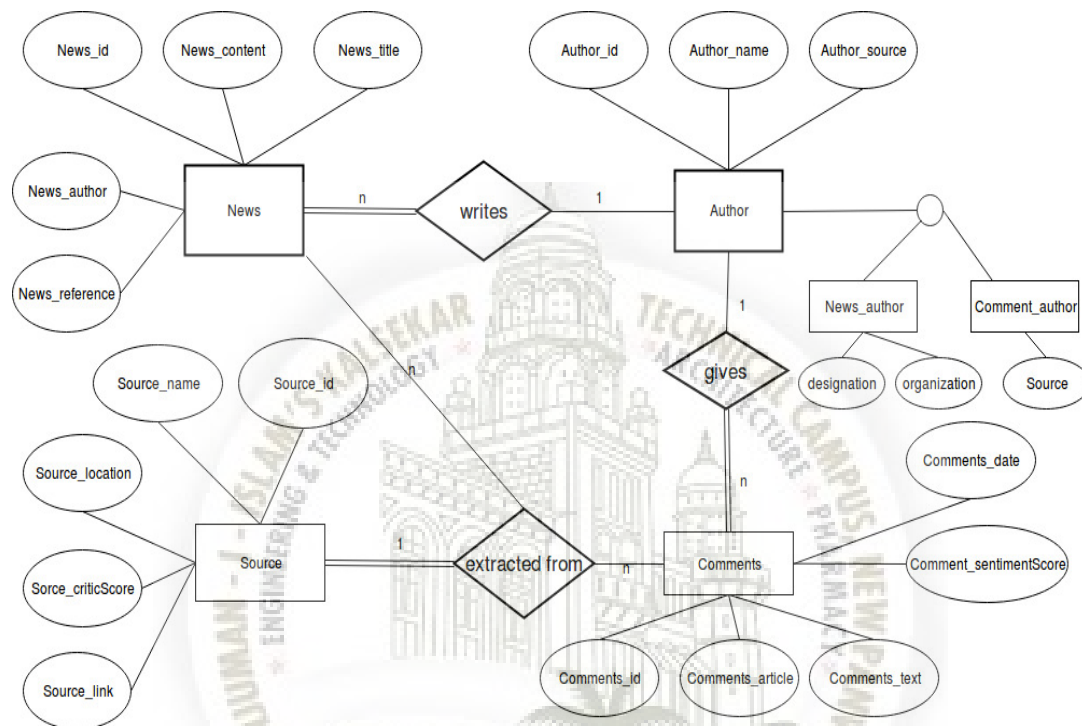


Figure 5.5: Entity Relationship for News Validation System

This is the E-R diagram of the system in which the modules which will be there after the deployment are shown. It is shown in a very easy way to have a brief overview of project.

## 5.2 System Architecture Design

The system architecture is mainly divided in two sub modules, the first sub module consists of client side and the main module is news validation server.

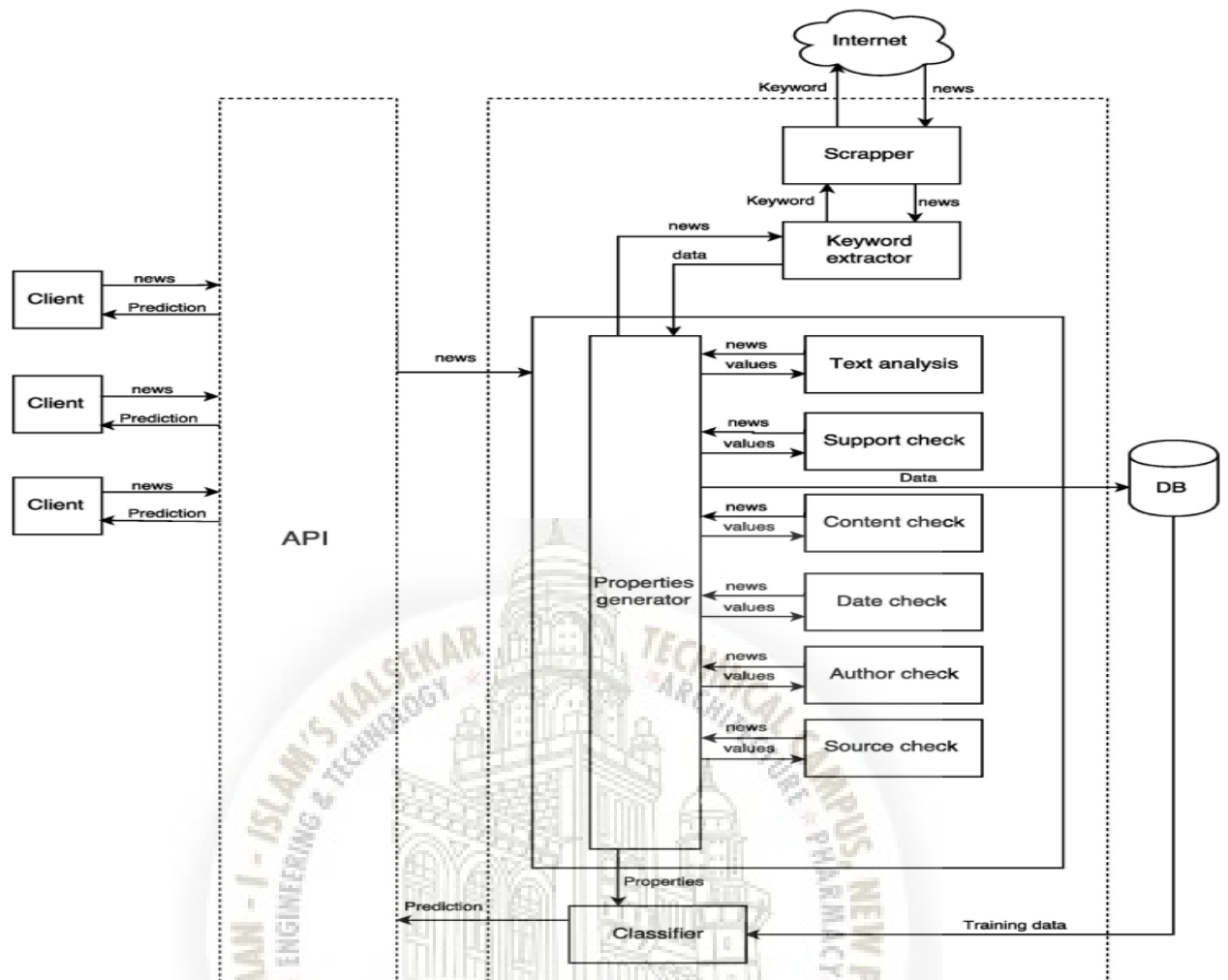


Figure 5.6: System Architecture for News Validation System

## 5.3 Sub-system Development

There are total four main modules in system architecture namely scrapper,keyword extractor,properties generator and classifier. scrapper will scrap the news from various news website using api ,keyword extractor will extract the keywords from the scrap data to process it,properties generator will generate a properties of news based on the extracted keywords,based on properties generated from properties generator and classifier the system will predict the prediction.

### 5.3.1 Authentic Source Identification

Data from authentic and authorized source plays very important role in news validation. Authenticity of source and author of training data set can increase the accuracy of results. This technique will be implemented in Source and author check sub-module in validator module of news validation system.

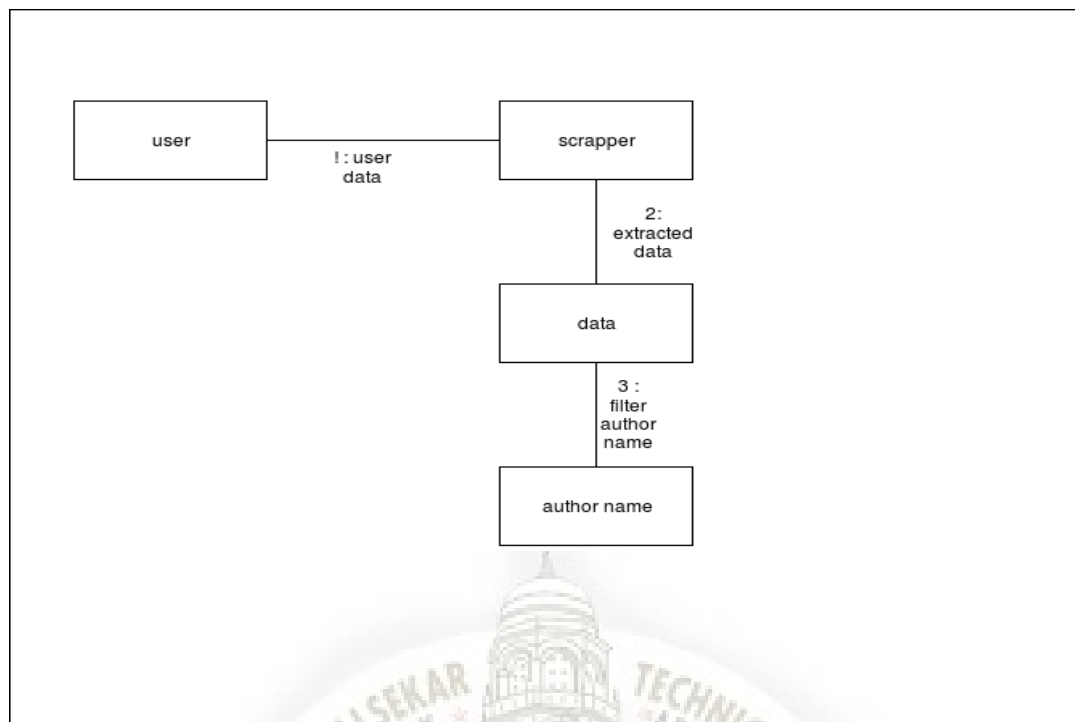


Figure 5.7: Modular diagram for authentic source identification

### 5.3.2 Date Validation

Sometimes it happens like old news get recreated and shared intentionally or some old news is shared without checking its date, so to tackle this problem date validation is very important. there are chances that some old news is updated by positive results but not mentioned in news, hence it is very important to check date. This technique will be implemented in date check sub-module in validator module of news validation system.

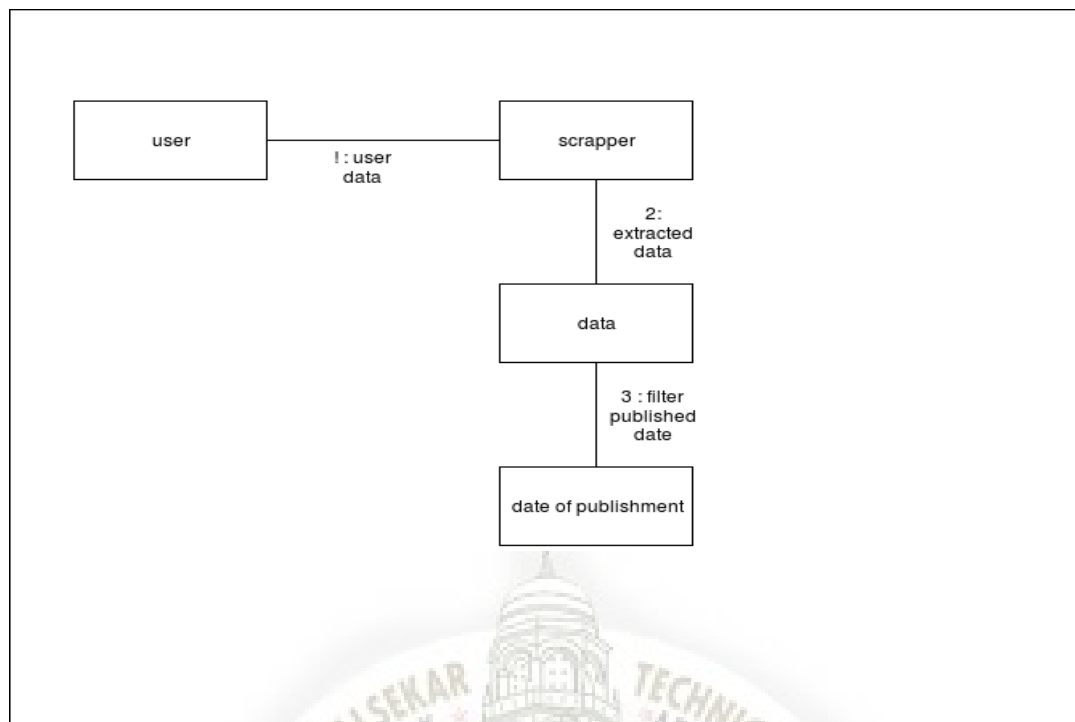


Figure 5.8: Modular diagram for date validation

### 5.3.3 Data collection and comparison

Collection of articles from authentic and genuine source and analyzing that plays a very important role in news validation. Data is collected based on keyword extracted from heading and paragraph of target news article and comparison is done between scrapped articles and target news articles. Matching rate of both articles is inversely proportional to chances of news to be fake, more match rate can result to less chance of news article to be fake. This technique will be implemented in content check and pattern match sub-module in validator module of news validation system. Grammatical leakage is also analyzed in this module.

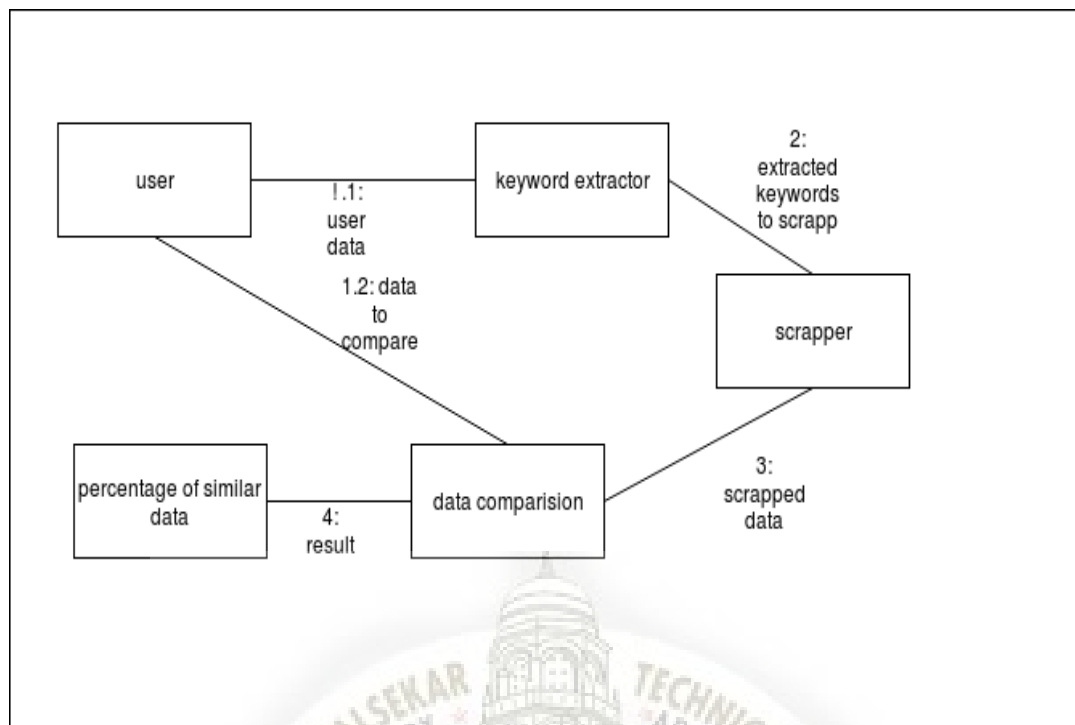


Figure 5.9: Modular diagram for data collection and comparison

### 5.3.4 Support Analysis

Reaction of users on a news also plays an important role in defining its validity. Reaction will be analyzed using sentimental analysis on comments of users on related news articles fetched from genuine news websites. Design is the first step in the development phase for any techniques and principles for the purpose of defining a device, a process or system in sufficient detail to permit its physical realization. Once the software requirements have been analyzed and specified the software design, coding, implementation and testing that are required to build and verify the software.

The design activities are of main importance in this phase, because in this activity decision ultimately affecting the success of the software implementation and its ease of maintenance are made. These decisions have the final bearing upon reliability and maintainability of the system. Design is the place where quality is fostered in development. Software design is the process through which requirements are translated into a representation of software. Software design is conducted in two steps.

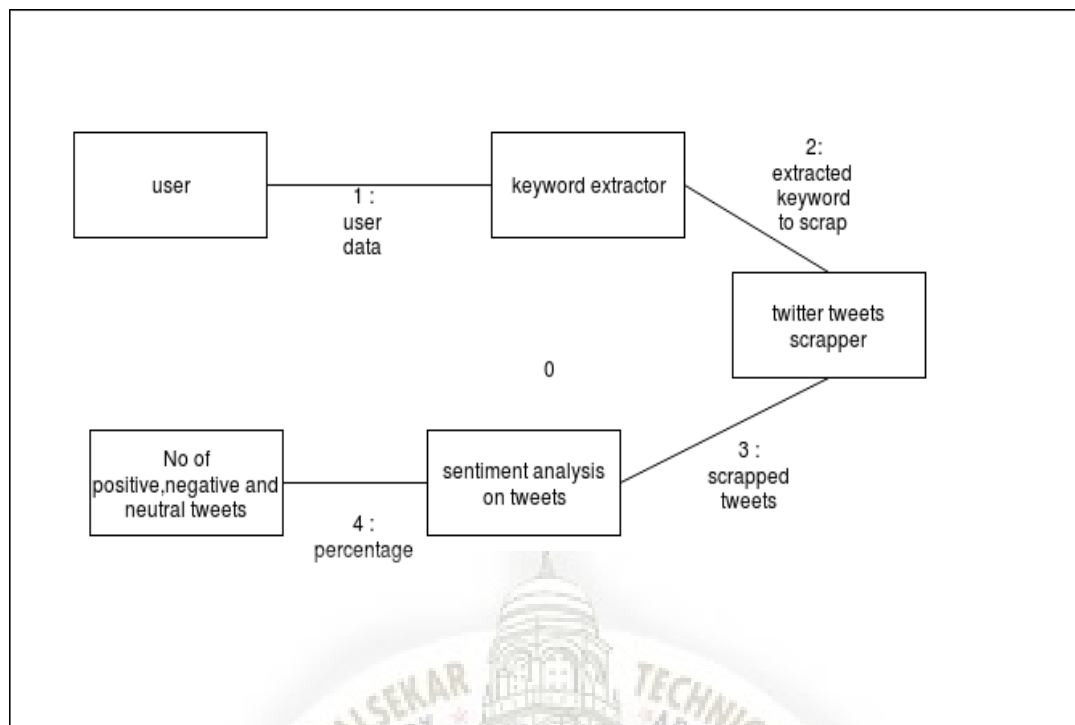


Figure 5.10: Modular diagram for support analysis

## 5.4 Systems Integration

There are mainly five main module in our system. Scrapper this module extract the data from the various news sources, Keyword extractor this module will extract the keyword from the given input data, properties generator it will generate the properties depends of the given data, classifier in this module by using the data from properties generator and training data the system will predict whether the news is correct or not.

### 5.4.1 Class Diagram

This is the Class diagram of the system in which the modules which will be there after the deployment are shown. This class diagram is an illustration of the relationships and source code dependencies among classes in the Unified Modeling Language (UML). In this context, a class defines the methods and variables in an object, which is a specific entity in a program or the unit of code representing that entity.

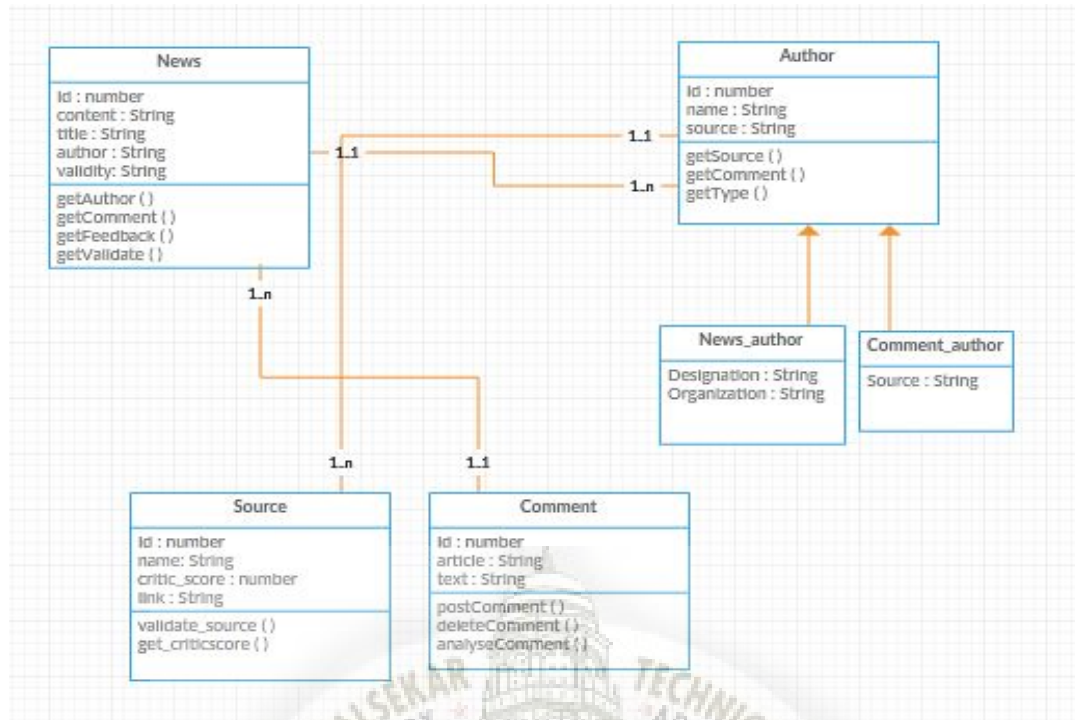


Figure 5.11: Class diagram for News Validation System

## 5.4.2 Sequence Diagram

To understand what a sequence diagram is, it's important to know the role of UML. UML, or the Unified Modeling Language, is a modeling toolkit that guides the creation and notation of many types of diagrams, including behavior diagrams, interaction diagrams, and structure diagrams. Sequence diagrams are a kind of interaction diagram, because they describe how and in what order a group of objects works together. These diagrams are used by software developers and business people alike to understand requirements for a new system or to document an existing process. Sequence diagrams are sometimes known as event diagrams or event scenarios.



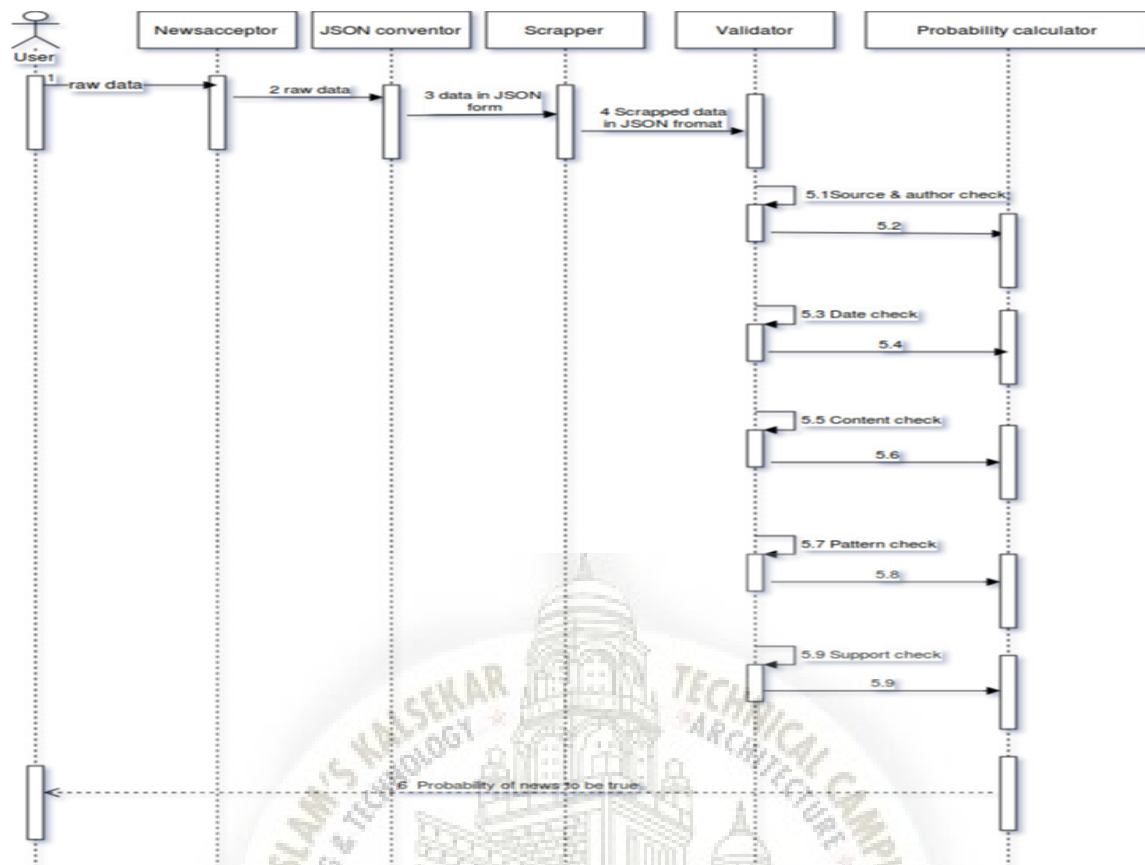


Figure 5.12: Sequence Diagram for News Validation System

### 5.4.3 Component Diagram

Component diagrams are different in terms of nature and behavior. Component diagrams are used to model the physical aspects of a system. Now the question is, what are these physical aspects? Physical aspects are the elements such as executables, libraries, files, documents, etc. which reside in a node. Component diagrams are used to visualize the organization and relationships among components in a system. These diagrams are also used to make executable systems.

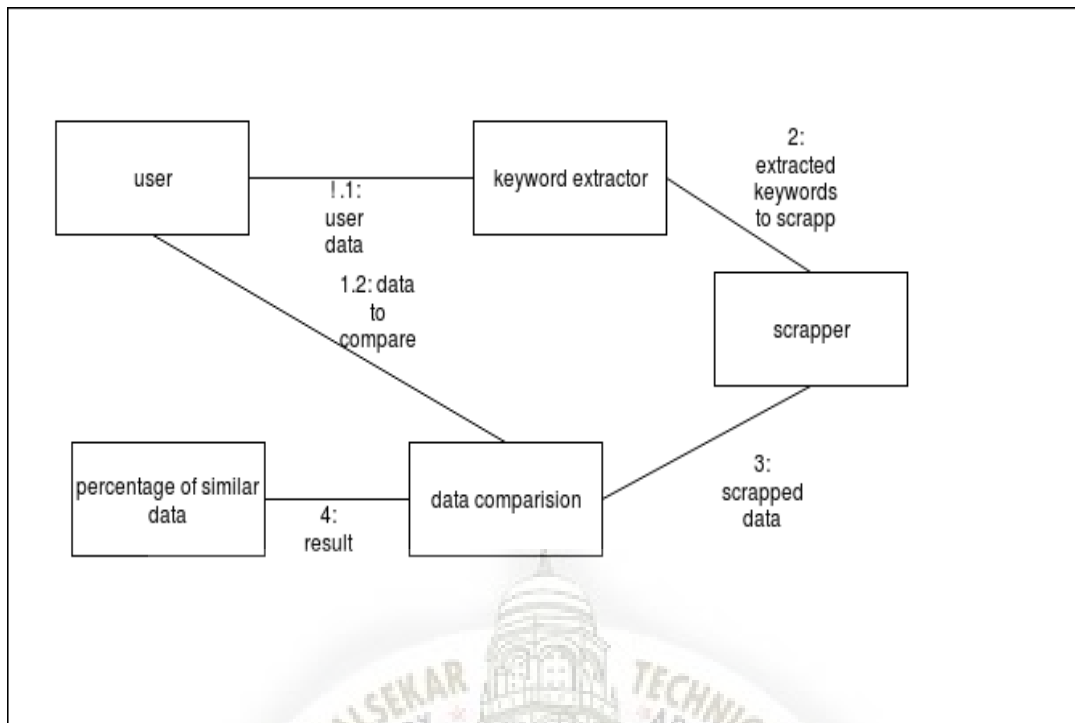


Figure 5.13: Component diagram for News Validation System

### 5.4.4 Deployment Diagram

Deployment diagram is a structure diagram which shows architecture of the system as deployment (distribution) of software artifacts to deployment targets. The below figure describe the Deployment Diagram of our system.

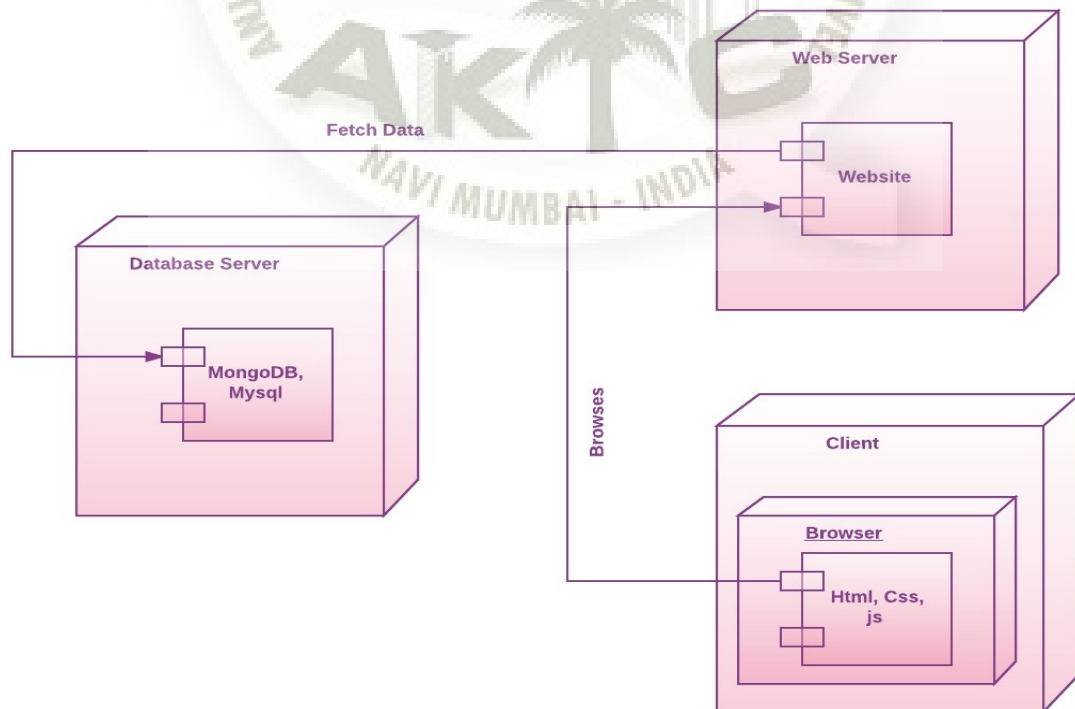


Figure 5.14: Deployment diagram for News Validation System

# Chapter 6

## Implementation

### 6.1 Text Analysis Module

This module will analyse the text which is given input to it. It will check the spelling of input data ,grammar of the input data and also punctuation of the input data.

#### Grammer Check

```

1 """
2 this sub-module checks for grammar errors.
3 """
4 import sys
5 import urllib.parse
6 import urllib.request
7 from urllib.error import HTTPError
8 from urllib.error import URLError
9 import json
10
11 def get_ginger_url(text):
12     """Get URL for checking grammar using Ginger.
13     @param text English text
14     @return URL
15     """
16     API_KEY = "6ae0c3a0-afdc-4532-a810-82ded0054236"
17
18     scheme = "http"
19     netloc = "services.gingersoftware.com"
20     path = "/Ginger/correct/json/GingerTheText"
21     params = ""
22     query = urllib.parse.urlencode([
23         ("lang", "US"),
24         ("clientVersion", "2.0"),
25         ("apiKey", API_KEY),
26         ("text", text)])
27     fragment = ""
28
29     return(urllib.parse.urlunparse((scheme, netloc, path, params, query,
30         fragment)))
31
32 def get_ginger_result(text):
33     """Get a result of checking grammar.
34     @param text English text
35     @return result of grammar check by Ginger

```

```

35 """
36 url = get_ginger_url(text)
37
38 try:
39     response = urllib.request.urlopen(url)
40 except HTTPError as e:
41     print("HTTP Error:", e.code)
42     quit()
43 except URLError as e:
44     print("URL Error:", e.reason)
45     quit()
46
47 try:
48     result = json.loads(response.read().decode('utf-8'))
49 except ValueError:
50     print("Value Error: Invalid server response.")
51     quit()
52
53 return(result)
54
55 get_ginger_result("this sub-module checks for gramar errrs???)

```

```

1 """
2 this sub module checks for misspelled words in article
3 """
4 from nltk import word_tokenize
5 import enchant
6 import re
7 article = "this sub module checks for misspelled words in article"
8
9 d = enchant.Dict("en_US")
10
11 non_dict_words = list(set([word.encode('ascii', 'ignore') for word in
12     word_tokenize(article) if d.check(word) is False and re.match('[a-zA-Z ]*$',
13     ,word)] ))
14
15 non_dict_words

```

```

1 """
2 this sub module will check sentiments of text in between dialogues.
3 """
4 import re
5 from textblob import TextBlob
6
7 content = 'he exclaimed "what the FUCK"'
8
9 text = re.findall('"([\^"]*)"', content)
10 blob = TextBlob(text[0])
11
12 blob.tags
13
14 for sentence in blob.sentences:
15     print(sentence)
16     print(sentence.sentiment.polarity)

```

## 6.2 Source, Author , Date Check Module

This module will extract the date of publish, name of author, name of source of the given news from different news website.

```

1  '''This module will give author,source and date of the data
2  '''
3  # -*- coding: utf-8 -*-
4  print("Initializing...:")
5
6  import requests
7  from goose import Goose
8  from nltk import word_tokenize , pos_tag
9  from nltk.corpus import wordnet as wn
10 from rake_nltk import Rake
11
12
13 authenticate_sources = ["NewsExpress"]
14 # semantic similarity definitions
15
16
17 def penn_to_wn(tag):
18     """ Convert between a Penn Treebank tag to a simplified Wordnet tag """
19     if tag.startswith('N'):
20         return 'n'
21
22     if tag.startswith('V'):
23         return 'v'
24
25     if tag.startswith('J'):
26         return 'a'
27
28     if tag.startswith('R'):
29         return 'r'
30
31     return None
32
33
34 def tagged_to_synset(word, tag):
35     wn_tag = penn_to_wn(tag)
36     if wn_tag is None:
37         return None
38
39     try:
40         return wn.synsets(word, wn_tag)[0]
41     except BaseException:
42         return None
43
44
45 def sentence_similarity(sentence1 , sentence2):
46     """ compute the sentence similarity using Wordnet """
47     # Tokenize and tag
48     sentence1 = pos_tag(word_tokenize(sentence1))
49     sentence2 = pos_tag(word_tokenize(sentence2))
50
51     # Get the synsets for the tagged words
52     synsets1 = [tagged_to_synset(*tagged_word) for tagged_word in sentence1]
53     synsets2 = [tagged_to_synset(*tagged_word) for tagged_word in sentence2]
54

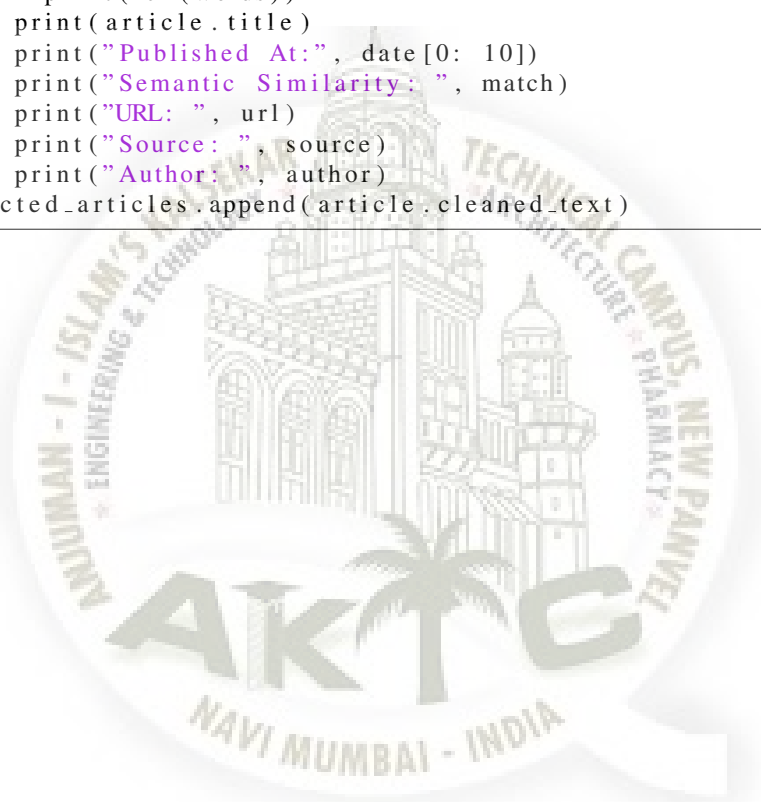
```

```

55     # Filter out the Nones
56     synsets1 = [ss for ss in synsets1 if ss]
57     synsets2 = [ss for ss in synsets2 if ss]
58
59     score, count = 0.0, 0
60
61     # For each word in the first sentence
62     for synset in synsets1:
63         # Get the similarity value of the most similar word in the other
64         # sentence
65         best_score = max([synset.path_similarity(ss) for ss in synsets2])
66
67         # Check that the similarity could have been computed
68         if best_score is not None:
69             score += best_score
70             count += 1
71
72     # Average the values
73     score /= count
74     return score
75
76 # taking input from user
77
78
79 input_news_content = "Late last year, lawyers for President Donald Trump
    expressed optimism that special counsel Robert Mueller was nearing the end
    of his probe of Russia's interference in the 2016 election. But if there was
    hope in the White House that Trump might be moving past an investigation
    that has dogged his presidency from the start, 2018 is beginning without
    signs of abatement. In fact, the new year set off a flurry of developments
    in the probes by Mueller and Congress"
80 #input_news_content = "Forensic doctors in Dubai concluded that Sridevi died of
    a heart attack and added there is nothing suspicious about the way the
    superstar passed away , official sources in Dubai said."
81 #input_news_content = "Shwag apologises for his tweet that gave Kerala lynching
    a communal colour"
82 #input_news_content = "Scientific research ascertains mercury toxicity but
    Sadhguru continues to endorse it for Indian traditional medicines"
83 rake = Rake()
84 rake.extract_keywords_from_text(input_news_content)
85 keyword = rake.get_ranked_phrases()
86 print("Extracting Keywords")
87 print(keyword)
88
89 # newsapi starting link declaration
90 news_api_link = "https://newsapi.org/v2/everything?q="
91
92 news_api_query = keyword[0]
93 print("Getting news data:")
94
95 news_data_request = requests.get(
96     news_api_link +
97     input_news_content +
98     "&sortBy=publishedAt&apiKey=de67235049564f0f8d206f8aff2476a1")
99 news_api_data = news_data_request.json()
100
101 if len(news_api_data['articles']) == 0:
102     print("No Authenticate news sources found.")
103 else:
104
105     g = Goose()

```

```
106 match = 0
107 url = ""
108 source = ""
109 for i in range(0, len(news_api_data['articles'])):
110     url = news_api_data['articles'][i]['url']
111     article = g.extract(url=url)
112     article_text = article.cleaned_text
113     if article_text != '':
114         similarity = sentence_similarity(article_text, input_news_content)
115         if match < similarity:
116             match = similarity
117             matched_news = article.title
118             date = str(news_api_data['articles'][i]['publishedAt'])
119             url = news_api_data['articles'][i]['url']
120             source = news_api_data['articles'][i]['source']['name']
121             author = news_api_data['articles'][i]['author']
122             # print(len(words))
123             print(article.title)
124             print("Published At:", date[0: 10])
125             print("Semantic Similarity:", match)
126             print("URL: ", url)
127             print("Source: ", source)
128             print("Author: ", author)
129 # extracted_articles.append(article.cleaned_text)
```





### 6.3 Support Check Module.

This Module will extract the twitter tweets of verified user on particular topic and return the numbers of positive tweets, negative tweets and neutral tweets.

```

1 '''
2 python3 tweep.py
3
4 Twitter sentiment analysis
5
6 '''
7
8 #!/usr/bin/python3
9 from bs4 import BeautifulSoup
10 from time import gmtime, strftime
11 import argparse
12 import aiohttp
13 import asyncio
14 import async_timeout
15 import csv
16 import datetime
17 import json
18 import re
19 import sys
20 import json
21 from textblob import TextBlob
22
23
24 async def getUrl(init):
25     '''
26     URL Descision:
27     Tweep utilizes positions of Tweet's from Twitter's search feature to
28     iterate through a user's Twitter feed. This section decides whether
29     this is the first URL request or not and develops the URL based on the
30     args given.
31
32     Returns complete URL.
33
34     Todo: Make everything URL encoded at the end.
35     '''
36     if init == -1:
37         url = "https://twitter.com/search?f=tweets&vertical=default&lang=en&q="
38     else:
39         url = "https://twitter.com/i/search/timeline?f=tweets&vertical=default"
40         url+= "&lang=en&include_available_features=1&include_entities=1&reset_"
41         url+= "error_state=false&src=typd&max_position={}&q=".format(init)
42
43     if arg.u != None:
44         url+= "from%3A{0.u}".format(arg)
45     if arg.g != None:
46         arg.g = arg.g.replace(" ", "")
47         url+= "geocode%3A{0.g}".format(arg)
48     if arg.s != None:
49         arg.s = arg.s.replace(" ", "%20").replace("#", "%23")
50         url+= "%20{0.s}".format(arg)
51     if arg.year != None:
52         url+= "%20until%3A{0.year}-1-1".format(arg)
53     if arg.since != None:
54         url+= "%20since%3A{0.since}".format(arg)

```

```

55     if arg.fruit:
56         url+= "%20myspace.com%20OR%20last.fm%20OR"
57         url+= "%20mail%20OR%20email%20OR%20gmail%20OR%20e-mail"
58         url+= "%20OR%20phone%20OR%20call%20me%20OR%20text%20me"
59         url+= "%20OR%20keybase"
60     if arg.verified:
61         url+= "%20filter%3Averified"
62
63     return url
64
65 async def fetch(session, url):
66     """
67     Basic aiohttp request with a 30 second timeout.
68     """
69     with async_timeout.timeout(30):
70         async with session.get(url) as response:
71             return await response.text()
72
73 async def initial(response):
74     """
75     Initial response parsing and collecting the position ID
76     """
77     soup = BeautifulSoup(response, "html.parser")
78     feed = soup.find_all("li", "js-stream-item")
79     init = "TWEET-{}-{}".format(feed[-1]["data-item-id"], feed[0]["data-item-id"]
80                                ")")
81
82     return feed, init
83
84 async def cont(response):
85     """
86     Regular json response parsing and collecting Position ID
87     """
88     json_response = json.loads(response)
89     html = json_response["items_html"]
90     soup = BeautifulSoup(html, "html.parser")
91     feed = soup.find_all("li", "js-stream-item")
92     split = json_response["min_position"].split("-")
93     split[1] = feed[-1]["data-item-id"]
94     init = "-".join(split)
95
96     return feed, init
97
98 async def getFeed(init):
99     """
100     Parsing Descision:
101     Responses from requests with the position id's are JSON,
102     so this section decides whether this is an initial request
103     or not to use the appropriate response reading for parsing
104     with BeautifulSoup4.
105
106     Returns html for Tweets and position id.
107     """
108     async with aiohttp.ClientSession() as session:
109         response = await fetch(session, await getUrl(init))
110         feed = []
111         try:
112             if init == -1:
113                 feed, init = await initial(response)
114             else:
115                 feed, init = await cont(response)

```

```

115     except:
116         # Tweep will realize that it's done scraping.
117         pass
118
119     return feed, init
120
121 def outTweet(tweet, dat):
122     '''
123     Parsing Section:
124     This function will create the desired output string and
125     write it to a file or csv if specified.
126
127     Returns output.
128     '''
129     #print("hello ", dat)
130     tweetid = tweet["data-item-id"]
131     global lst
132     # Formatting the date & time stamps just how I like it.
133     datestamp = tweet.find("a", "tweet-timestamp")["title"].rpartition(" - ")[-1]
134     d = datetime.datetime.strptime(datestamp, "%d %b %Y")
135     date = d.strftime("%Y-%m-%d")
136     timestamp = str(datetime.timedelta(seconds=int(tweet.find("span", "_timestamp")["data-time"])).rpartition(", ")[-1])
137     t = datetime.datetime.strptime(timestamp, "%H:%M:%S")
138     time = t.strftime("%H:%M:%S")
139     # The @ in the username annoys me.
140     username = tweet.find("span", "username").text.replace("@", "")
141     timezone = strftime("%Z", gmtime())
142     # The context of the Tweet compressed into a single line.
143     text = tweet.find("p", "tweet-text").text.replace("\n", "").replace("http",
144         " http").replace("pic.twitter", " pic.twitter")
145     # Regex for gathering hashtags
146     hashtags = ", ".join(re.findall(r'(?i)\#\w+', text, flags=re.UNICODE))
147     replies = tweet.find("span", "ProfileTweet-action--reply u-hiddenVisually").find("span")["data-tweet-stat-count"]
148     retweets = tweet.find("span", "ProfileTweet-action--retweet u-hiddenVisually").find("span")["data-tweet-stat-count"]
149     likes = tweet.find("span", "ProfileTweet-action--favorite u-hiddenVisually").find("span")["data-tweet-stat-count"]
150     '''
151     This part tries to get a list of mentions.
152     It sometimes gets slow with Tweets that contain
153     40+ mentioned people.. rather than just appending
154     the whole list to the Tweet, it goes through each
155     one to make sure there aren't any duplicates.
156     '''
157     try:
158         mentions = tweet.find("div", "js-original-tweet")["data-mentions"].split(" ")
159         for i in range(len(mentions)):
160             mention = "@{}".format(mentions[i])
161             if mention not in text:
162                 text = "{}".format(text)
163     except:
164         pass
165
166     # Preparing to output
167     '''
168     There were certain cases where I used Tweep

```

```

169     to gather a list of users and then fed that
170     generated list into Tweep. That's why these
171     modes exist.
172     '''
173     if arg.users:
174         output = username
175     elif arg.tweets:
176         output = tweets
177     else:
178         '''
179         The standard output is how I like it, although
180         this can be modified to your desire. Uncomment
181         the bottom line and add in the variables in the
182         order you want them or how you want it to look.
183         '''
184         output = ""
185         output = "{}".format(text)
186
187
188     if arg.o == None:
189         dat=[text]
190
191     return dat
192
193 async def getTweets(init, var):
194     '''
195     This function uses the html responses from getFeed()
196     and sends that info to the Tweet parser outTweet() and
197     outputs it.
198
199     Returns response feed, if it's first-run, and Tweet count.
200     '''
201     tweets, init = await getFeed(init)
202     count = 0
203     dat=[]
204     for tweet in tweets:
205         '''
206         Certain Tweets get taken down for copyright but are still
207         visible in the search. We want to avoid those.
208         '''
209         copyright = tweet.find("div", "StreamItemContent—withheld")
210         if copyright is None:
211             count +=1
212             var.append(outTweet(tweet, dat))
213             if(count==20):
214                 break
215             # print(var)
216     return count, var
217
218 async def getUsername():
219     '''
220     This function uses a Twitter ID search to resolve a Twitter User
221     ID and return it's corresponding username.
222     '''
223     async with aiohttp.ClientSession() as session:
224         r = await fetch(session, "https://twitter.com/intent/user?user_id={0.
225             userid}".format(arg))
226         soup = BeautifulSoup(r, "html.parser")
227         return soup.find("a", "fn_url alternate-context")["href"].replace("/", "")
228
229 async def main():

```

```

229     '''
230     Putting it all together.
231     '''
232     if arg.userid is not None:
233         arg.u = await getUsername()
234
235     feed = [-1]
236     init = -1
237     num = 0
238     var=[]
239     neu=pos=neg=0
240     tweets_senti=[]
241     if len(feed) > 0:
242         num,var = await getTweets(init , var)
243         for i in range (num):
244             tweets=""
245             tweets="" .join (var [i])
246             tweets1=tweets .split ('http')
247             #print (tweets1 [0])
248             tweets_senti .append ([ TextBlob (tweets1 [0]) .polarity ])
249             if (TextBlob (tweets1 [0]) .polarity >0):
250                 pos+=1
251             elif (TextBlob (tweets1 [0]) .polarity <0):
252                 neg+=1
253             else:
254                 neu+=1
255             #print (tweets_senti)
256             print ("positive tweets ",pos," negative tweets ",neg," neutral
                tweets ",neu)
257
258
259 def Error(error , message):
260     # Error formatting
261     print ("[-] {}: {}".format (error , message))
262     sys.exit (0)
263
264 def check():
265     # Performs main argument checks so nothing unintended happens.
266     if arg.u is not None:
267         if arg.users:
268             Error ("Contradicting Args", "Please use --users in combination with
                -s.")
269         if arg.verified:
270             Error ("Contradicting Args", "Please use --verified in combination
                with -s.")
271         if arg.userid:
272             Error ("Contradicting Args", "--userid and -u cannot be used together
                .")
273     if arg.tweets and arg.users:
274         Error ("Contradicting Args", "--users and --tweets cannot be used
                together.")
275     if arg.csv and arg.o is None:
276         Error ("Error", "Please specify an output file (Example: -o file.csv)")
277     if arg.u is None and arg.s is None and arg.userid is None and arg.g is None:
278         Error ("Error", "Please specify a username, user id, search or geotag.")
279
280 def check_support (user_data , limit , verified):
281
282 #if __name__ == "__main__":
283     ap = argparse.ArgumentParser (prog="tweep.py", usage="python3 %(prog)s [
        options]", description="tweep.py - An Advanced Twitter Scraping Tool")

```

```
284 ap.add_argument("-u", help="User's Tweets you want to scrape.")
285 ap.add_argument("-s", help="Search for Tweets containing this word or phrase
    .", default=user_data)
286 ap.add_argument("-o", help="Save output to a file.")
287 ap.add_argument("-g", help="Search for geocoded tweets.")
288 ap.add_argument("--year", help="Filter Tweets before specified year.")
289 ap.add_argument("--since", help="Filter Tweets sent since date (Example:
    2017-12-27).")
290 ap.add_argument("--fruit", help="Display 'low-hanging-fruit' Tweets.",
    action="store_true")
291 ap.add_argument("--tweets", help="Display Tweets only.", action="store_true"
    )
292 ap.add_argument("--verified", help="Display Tweets only from verified users
    (Use with -s).", action="store_true", default=verified)
293 ap.add_argument("--users", help="Display users only (Use with -s).", action=
    "store_true")
294 ap.add_argument("--csv", help="Write as .csv file.", action="store_true")
295 ap.add_argument("--hashtags", help="Output hashtags in seperate column.",
    action="store_true")
296 ap.add_argument("--userid", help="Twitter user id")
297 ap.add_argument("--limit", help="Number of Tweets to pull (Increments of 20)
    .", default=limit)
298 ap.add_argument("--count", help="Display number Tweets scraped at the end of
    session.", action="store_true")
299 ap.add_argument("--stats", help="Show number of replies, retweets, and likes
    ", action="store_true")
300 global arg
301 #global var
302 arg = ap.parse_args()
303 #print(arg.verified)
304 check()
305
306 loop = asyncio.get_event_loop()
307 loop.run_until_complete(main())
308 polar=[]
309
310 check_support('Salman khan after black buck case',20,True)
```

## 6.4 Content Check Module

This module will basically check the content of the given data input.

```

1 '''
2 Content Check Module
3 '''
4 from textblob import TextBlob
5 from gingerit.gingerit import GingerIt
6 import re
7
8 class Linguistic:
9     def __init__(self):
10         self.grammar_mistakes_count = 0
11         self.content_sentiment = {}
12         self.quote_sentiment = {}
13         self.caps_sentiment = {}
14
15     def grammar_mistake(self, content):
16         parser = GingerIt()
17         result = parser.parse(content)
18         return len(result['corrections'])
19
20     def get_quotes_sentiment(self, content):
21         quotes = re.findall(r'"([\^"]*)"', content)
22         quotes_string = ' '.join(quotes)
23
24         self.quote_sentiment = self.sentiment_analysis(quotes_string)
25
26     def get_caps_sentiment(self, content):
27         caps = re.findall('\w+[A-Z]', content)
28         caps_string = ' '.join(caps)
29         #
30         self.caps_string_length=len(' '.join(caps))
31         #
32         self.caps_sentiment = self.sentiment_analysis(caps_string)
33
34     def sentiment_analysis(self, content):
35         blob = TextBlob(content)
36         sentiments = {'negative': 0, 'positive': 0, 'neutral': 0}
37         for sentence in blob.sentences:
38             if sentence.sentiment.polarity < 0:
39                 sentiments['negative'] += 1
40             elif sentence.sentiment.polarity > 0:
41                 sentiments['positive'] += 1
42             elif sentence.sentiment.polarity == 0:
43                 sentiments['neutral'] += 1
44         return sentiments
45
46     def get_linguistic_features(self, input_news_content):
47
48         self.grammar_mistakes_count = self.grammar_mistake(input_news_content)
49         self.content_sentiment = self.sentiment_analysis(input_news_content)
50         self.get_quotes_sentiment(input_news_content)
51         self.get_caps_sentiment(input_news_content)
52
53         return [self.grammar_mistakes_count, self.caps_string_length, self.
54                 content_sentiment, self.quote_sentiment, self.caps_sentiment]

```



## 6.5 Classifier

In our system we have used two classifier ID3 classifier and KNN classifier.

```

1 '''
2 ID3
3 '''
4 import pandas as pd
5 import numpy as np
6
7 df = pd.read_csv('kyphosis.csv')
8
9 from sklearn.model_selection import train_test_split
10
11 X = df.drop('Kyphosis', axis=1)
12 y = df['Kyphosis']
13
14 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.30)
15
16 from sklearn.tree import DecisionTreeClassifier
17
18 dtree = DecisionTreeClassifier()
19
20 dtree.fit(X_train, y_train)
21
22 predictions = dtree.predict(X_test)
23
24 from sklearn.metrics import classification_report, confusion_matrix
25
26 print(classification_report(y_test, predictions))
27
28 print(confusion_matrix(y_test, predictions))
29
30
31 #random forest
32
33 from sklearn.ensemble import RandomForestClassifier
34
35 rfc = RandomForestClassifier(n_estimators=100)
36
37 rfc.fit(X_train, y_train)
38
39 rfc_pred = rfc.predict(X_test)
40
41 print(confusion_matrix(y_test, rfc_pred))
42
43 print(classification_report(y_test, rfc_pred))

```

```

1 '''
2 Knn
3 '''
4 import pandas as pd
5 import numpy as np
6 import matplotlib.pyplot as plt
7 import numpy as np
8
9 df = pd.read_csv('kyphosis.csv')
10

```

```

11 from sklearn.preprocessing import StandardScaler
12
13 scaler = StandardScaler()
14
15 scaler.fit(df.drop('Kyphosis',axis=1))
16
17 scaled_features = scaler.transform(df.drop('Kyphosis',axis=1))
18
19 df_feat = pd.DataFrame(scaled_features, columns=df.columns[:-1])
20
21 from sklearn.model_selection import train_test_split
22
23 X_train, X_test, y_train, y_test = train_test_split(scaled_features,df['Kyphosis
24     '],
25                                                     test_size=0.30)
26
27 from sklearn.neighbors import KNeighborsClassifier
28 knn = KNeighborsClassifier(n_neighbors=1)
29 knn.fit(X_train, y_train)
30 pred = knn.predict(X_test)
31
32 from sklearn.metrics import classification_report, confusion_matrix
33
34 print(confusion_matrix(y_test, pred))
35
36 print(classification_report(y_test, pred))
37
38 error_rate = []
39
40 # Will take some time
41 for i in range(1, 40):
42     knn = KNeighborsClassifier(n_neighbors=i)
43     knn.fit(X_train, y_train)
44     pred_i = knn.predict(X_test)
45     error_rate.append(np.mean(pred_i != y_test))
46
47 plt.figure(figsize=(10,6))
48 plt.plot(range(1,40), error_rate, color='blue', linestyle='dashed', marker='o',
49         markerfacecolor='red', markersize=10)
50 plt.title('Error Rate vs. K Value')
51 plt.xlabel('K')
52 plt.ylabel('Error Rate')
53 plt.show()
54 knn = KNeighborsClassifier(n_neighbors=10)
55
56 knn.fit(X_train, y_train)
57
58 print(type(X_test))
59 from numpy import array
60 #a = array([[1.46067113, -0.65203532, 1.34045062]])
61 #a.reshape(-1,1)
62 pred = knn.predict(X_test)
63 print(pred)
64 print('WITH K=23')
65 print('\n')
66 print(confusion_matrix(y_test, pred))
67 print('\n')
68 print(classification_report(y_test, pred))

```

# Chapter 7

## System Testing

The system testing is done as per as following criterias

### 7.1 Test Cases and Test Results

Test ID	Test Case Title	Test Condition	System Behavior	Expected Result
T01	Login	Should be registered user	User should login into our system	User should login into our system
T02	Search data	Data must be in correct form	Display the probability	Display the probability
T03	Search non existing data	Data must be in correct form	Output should be null	Output should be null

### 7.2 Sample of a Test Case

**Title:** Login Page – Authenticate Successfully on our system.

**Description:** A registered user should be able to successfully login at our system.

*Precondition:* The user must already be registered with an name, email address and password.

*Assumption:* A supported browser is being used.

**Test Steps:**

1. Navigate to register.
2. Register on our system
3. In the 'email' field, enter the email of the registered user.

4. Enter the password of the registered user
5. Click 'Login'

**Expected Result:** After the user has been login successfully the user will get its API key.

**Actual Result:** After user has login a random and unique key will be generated for further authentication of user. As expected the user will get its API key.

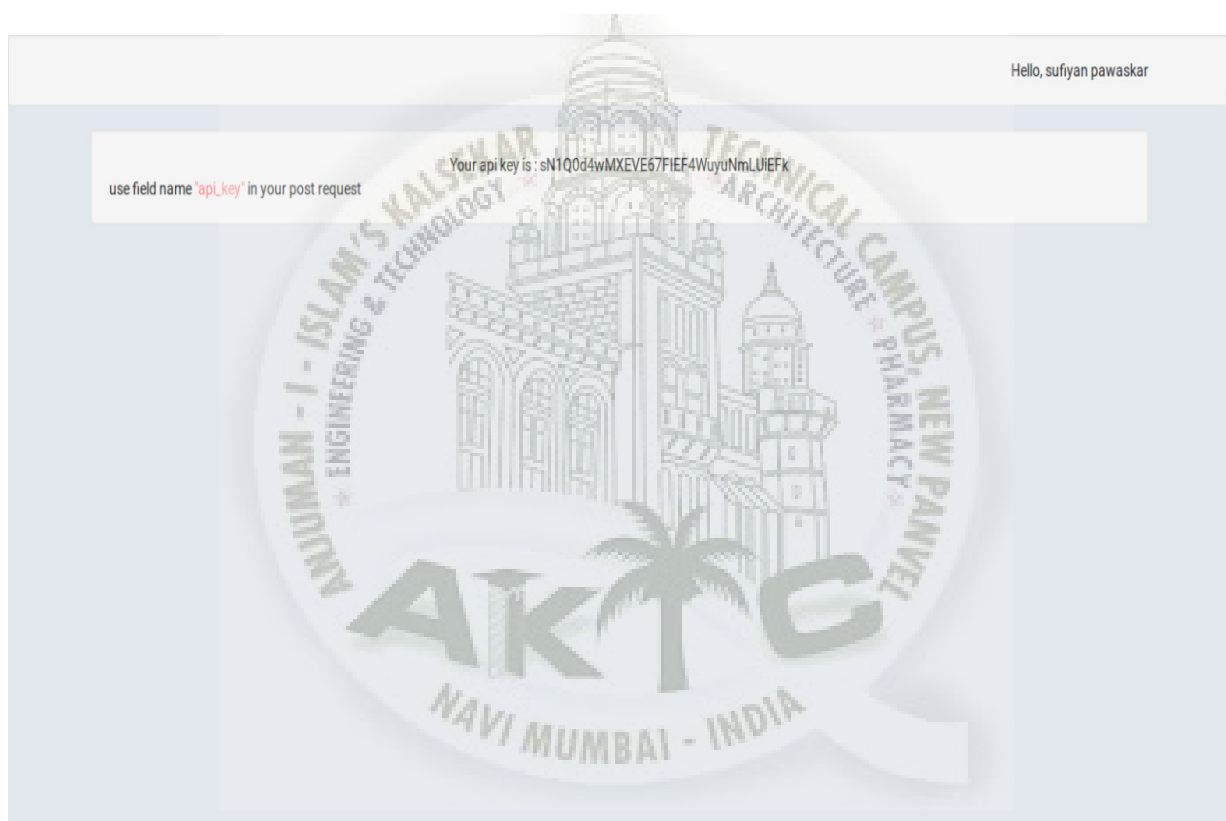


Figure 7.1: Key generation after login

### 7.2.1 Software Quality Attributes

1. Availability: The system should not be down, whenever the user use the system the specific data should be available to the user.
2. Correctness: As per as the user search correct data should be shown to user.

3. **Maintainability:**The administrator of the system should maintain the system.
4. **Reliability:** The system should be reliable for producing correct output so that user can reliable on system.
5. **Extensibility:** The system is capable to be modified by changing some modules or by adding some features to the existing system.



# Chapter 8

## Screenshots of Project

### 8.1 Registration of client



**Figure 8.1: Register**

## 8.2 Login of client

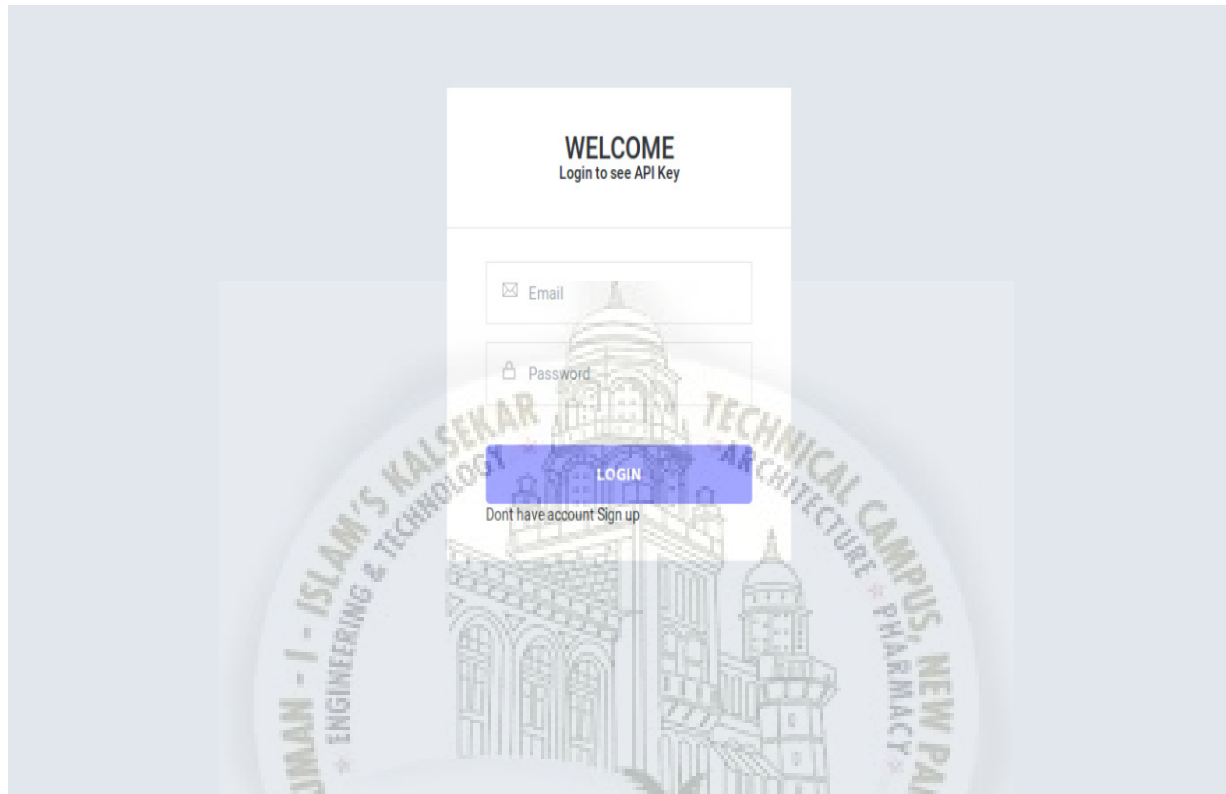


Figure 8.2: Login



### 8.3 Api Key Generation for client

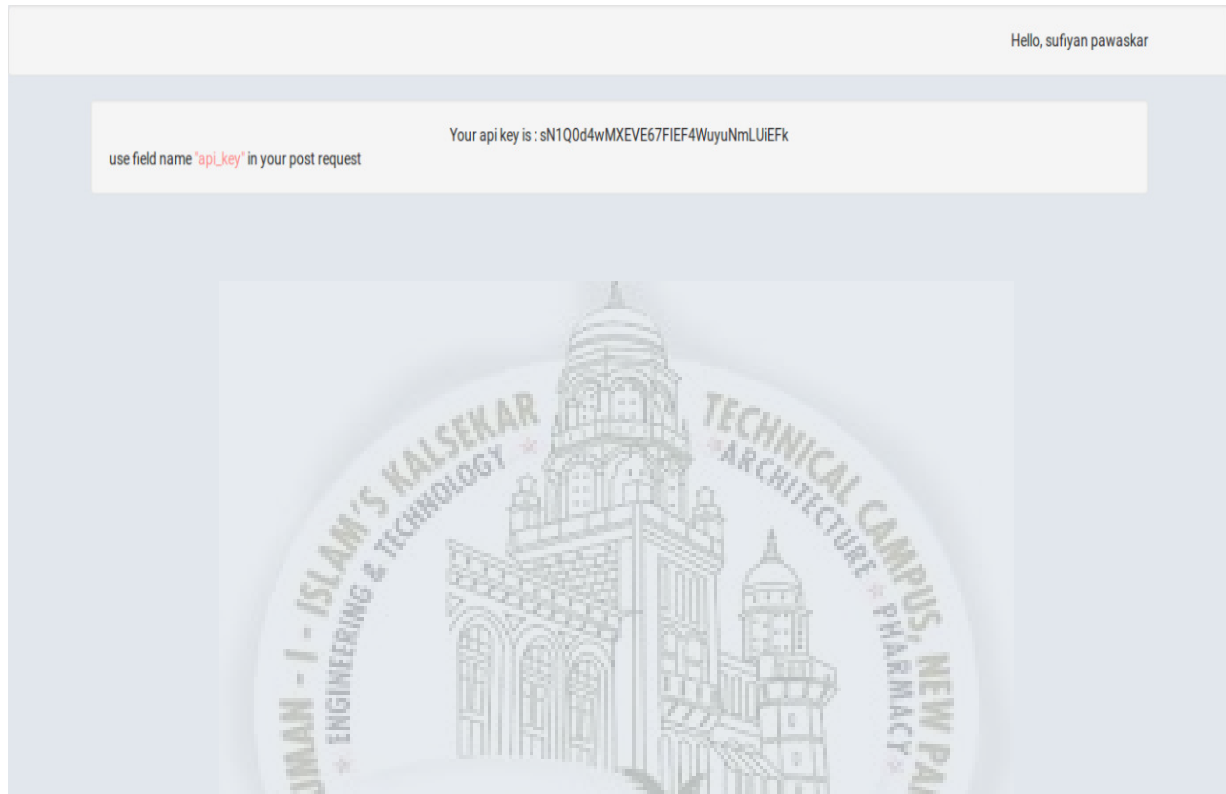
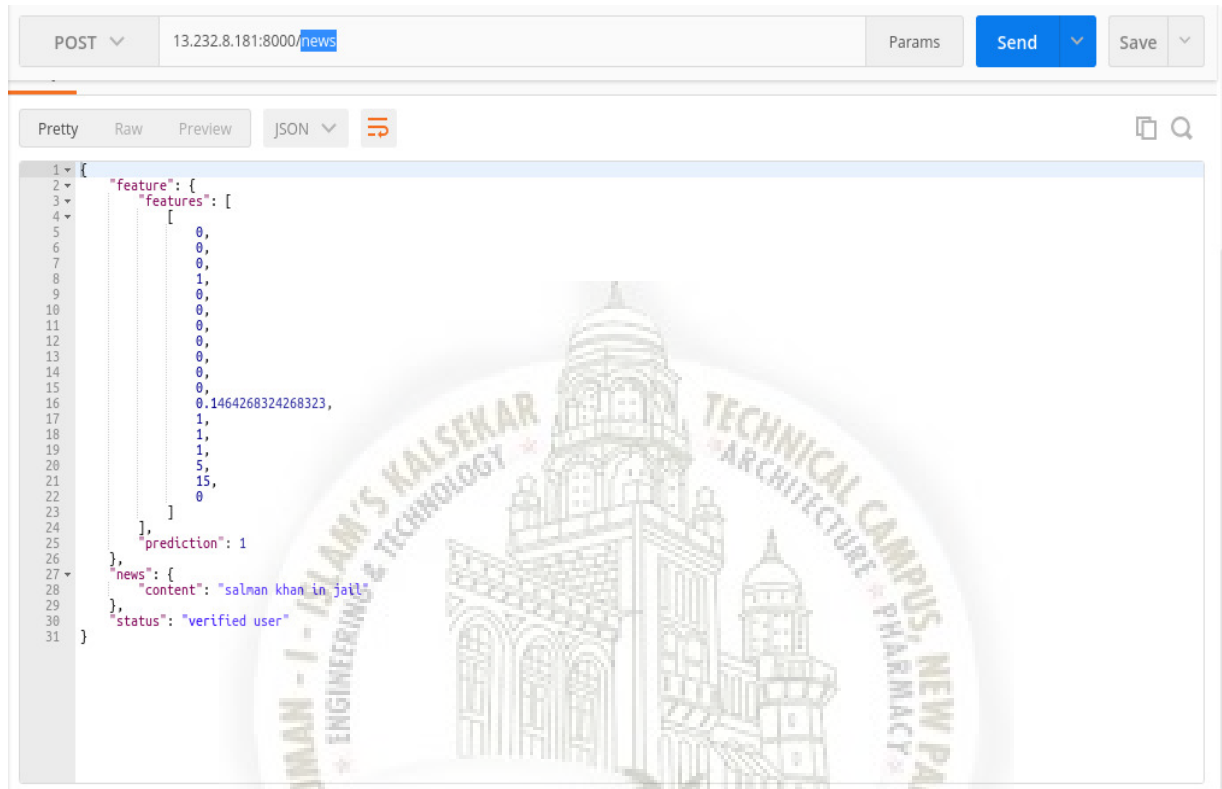


Figure 8.3: Key generation after login

## 8.4 Output for the search data



```
POST 13.232.8.181:8000/news Params Send Save

Pretty Raw Preview JSON

1 {
2   "feature": {
3     "features": [
4       [
5         0,
6         0,
7         0,
8         1,
9         0,
10        0,
11        0,
12        0,
13        0,
14        0,
15        0,
16        0.1464268324268323,
17        1,
18        1,
19        1,
20        5,
21        15,
22        0
23      ]
24    },
25    "prediction": 1
26  },
27  "news": {
28    "content": "salman khan in jail",
29  },
30  "status": "verified user"
31 }
```

Figure 8.4: Output

## Chapter 9

# Conclusion and Future Scope

### 9.1 Conclusion

Information in internet should be correct and reliable. But currently our society lacks a system that can check whether if some information implied in news is true or not. Developing such system will make a difference in society and will help to maintain peace.

### 9.2 Future Scope

Currently, We are targeting only global news, In future targeting the local news can be implemented.

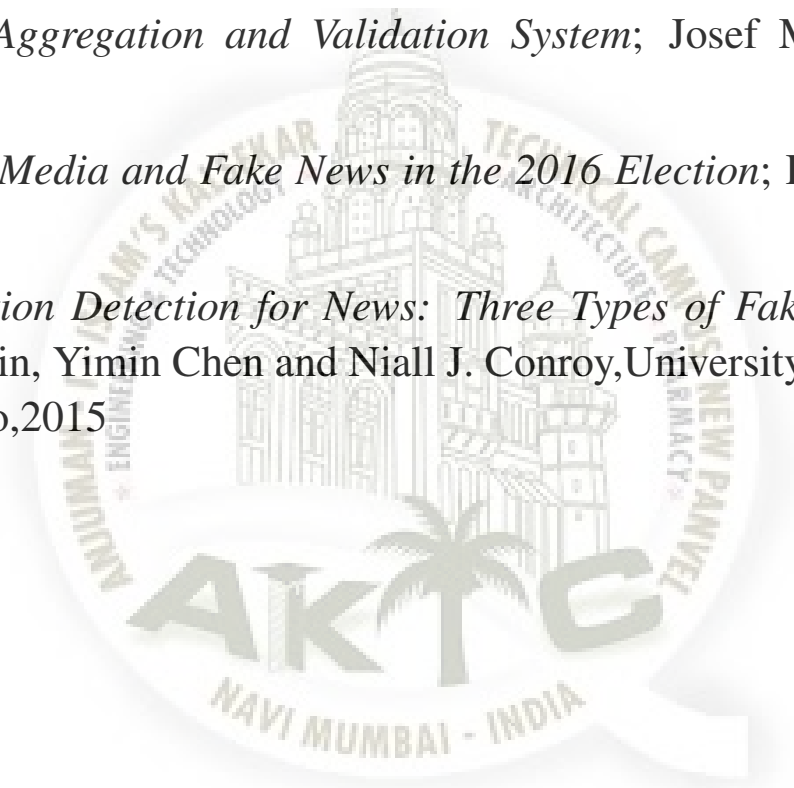
Also geographical area wise analysis can be implemented that would make the results more efficient.

Also, our system can be implemented on social media platform, so that sharing of any uncertain news can be evaluated.

Only text can be put as input in our system, in future taking input in the form of image and videos can make our system more wide in scope of news.

## References

- [1] *Automatic deception detection: Methods for finding fake news*; N. J. Conroy, V. L. Rubin, and Y. Chen, Proceedings of the Association for Information Science and Technology, 2015
- [2] *News Aggregation and Validation System*; Josef Moucachein, 2017
- [3] *Social Media and Fake News in the 2016 Election*; Hunt Alcott, 2017
- [4] *Deception Detection for News: Three Types of Fakes*; Victoria L. Rubin, Yimin Chen and Niall J. Conroy, University of Western Ontario, 2015



# Achievements

## 1. Conferences

- (a) *News Validation System*; Sufiyan Pawaskar, Nooralam Shaikh, Pradnyesh Rane, Prof. Tabrez Khan, Conference on Recent Trends in Computer Engineering, February and 2018 of attend (Venue : Thakur College of Engineering. )

## 2. Project Competitions

- (a) *News Validation System*; Sufiyan Pawaskar, Nooralam Shaikh, Pradnyesh Rane, Prof. Tabrez Khan, 4th National Level Project Exhibition Cum Poster Presentation, 9th March and 2018 of attend (Venue : Universal College of Engineering ) **Secured 2nd position.**

# MULTICON-W 2018

9<sup>th</sup> National & International Conferences & Workshops



**Platinum**



**Diamond**



**Patronage Partners**



**Gold Partners**



**Other Partners**



**THAKUR COLLEGE OF ENGINEERING & TECHNOLOGY**  
 (Approved by AICTE, Govt. of Maharashtra & Affiliated to University of Mumbai)  
 Thakur Village, Kandivall (E), Mumbai - 400101.

*Saylu Singh Charitable Trust's (Regd.)*

## Certificate of Appreciation

This is to certify that Dr./Mr./Ms. Sufiyan Pawaskar has presented/published a Short length paper with the title News Validation System in the Conference on Recent Trends in Computer Engineering (CRITCE 2018) organized during February, 23<sup>rd</sup> & 24<sup>th</sup>, 2018 at Thakur College of Engineering and Technology.

*(Signature)*  
(Dr. Sheetal Rathi)  
CONVENOR

*(Signature)*  
(Dr. R. R. Sedamkar)  
Technical Chair

*(Signature)*  
(Dr. B. K. Mishra)  
Principal & Programme Chair





ESTD 2005  
905 9905 : 2015 Certified  
MSA and NAAC Accredited

Department of Computer Engineering, AIKTC, New Panvel, Navi Mumbai

54











# Universal College of Engineering

**DTE Code: 3460**

(Permanently Unaided | Approved by AICTE, DTE & Affiliated to University of Mumbai)  
Near Bhajansons and Pinyadham, Kaman Bhiwandi Road, Vasai  
in association with I.E.T.E. - I.S.F. G.S.I. & I.S.A.



## 4<sup>th</sup> National Level Project Exhibition cum Poster Presentation

### Certificate of Merit


This is to certify that Ms./Mr. RADNYESH RANE

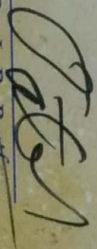
of A.I.K.T.C

College has won ~~First Prize~~ / ~~1st Runner up~~ / ~~2nd Runner up~~ / ~~Consolation Prize~~ in the

"4<sup>th</sup> National Level Project Exhibition cum Poster Presentation" 2018.

Date: 9<sup>th</sup> March 2018

  
Dr. Ajay Kumar  
(Principal)

  
Dr. I.B. Patil  
(Campus Director)