# SPAM DETECTION AND FAKE USER IDENTIFICATION NETWORKS

Submitted in partial fulfillment of the requirements for
the award of
Bachelor of Engineering degree in Computer Science and Engineering

by

## GUNDA AAKASH (Reg. No.37110001)



**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING
SCHOOL OF COMPUTING**

**SATHYABAMA**
**INSTITUTE OF SCIENCE AND TECHNOLOGY**
**(DEEMED TO BE UNIVERSITY)**
**Accredited with Grade "A" by NAAC**
**JEPPIAAR NAGAR, RAJIV GANDHI SALAI,**
**CHENNAI – 600 119**

**MARCH-2021**

# SATHYABAMA

## INSTITUTE OF SCIENCE AND TECHNOLOGY
### (DEEMED TO BE UNIVERSITY)
Accredited with "A" grade by NAAC
Jeppiaar Nagar, Rajiv Gandhi Salai, Chennai – 600 119
www.sathyabama.ac.in

DNV·GL
ISO 9001:2015

## DEPARTMENT OF COMPUTER SCIENCE AND TECHNOLOGY

## BONAFIDE CERTIFICATE

This is to certify that this project report is the bonafide work of **GUNDA AAKASH (Reg. No. 37110001)** who carried out the project entitled " **SPAM DETECTION AND FAKE USER IDENTIFICATION ON SOCIAL NETWORKS** " under my supervision from **August 2020** to **March 2021**.

### Internal Guide
Mrs. A.C.Santha Sheela,M.E.,(ph.d)

## Head of the Department

_____

**Submitted for Viva voce Examination held on** _____

**Internal Examiner**                                          **External Examiner**

# DECLARATION

I **GUNDA AAKASH** here by declare that the Project Report entitled "**SPAM DETECTION AND FAKE USER IDENTIFICATION ON SOCIA NETWORKS**" Is done by me under the guidance of **Mrs.santha Sheela,M,E.,PH.D.,**Department of Computer Science and Engineering at Sathyabama Institute of Science and Technology is submitted in partial fulfillment of the requirements for the award of Bachelor of Engineering degree in Computer Science and Engineering**.**

**DATE:**

**PLACE: CHENNAI**                                        **SIGNATURE OF THE CANDIDATE**

# ACKNOWLEDGEMENT

I am pleased to acknowledge my sincere thanks to **Board of Management** of **SATHYABAMA** for their kind encouragement in doing this project and for completing it successfully. I am grateful to them.

I convey my thanks to **Dr. T. Sasikala, M.E., Ph.D., Dean**, School of Computing, **Dr. S. Vigneswari, M.E., Ph.D., and Dr. L. Lakshmanan, M.E., Ph.D., Heads Of the Department** of Computer Science and Engineering for providing me necessary support and details at the right time during the progressive reviews.

I would like to express my sincere and deep sense of gratitude to my Project Guide **Mrs.A.c.Santha sheela,M.E.,Ph.d.,Assistant Professor**, for her valuable guidance,suggestions and constant encouragement paved way for the successful completion of my project work.

I wish to express my thanks to all Teaching and Non-teaching staff members of the **Department of Computer Science and Engineering** who were helpful in many ways for the completion of the project.

# ABSTRACT

Crime which is an unlawful act is increasing in our society day by day. With the advancement of the technology the criminals are also getting new ways of doing crimes. Crime takes place due to few common reasons that is money, imbalance mentality and emotions. After the crime takes place victims needs to go through a very complicated and a lengthy process for reporting the crime in police station. It"s also a very hectic process for the Crime branch to do it manually and maintain the records. So the crime Reporting system is the solution for all the victims and for the crime department. This will not only make the work easy but also it will help the users to access many features like news feed and all the updates regarding crime taking place in our locality. It will bring the police and the victims more Close and hence increasing the security. This makes the FIR registration Simple and easy hence making it time efficient. The system will help the crime department to take action as quick as possible and maintain the Database efficiently. The police can even update an alert to the citizens regarding the most wanted persons, lost belongings and any kind of emergency through this system. As a result this system will give a sustainable solution to the users, police and victims for managing the crime in a better and a structured way.

# TABLE OF CONTENTS

# LIST OF FIGURES

# CHAPTER 1
# INTRODUCTION

## 1.1 INTRODUCTION

It has become quite unpretentious to obtain any kind of information from any source across the world by using the Internet. The increased demand of social sites permits users to collect abundant amount of information and data about users. Huge volumes of data available on these sites also draw the attention of fake users. Twitter has rapidly become an online source for acquiring real-time information about users. Twitter is an Online Social Network(OSN) where users can share anything and everything, such as news, opinions, and even their moods. Several arguments can be held over different topics ,such as politics ,current affairs ,and important events. When a user tweets something, it is instantly conveyed to his/her followers, allowing them to outspread the received information at a much broader level [2]. With the evolution of OSNs, the need to study and analyze users' behaviors in online social platforms has intensified. Many people who do no have much information regarding the OSN scan easily be tricked by the fraudsters. There is also a demand to combat and place a control on the people who use OSNs only for advertisements and thus spam other people's accounts. Recently, the detection of spam in social networking sites attracted the attention of researchers. Spam detection is a difficult task in maintaining the security of social networks. It is essential to recognize spams in the OSN sites to save users from various kinds of malicious attacks and to preserve their security and privacy. These hazardous maneuvers adopted by spammers cause massive destruction of the community in the real world. Twitter spammers have various objectives, such as spreading invalid information, fake news, rumors, and spontaneous messages. Spammers achieve their malicious objectives through advertisements and several other means where they support different mailing lists and subsequently dispatch spam messages randomly to broadcast their interests. These activities cause disturbance to the original users who are known as non-spammers. In addition, it also decreases the repute of the OSN platforms. Therefore, it is essential to design a scheme to spot spammers so that corrective efforts can be taken to counter their malicious activities. Several research works have been carried out in the domain of Twitter spam detection. To encompass the existing state-of the-art, a few surveys have also been carried out on fake user identification from Twitter. Tin min et al. provide a survey of new methods and techniques to

identify Twitter spam detection. The above survey presents a comparative study of the current approaches. On the other hand, the authors in conducted a survey on different behaviors exhibited by spammers on Twitter social network. The study also provides a literature review that recognizes the existence of spammers on Twitter social network. Despite all the existing studies, there is still a gap in the existing literature. Therefore, to bridge the gap, we review state-of-the-art in the spammer detection and fake user identification on Twitter. Moreover, this Survey presents at taxonomy of the Twitter spam detection approaches and attempts to offer a detailed description of recent developments in the domain. The aim of this paper is to identify different approaches of spam detection on Twitter and to present a taxonomy by classifying these approaches into several categories. For classification, we have identified four means of reporting spammers that can be helpful in identifying fake identities of users. Spammers can be identified based on: (i) fake content, (ii) URL based spam detection, (iii) detecting spam in trending topics, and(iv)fakeuseridentification.Table1provides a comparison of existing techniques and helps users to recognize the significance and effectiveness of the proposed methodologies in addition to providing a comparison of their goals and results.Table2 compares different features that are used for identifying spam on Twitter. We anticipate that this survey will help readers find diverse information on spammer detection techniques at a single point. This article is structured such that Section II presents the taxonomy for the spammer detection techniques on Twitter. The comparison of proposed methods for detecting spammers on Twitter is discussed in Section III. Section IV presents an overall analysis and discussion, whereas Section V concludes the paper and highlights some directions for future work.

## 2.2 DOMAIN INTRODUCTION

**What Is A Social Network?**

Wikipedia defines a social network service as a service which "focuses on the building and verifying of online social networks for communities of people who share interests and activities, or who are interested in exploring the interests and activities of others, and which necessitates the use of software."

A report published by OCLC provides the following definition of social networking sites: "Web sites primarily designed to facilitate interaction between users who share interests, attitudes and activities, such as Facebook, Mix i and My Space."

**What Can Social Networks Be Used For?**

Social networks can provide a range of benefits to members of an organisation:

**Support for learning**: Social networks can enhance informal learning and support social connections within groups of learners and with those involved in the support of learning.

**Support for members of an organisation**: Social networks can potentially be used my all members of an organisation, and not just those involved in working with students. Social networks can help the development of communities of practice.

**Engaging with others**: Passive use of social networks can provide valuable business intelligence and feedback on institutional services (although this may give rise to ethical concerns).

**Ease of access to information and applications**: The ease of use of many social networking services can provide benefits to users by simplifying access to other tools and applications. The Facebook Platform provides an example of how a social networking service can be used as an environment for other tools.

**Common interface**: A possible benefit of social networks may be the common interface which spans work / social boundaries. Since such services are often used in a personal capacity the interface and the way the service works may be familiar, thus minimising training

and support needed to exploit the services in a professional context. This can, however, also be a barrier to those who wish to have strict boundaries between work and social activities.

**Examples of Social Networking Services**

Examples of popular social networking services include:

**Facebook**: Facebook is a social networking Web site that allows people to communicate with their friends and exchange information. In May 2007 Facebook launched the Facebook Platform which provides a framework for developers to create applications that interact with core Facebook features

**My Space**: My Space is a social networking Web site offering an interactive, user-submitted network of friends, personal profiles, blogs and groups, commonly used for sharing photos, music and videos.

**Ning**: An online platform for creating social Web sites and social networks aimed at users who want to create networks around specific interests or have limited technical skills.

**Twitter**: Twitter is an example of a micro-blogging service. Twitter can be used in a variety of ways including sharing brief information with users and providing support for one's peers.

**2.2 OPPURTUNITIES AND CHALLENGES**

The popularity and ease of use of social networking services have excited institutions with their potential in a variety of areas. However effective use of social networking services poses a number of challenges for institutions including long-term sustainability of the services; user concerns over use of social tools in a work or study context; a variety of technical issues and legal issues such as copyright, privacy, accessibility; etc.

Institutions would be advised to consider carefully the implications before promoting significant use of such services.

**What is Secure Computing?Computer security** (Also known as cyber security or IT Security) is information security as applied to computers and networks. The field covers all the processes

and mechanisms by which computer-based equipment, information and services are protected from unintended or unauthorized access, change or destruction. Computer security also includes protection from unplanned events and natural disasters. Otherwise, in the computer industry, the term security -- or the phrase computer security -- refers to techniques for ensuring that data stored in a computer cannot be read or compromised by any individuals without authorization. Most computer security measures involve data encryption and passwords. Data encryption is the translation of data into a form that is unintelligible without a deciphering mechanism. A password is a secret word or phrase that gives a user access to a particular program or system.



Diagram clearly explain the about the secure computing

**Working conditions and basic needs in the secure computing:**

If you don't take basic steps to protect your work computer, you put it and all the information on it at risk. You can potentially compromise the operation of other computers on your organization's network, or even the functioning of the network as a whole.

**1. Physical security:**

Technical measures like login passwords, anti-virus are essential.    (More about those below)  However, a secure physical space is the first and more important line of defense.

Is the place you keep your workplace computer secure enough to prevent theft or access to it while you are away? While the Security Department provides coverage across the Medical

center, it only takes seconds to steal a computer, particularly a portable device like a laptop or a PDA.   A computer should be secured like any other valuable possession when you are not present.

Human threats are not the only concern. Computers can be compromised by environmental mishaps (e.g., water, coffee) or physical trauma.   Make sure the physical location of your computer takes account of those risks as well.

## 2.  Access passwords:

The University's networks and shared information systems are protected in part by login credentials (user-IDs and passwords). Access passwords are also an essential protection for personal computers in most circumstances. Offices are usually open and shared spaces, so physical access to computers cannot be completely controlled.

To protect your computer, you should consider setting   passwords   for   particularly sensitive applications resident on the computer (e.g., data analysis software), if the software provides that capability.

## 3.  Prying eye protection:

Because we deal with all facets of clinical, research, educational and administrative data here on the medical campus, it is important to do everything possible to minimize exposure of data to unauthorized individuals.

## 4.  Anti-virus software:

Up-to-date, properly configured anti-virus software is essential. While we have server-side anti-virus software on our network computers, you still need it on the client side (your computer).

## 5.  Firewalls:

Which are being secure by the client side and also on the server side. It is an important to keep secure files by firewalls.

**6.** Anti-virus products inspect files on your computer and in email. Firewall software and hardware monitor communications between your computer and the outside world. That is essential for any networked computer.

**7. Software updates:**

It is critical to keep software up to date, especially the operating system, anti-virus and anti-spyware, email and browser software. The newest versions will contain fixes for discovered vulnerabilities.

Almost all anti-virus have automatic update features (including SAV). Keeping the "signatures" (digital patterns) of malicious software detectors up-to-date is essential for these products to be effective.

**8. Keep secure backups:**

Even if you take all these security steps, bad things can still happen. Be prepared for the worst by making backup copies of critical data, and keeping those backup copies in a separate, secure location. For example, use supplemental hard drives, CDs/DVDs, or flash drives to store critical, hard-to-replace data.

**9. Report problems:**

If you believe that your computer or any data on it has been compromised, your should make a information security incident report. That is required by University policy for all data on our systems, and legally required for health, education, financial and any other kind of record containing identifiable personal information.

**Benefits of secure computing:**

- **Protect-yourself-Civilliability**:

  You may be held legally liable to compensate a third party should they experience financial damage or distress as a result of their personal data being stolen from you or leaked by you.

- **Protect-your-credibility-Compliance**:

  You may require compliancy with the Data Protection Act, the FSA, SOX or other regulatory standards. Each of these bodies stipulates that certain measures be taken to protect the data on your network.

- **Protect-your-reputation-Spam:**

  A common use for infected systems is to join them to a botnet (a collection of infected machines which takes orders from a command server) and use them to send out spam. This spam can be traced back to you, your server could be blacklisted and you could be unable to send email.

- **Protect-your-income-advantage:**

  There are a number of "hackers-for-hire" advertising their services on the internet selling their skills in breaking into company's servers to steal client databases, proprietary software, merger and acquisition information, personnel details *et al*.

- **Protect-your-business-Blackmail**:

  A seldom-reported source of income for "hackers" is to·break into your server, change all your passwords and lock you out of it. The password is then sold back to you. Note: the "hackers" may implant a backdoor program on your server so that they can repeat the exercise at will.

- **Protect-your-investment-Free-storage:**

  Your server's harddrive space is used (or sold on) to house the hacker's video clips, music collections, pirated software or worse. Your server or computer then becomes continuously slow and your internet connection speeds deteriorate due to the number of people connecting to your server in order to download the offered wares.

# CHAPTER 2
# LITERATURE SURVEY

## 2.1 STATISTICAL FEATURES

**AUTHORS:** C. Chen, Y. Wang, J. Zhang, Y. Xiang, W. Zhou, and G. Min

Twitter spam has become a critical problem nowadays. Recent works focus on applying machine learning techniques for Twitter spam detection, which make use of the statistical features of tweets. In our labeled tweets data set, however, we observe that the statistical properties of spam tweets vary over time, and thus, the performance of existing machine learning-based classifiers decreases. This issue is referred to as "Twitter Spam Drift". In order to tackle this problem, we first carry out a deep analysis on the statistical features of one million spam tweets and one million non-spam tweets, and then propose a novel L fun scheme. The proposed scheme can discover "changed" spam tweets from unlabeled tweets and incorporate them into classifier's training process. A number of experiments are performed to evaluate the proposed scheme. The results show that our proposed L fun scheme can significantly improve the spam detection accuracy in real-world scenarios.

## 2) Automatically identifying fake news in popular Twitter threads

**AUTHORS:** C. Buntain and J. Gol beck

Information quality in social media is an increasingly important issue, but web-scale data hinders experts' ability to assess and correct much of the inaccurate content, or "fake news," present in these platforms. This paper develops a method for automating fake news detection on Twitter by learning to predict accuracy assessments in two credibility-focused Twitter datasets:

CREDBANK, a crowdsourced dataset of accuracy assessments for events in Twitter, and PHEME, a dataset of potential rumors in Twitter and journalistic assessments of their accuracies. We apply this method to Twitter content sourced from BuzzFeed's fake news dataset and show models trained against crowdsourced workers outperform models based on journalists' assessment and models trained on a pooled dataset of both crowdsourced workers and journalists. All three datasets, aligned into a uniform format, are also publicly available. A feature analysis then identifies features that are most predictive for crowdsourced and journalistic accuracy assessments, results of which are consistent with prior work. We close with a discussion contrasting accuracy and credibility and why models of non-experts outperform models of journalists for fake news detection in Twitter.

## 3) A performance evaluation of machine learning-based streaming spam tweets detection

**AUTHORS:** C. Chen, J. Zhang, Y. Xie, Y. Xiang,W. Zhou, M. M. Hassan, A. AlElaiwi, and M. Alrubaian

The popularity of Twitter attracts more and more spammers. Spammers send unwanted tweets to Twitter users to promote websites or services, which are harmful to normal users. In order to stop spammers, researchers have proposed a number of mechanisms. The focus of recent works is on the application of machine learning techniques into Twitter spam detection. However, tweets are retrieved in a streaming way, and Twitter provides the Streaming API for developers and researchers to access public tweets in real time. There lacks a performance evaluation of existing machine learning-based streaming spam detection methods. In this paper, we bridged the gap by carrying out a performance evaluation, which was from three different aspects of data, feature, and model. A big ground-truth of over 600 million public tweets was created by using a commercial URL-based security tool. For real-time spam detection, we further extracted 12 lightweight features for tweet representation. Spam detection was then transformed to a binary classification problem in the feature space and can be solved by conventional machine learning algorithms. We evaluated the impact of different factors to the spam detection performance, which included spam to nonspam ratio, feature discretization, training data size, data sampling,

time-related data, and machine learning algorithms. The results show the streaming spam tweet detection is still a big challenge and a robust detection technique should take into account the three aspects of data, feature, and model.

## 4) A model-based approach for identifying spammers in social networks

**AUTHORS:** F. Fathaliani and M. Bouguessa

In this paper, we view the task of identifying spammers in social networks from a mixture modeling perspective, based on which we devise a principled unsupervised approach to detect spammers. In our approach, we first represent each user of the social network with a feature vector that reflects its behaviour and interactions with other participants. Next, based on the estimated users feature vectors, we propose a statistical framework that uses the Dirichlet distribution in order to identify spammers. The proposed approach is able to automatically discriminate between spammers and legitimate users, while existing unsupervised approaches require human intervention in order to set informal threshold parameters to detect spammers. Furthermore, our approach is general in the sense that it can be applied to different online social sites. To demonstrate the suitability of the proposed method, we conducted experiments on real data extracted from Instagram and Twitter.

## 5) Spam detection of Twitter traffic: A framework based on random forests and non-uniform feature sampling

**AUTHORS:** C. Meda, E. Ragusa, C. Gianoglio, R. Zunino, A. Ottaviano, E. Scillia, and R. Surlinelli

Law Enforcement Agencies cover a crucial role in the analysis of open data and need effective techniques to filter troublesome information. In a real scenario, Law Enforcement Agencies analyze Social Networks, i.e. Twitter, monitoring events and profiling accounts. Unfortunately,

between the huge amount of internet users, there are people that use microblogs for harassing other people or spreading malicious contents. Users' classification and spammers' identification is a useful technique for relieve Twitter traffic from uninformative content. This work proposes a framework that exploits a non-uniform feature sampling inside a gray box Machine Learning System, using a variant of the Random Forests Algorithm to identify spammers inside Twitter traffic. Experiments are made on a popular Twitter dataset and on a new dataset of Twitter users. The new provided Twitter dataset is made up of users labeled as spammers or legitimate users, described by 54 features. Experimental results demonstrate the effectiveness of enriched feature sampling method

## 2.1.1 EXISTING SYSTEM

❖ Tingmin*et al.* provide a survey of new methods and techniques to identify Twitter spam detection. The above survey presents a comparative study of the current approaches.

❖ On the other hand, S. J. Somanet. al. conducted a survey on different behaviors exhibited by spammers on Twitter social network. The study also provides a literature review that recognizes the existence of spammers on Twitter social network.

❖ Despite all the existing studies, there is still a gap in the existing literature. Therefore, to bridge the gap, we review state-of-the-art in the spammer detection and fake user identification on Twitter

## DISADVANTAGES EXISTING SYSTEM

❖ No efficient methods used.

❖ No real time data's used.

❖ More complex

# CHAPTER 3

# SYSTEM DESIGN & METHODOLOGY

## 3.1 PROPOSED SYSTEM

❖ The aim of this paper is to identify fake user detection on Twitter and to present a framework by classifying these approaches into several categories. For classification, we have identified four means of reporting spammers that can be helpful in identifying fake identities of users. Spammers can be identified based on: (i) fake content, (ii) URL based spam detection, (iii) detecting spam in trending topics, and (iv) fake user identification.

❖ Moreover, the analysis also shows that machine learning-based techniques can be effective for identifying fake user on Twitter. However, the selection of the most feasible techniques and methods is highly dependent on the available data.

## 3.2 ADVANTAGES OF PROPOSED SYSTEM

❖ This study includes machine learning methodology proposed using real time datasets and with different characteristics and accomplishments.

❖ The proposed system is more effective and accurate than other existing systems.

❖ Tested with real time data's.

## 3.3 ARCHITECTURE DIAGRAM

Generate tweets

SPAMMER

TWITTER

Get sample tweets from Twitter

Pre-processing

ADMIN

Analysis Spam tweets from the Twitter

**3.4 SYSTEM REQUIREMENTS:**

**HARDWARE REQUIREMENTS:**

- System : Pentium Dual Core.
- Hard Disk : 120 GB.
- Monitor : 15'' LED
- Input Devices : Keyboard, Mouse
- Ram : 4 GB.

**SOFTWARE REQUIREMENTS:**

- Operating system: Windows 7/10.
- Coding Language :Python
- Tool : Pi-champ
- Database : MYSQL

# CHAPTER 4

# SOFTWARE ENVIRONMENT

## 4.1 SOFTWARE LANGUAGE

**Python:**

Python is a high-level, interpreted, interactive and object-oriented scripting language. Python is designed to be highly readable. It uses English keywords frequently where as other languages use punctuation, and it has fewer syntactical constructions than other languages.

- **Python is Interpreted** − Python is processed at runtime by the interpreter. You do not need to compile your program before executing it. This is similar to PERL and PHP.

- **Python is Interactive** − You can actually sit at a Python prompt and interact with the interpreter directly to write your programs.

- **Python is Object-Oriented** − Python supports Object-Oriented style or technique of programming that encapsulates code within objects.

- **Python is a Beginner's Language** − Python is a great language for the beginner-level programmers and supports the development of a wide range of applications from simple text processing to WWW browsers to games.

**History of Python**

Python was developed by Guido van Rossum in the late eighties and early nineties at the National Research Institute for Mathematics and Computer Science in the Netherlands.

Python is derived from many other languages, including ABC, Modula-3, C, C++, Algol-68, SmallTalk, and Unix shell and other scripting languages.

Python is copyrighted. Like Perl, Python source code is now available under the GNU General Public License (GPL).

Python is now maintained by a core development team at the institute, although Guido van Rossum still holds a vital role in directing its progress.

**Python Features**

Python's features include −

- **Easy-to-learn** − Python has few keywords, simple structure, and a clearly defined syntax. This allows the student to pick up the language quickly.

- **Easy-to-read** − Python code is more clearly defined and visible to the eyes.

- **Easy-to-maintain** − Python's source code is fairly easy-to-maintain.

- **A broad standard library** − Python's bulk of the library is very portable and cross-platform compatible on UNIX, Windows, and Macintosh.

- **Interactive Mode** − Python has support for an interactive mode which allows interactive testing and debugging of snippets of code.

- **Portable** − Python can run on a wide variety of hardware platforms and has the same interface on all platforms.

- **Extendable** − You can add low-level modules to the Python interpreter. These modules enable programmers to add to or customize their tools to be more efficient.

- **Databases** − Python provides interfaces to all major commercial databases.

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

- **Scalable** − Python provides a better structure and support for large programs than shell scripting.

Apart from the above-mentioned features, Python has a big list of good features, few are listed below −

- It supports functional and structured programming methods as well as OOP.

- It can be used as a scripting language or can be compiled to byte-code for building large applications.

- It provides very high-level dynamic data types and supports dynamic type checking.

- It supports automatic garbage collection.

- It can be easily integrated with C, C++, COM, ActiveX, CORBA, and Java.

Python is available on a wide variety of platforms including Linux and Mac OS X. Let's understand how to set up our Python environment.

**Getting Python**

The most up-to-date and current source code, binaries, documentation, news, etc., is available on the official website of Python https://www.python.org.

Windows Installation

Here are the steps to install Python on Windows machine.

- Open a Web browser and go to https://www.python.org/downloads/.

- Follow the link for the Windows installer python-XYZ.msifile where XYZ is the version you need to install.

- To use this installer python-XYZ.msi, the Windows system must support Microsoft Installer 2.0. Save the installer file to your local machine and then run it to find out if your machine supports MSI.

- Run the downloaded file. This brings up the Python install wizard, which is really easy to use. Just accept the default settings, wait until the install is finished, and you are done.

The Python language has many similarities to Perl, C, and Java. However, there are some definite differences between the languages.

**First Python Program**

Let us execute programs in different modes of programming.

**Interactive Mode Programming**

Invoking the interpreter without passing a script file as a parameter brings up the following prompt −

```
$ python

Python2.4.3(#1,Nov112010,13:34:43)

[GCC 4.1.220080704(RedHat4.1.2-48)] on linux2

Type"help","copyright","credits"or"license"for more information.

>>>
```

Type the following text at the Python prompt and press the Enter −

```
>>>print"Hello, Python!"
```

If you are running new version of Python, then you would need to use print statement with parenthesis as in **print ("Hello, Python!");**. However in Python version 2.4.3, this produces the following result −

```
Hello, Python!
```

**Script Mode Programming**

Invoking the interpreter with a script parameter begins execution of the script and continues until the script is finished. When the script is finished, the interpreter is no longer active.

Let us write a simple Python program in a script. Python files have extension **.py**. Type the following source code in a test.py file −

```
print"Hello, Python!"
```

We assume that you have Python interpreter set in PATH variable. Now, try to run this program as follows −

$ python test.py

This produces the following result −

Hello, Python!

**Flask Framework:**

Flask is a web application framework written in Python. Armin Ronacher, who leads an international group of Python enthusiasts named Pocco, develops it. Flask is based on Werkzeug WSGI toolkit and Jinja2 template engine. Both are Pocco projects.

Http protocol is the foundation of data communication in world wide web. Different methods of data retrieval from specified URL are defined in this protocol.

The following table summarizes different http methods −

| Sr.No | Methods & Description |
| --- | --- |
| 1 | **GET** <br><br> Sends data in unencrypted form to the server. Most common method. |

| 2 | **HEAD** |
|---|---|
| | Same as GET, but without response body |
| 3 | **POST** |
| | Used to send HTML form data to server. Data received by POST method is not cached by server. |
| 4 | **PUT** |
| | Replaces all current representations of the target resource with the uploaded content. |
| 5 | **DELETE** |
| | Removes all current representations of the target resource given by a URL |

By default, the Flask route responds to the **GET** requests. However, this preference can be altered by providing methods argument to **route()** decorator.

In order to demonstrate the use of **POST** method in URL routing, first let us create an HTML form and use the **POST** method to send form data to a URL.

Save the following script as login.html

```
<html>

<body>

<formaction="http://localhost:5000/login"method="post">

<p>Enter Name:</p>

<p><inputtype="text"name="nm"/></p>

<p><inputtype="submit"value="submit"/></p>
```

```
</form>

</body>

</html>
```

Now enter the following script in Python shell.

```python
from flask importFlask, redirect,url_for, request

app=Flask(_name_)

@app.route('/success/<name>')

def success(name):

return'welcome %s'% name

@app.route('/login',methods=['POST','GET'])

def login():

ifrequest.method=='POST':

user=request.form['nm']

return redirect(url_for('success',name= user))

else:

user=request.args.get('nm')

return redirect(url_for('success',name= user))

if__name__=='_main_':

app.run(debug =True)
```

After the development server starts running, open **login.html** in the browser, enter name in the text field and click **Submit**.



Form data is POSTed to the URL in action clause of form tag.

**http://localhost/login** is mapped to the **login()** function. Since the server has received data by **POST** method, value of 'nm' parameter obtained from the form data is obtained by −

```
user = request.form['nm']
```

It is passed to **'/success'** URL as variable part. The browser displays a **welcome** message in the window.

welcome mvl

Change the method parameter to **'GET'** in **login.html** and open it again in the browser. The data received on server is by the **GET** method. The value of 'nm' parameter is now obtained by –

```
User = request.args.get('nm')
```

Here, **args** is dictionary object containing a list of pairs of form parameter and its corresponding value. The value corresponding to 'nm' parameter is passed on to '/success' URL as before.

### 4.2 MODULES:

❖ Admin Module

❖ Data Collection

❖ Train and Test

❖ Machine Learning Technique

❖ Detection of Fake User

### 4.3 MODULE DESCRIPTIONS:

**Admin Module:** In the first module, we develop the Online Social Networking (OSN) system module. We build up the system with the feature of Online Social Networking System, Twitter. Where, this module is used for admin login with their authentication.

**Data Collection:**

We will be using a Python Library called *Tweepy* to connect to the Twitter API and collect the data. We download tweets containing certain key words, to incorporate the words or hash tags that contain relevant keyword related to fake users.

Some of the most important fields are:

- *text*, which contains the text included in the tweet.
- *created_at,* which is a timestamp of when the tweet was created.
- *user*, which contains information about the user that created the tweet, like the username and user id.

**Train and Test:**

- ❖ We present the proposed framework for metadata features are extracted from available additional information regarding the tweets of a user, whereas content-based features aim to observe the message posting behavior of a user and the quality of the text that the user uses in posts.

**4.4  Machine Learning Technique:**

- ❖ The number of features, which are associated with tweet content, and the characteristics of users are recognized for the detection of spammers. These features are considered as the characteristics of machine learning process for categorizing users, i.e., to know whether they are spammers or not.
- ❖ In order to recognize the approach for detecting spammers on Twitter, the labelled collection in pre-classification of fake user and legitimate user has been done. Next, those steps are taken which are needed for the construction of labeled collection and acquired various desired properties.

- ❖ In other words, steps which are essential to be examined to develop the collection of users that can be labelled as fake user or legitimate user. At the end, user attributes are identified based on their behavior, e.g., who they interact with and what is the frequency of their interaction.

- ❖ In order to confirm this instinct, features of users of the labelled collection has been checked. Two attribute sets are considered, i.e., content attributes and user behavior attributes, to differentiate one user from the other.

## 4.5 Detection of Fake User:

- ❖ In this module, we implement the collection of tweets with respect to trending topics on Twitter. After storing the tweets in a particular file format, the tweets are subsequently analyzed.

- ❖ Labelling of fake user is performed to check through all datasets that are available to detect the malignant.

- ❖ Feature extraction separates the characteristics construct based on the language model that uses language as a tool and helps in determining whether the user is fake or not.

- ❖ The classification of data set is performed by shortlisting the set of tweets that is described by the set of features provided to the classifier to instruct the model and to acquire the knowledge for spam detection.

- ❖ The fake user detection uses the classification technique to accept tweets as the input and classify the fake user and legitimate user.

## 4.6 REQUIREMENT ANALYSIS

Requirement analysis, also called requirement engineering, is the process of determining user expectations for a new modified product. It encompasses the tasks that determine the need for analysing, documenting, validating and managing software or system requirements. The requirements should be documentable, actionable, measurable, testable and traceable related to identified business needs or opportunities and define to a level of detail, sufficient for system design.

**4.7 FUNCTIONAL REQUIREMENTS**

It is a technical specification requirement for the software products. It is the first step in the requirement analysis process which lists the requirements of particular software systems including functional, performance and security requirements. The function of the system depends mainly on the quality hardware used to run the software with given functionality.

**Usability**

It specifies how easy the system must be use. It is easy to ask queries in any format which is short or long, porter stemming algorithm stimulates the desired response for user.

**Robustness**

It refers to a program that performs well not only under ordinary conditions but also under unusual conditions. It is the ability of the user to cope with errors for irrelevant queries during execution.

## Security

The state of providing protected access to resource is security. The system provides good security and unauthorized users cannot access the system there by providing high security.

**Reliability**

It is the probability of how often the software fails. The measurement is often expressed in MTBF (Mean Time Between Failures). The requirement is needed in order to ensure that the processes work correctly and completely without being aborted. It can handle any load and survive and survive and even capable of working around any failure.

**Compatibility**

It is supported by version above all web browsers. Using any web servers like localhost makes the system real-time experience.

**Flexibility**

The flexibility of the project is provided in such a way that is has the ability to run on different environments being executed by different users.

**Safety**

Safety is a measure taken to prevent trouble. Every query is processed in a secured manner without letting others to know one's personal information.

**4.7 NON- FUNCTIONAL REQUIREMENTS**

**Portability**

It is the usability of the same software in different environments. The project can be run in any operating system.

**Performance**

These requirements determine the resources required, time interval, throughput and everything that deals with the performance of the system.

**Accuracy**

The result of the requesting query is very accurate and high speed of retrieving information. The degree of security provided by the system is high and effective.

**Maintainability**

Project is simple as further updates can be easily done without affecting its stability. Maintainability basically defines that how easy it is to maintain the system. It means that how easy it is to maintain the system, analyse, change and test the application. Maintainability of this project is simple as further updates can be easily done without affecting its maintenance.

# CHAPTER 5

# SYSTEM DESIGN AND TESTING PLAN

## 5.1 INPUT DESIGN

The input design is the link between the information system and the user. It comprises the developing specification and procedures for data preparation and those steps are necessary to put transaction data in to a usable form for processing can be achieved by inspecting the computer to read data from a written or printed document or it can occur by having people keying the data directly into the system. The design of input focuses on controlling the amount of input required, controlling the errors, avoiding delay, avoiding extra steps and keeping the process simple. The input is designed in such a way so that it provides security and ease of use with retaining the privacy. Input Design considered the following things:

- ➢ What data should be given as input?
- ➢ How the data should be arranged or coded?
- ➢ The dialog to guide the operating personnel in providing input.
- ➢ Methods for preparing input validations and steps to follow when error occur.

## 5.2 OUTPUT DESIGN

A quality output is one, which meets the requirements of the end user and presents the information clearly. In any system results of processing are communicated to the users and to other system through outputs. In output design it is determined how the information is to be displaced for immediate need and also the hard copy output. It is the most important and direct source information to the user. Efficient and intelligent output design improves the system's relationship to help user decision-making.

The output form of an information system should accomplish one or more of the following objectives.

- ➢ Convey information about past activities, current status or projections of the
- ➢ Future.
- ➢ Signal important events, opportunities, problems, or warnings.
- ➢ Trigger an action.
- ➢ Confirm an action.

## FEASIBILITY STUDY

The feasibility of the project is analyzed in this phase and business proposal is put forth with a very general plan for the project and some cost estimates. During system analysis the feasibility study of the proposed system is to be carried out. This is to ensure that the proposed system is not a burden to the company. For feasibility analysis, some understanding of the major requirements for the system is essential.

Three key considerations involved in the feasibility analysis are

- ♦ ECONOMICAL FEASIBILITY
- ♦ TECHNICAL FEASIBILITY
- ♦ SOCIAL FEASIBILITY

## ECONOMICAL FEASIBILITY

This study is carried out to check the economic impact that the system will have on the organization. The amount of fund that the company can pour into the research and development of the system is limited. The expenditures must be justified. Thus the developed system as well within the budget and this was achieved because most of the technologies used are freely available. Only the customized products had to be purchased. This study is carried out to check the technical feasibility, that is, the technical requirements of the system. Any system
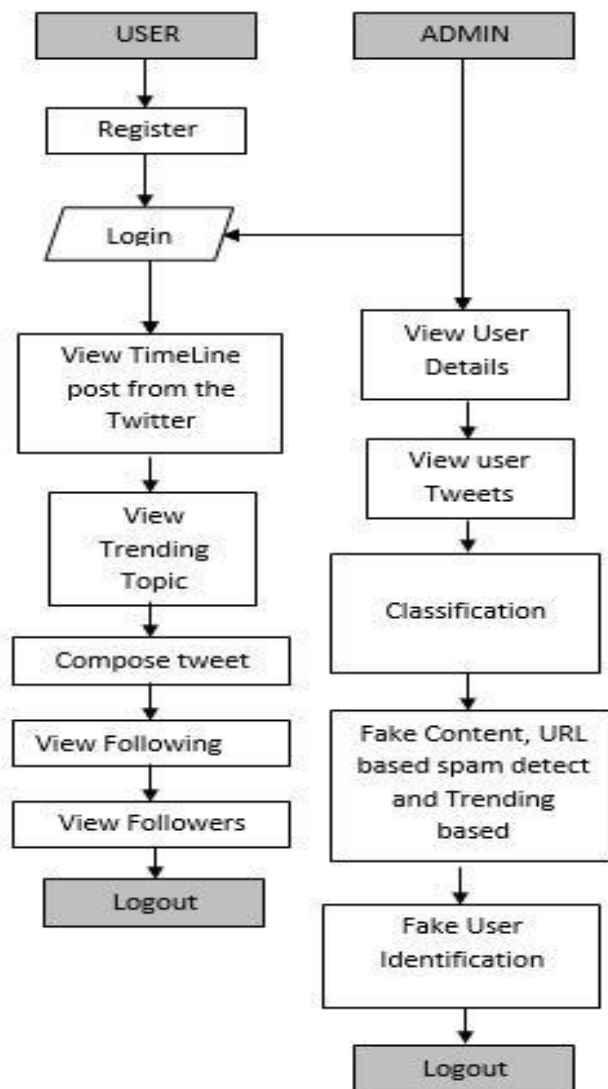
developed must not have a high demand on the available technical resources. This will lead to high demands on the available technical resources. This will lead to high demands being placed on the client. The developed system must have a modest requirement, as only minimal or null changes are required for implementing this system.

**SOCIAL FEASIBILITY**

The aspect of study is to check the level of acceptance of the system by the user. This includes the process of training the user to use the system efficiently. The user must not feel threatened by the system, instead must accept it as a necessity. The level of acceptance by the users solely depends on the methods that are employed to educate the user about the system and to make him familiar with it. His level of confidence must be raised so that he is also able to make some constructive criticism, which is welcomed, as he is the final user of the system.

**5.4 DATA FLOW DIAGRAM:**

1. The DFD is also called as bubble chart. It is a simple graphical formalism that can be used to represent a system in terms of input data to the system, various processing carried out on this data, and the output data is generated by this system.
2. The data flow diagram (DFD) is one of the most important modeling tools. It is used to model the system components. These components are the system process, the data used by the process, an external entity that interacts with the system and the information flows in the system.
3. DFD shows how the information moves through the system and how it is modified by a series of transformations. It is a graphical technique that depicts information flow and the transformations that are applied as data moves from input to output.
4. DFD is also known as bubble chart. A DFD may be used to represent a system at any level of abstraction. DFD may be partitioned into levels that represent increasing information flow and functional detail.

```
           USER                        ADMIN
            │                            │
            ▼                            │
       ┌─────────┐                       │
       │ Register │                      │
       └─────────┘                       │
            │                            │
            ▼                            │
       ╱─────────╱◄─────────────────────┘
      ╱  Login  ╱
     ╱─────────╱
            │                            │
            ▼                            ▼
  ┌──────────────┐              ┌──────────────┐
  │ View TimeLine │             │  View User   │
  │ post from the │             │   Details    │
  │   Twitter     │             └──────────────┘
  └──────────────┘                     │
            │                          ▼
            ▼                   ┌──────────────┐
  ┌──────────────┐             │  View user   │
  │    View      │             │   Tweets     │
  │  Trending    │             └──────────────┘
  │   Topic      │                    │
  └──────────────┘                    ▼
            │                  ┌──────────────┐
            ▼                  │Classification│
  ┌──────────────┐            └──────────────┘
  │ Compose tweet │                  │
  └──────────────┘                   ▼
            │                ┌──────────────────┐
            ▼                │ Fake Content, URL │
  ┌──────────────┐          │ based spam detect │
  │ View Following│         │   and Trending    │
  └──────────────┘          │      based        │
            │                └──────────────────┘
            ▼                        │
  ┌──────────────┐                   ▼
  │ View Followers│         ┌──────────────┐
  └──────────────┘          │  Fake User   │
            │               │Identification│
            ▼               └──────────────┘
      ┌─────────┐                   │
      │ Logout  │                   ▼
      └─────────┘            ┌─────────┐
                             │ Logout  │
                             └─────────┘
```

# CHAPTER 6

# UML DIAGRAMS

UML stands for Unified Modeling Language. UML is a standardized general-purpose modeling language in the field of object-oriented software engineering. The standard is managed, and was created by, the Object Management Group.

The goal is for UML to become a common language for creating models of object oriented computer software. In its current form UML is comprised of two major components: a Meta-model and a notation. In the future, some form of method or process may also be added to; or associated with, UML.

The Unified Modeling Language is a standard language for specifying, Visualization, Constructing and documenting the artifacts of software system, as well as for business modeling and other non-software systems.

The UML represents a collection of best engineering practices that have proven successful in the modeling of large and complex systems.

The UML is a very important part of developing objects oriented software and the software development process. The UML uses mostly graphical notations to express the design of software projects.
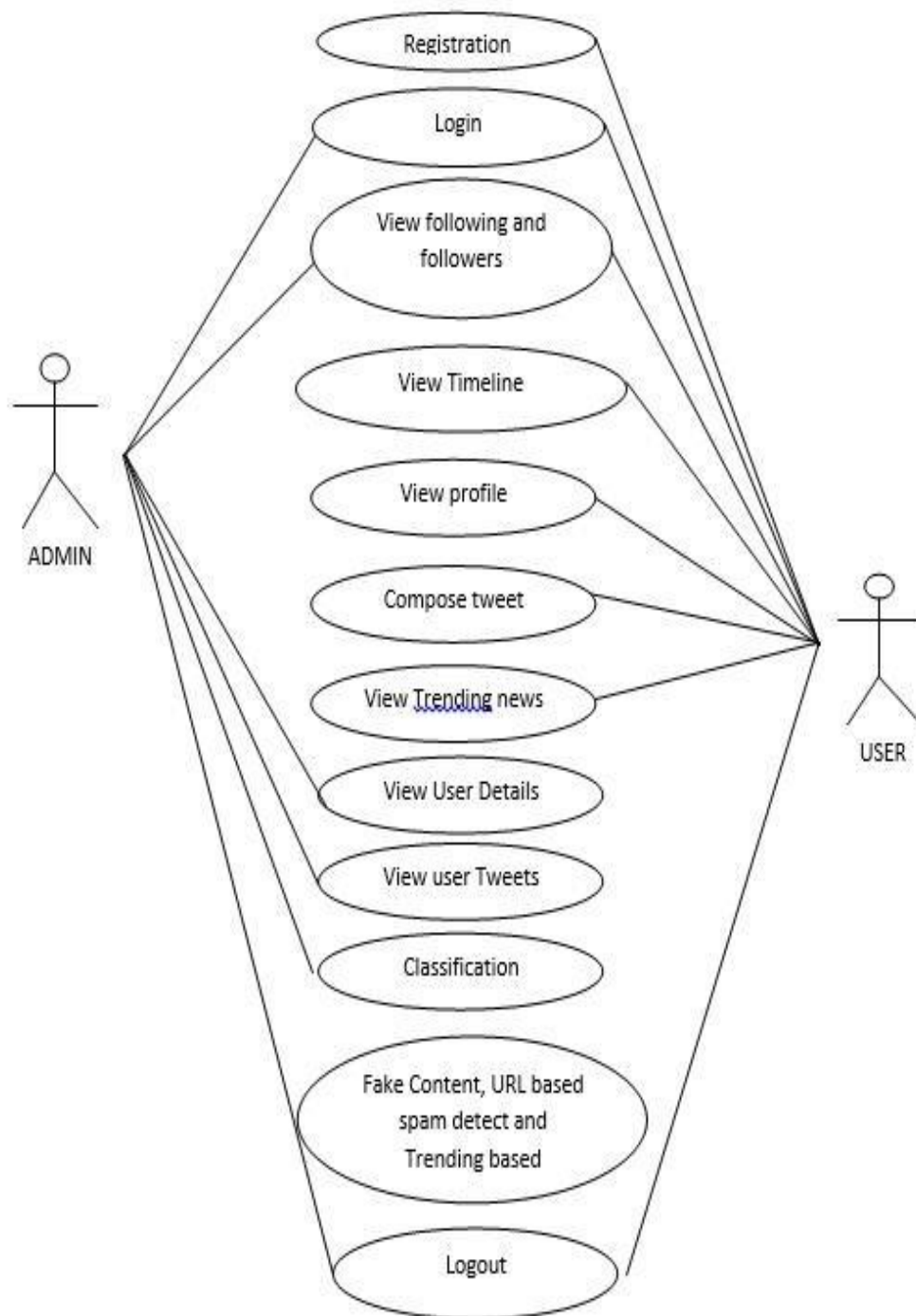
**GOALS:**

The Primary goals in the design of the UML are as follows:

1. Provide users a ready-to-use, expressive visual modeling Language so that they can develop and exchange meaningful models.
2. Provide extendibility and specialization mechanisms to extend the core concepts.
3. Be independent of particular programming languages and development process.

4. Provide a formal basis for understanding the modeling language.

5. Encourage the growth of OO tools market.

6. Support higher level development concepts such as collaborations, frameworks, patterns and components.
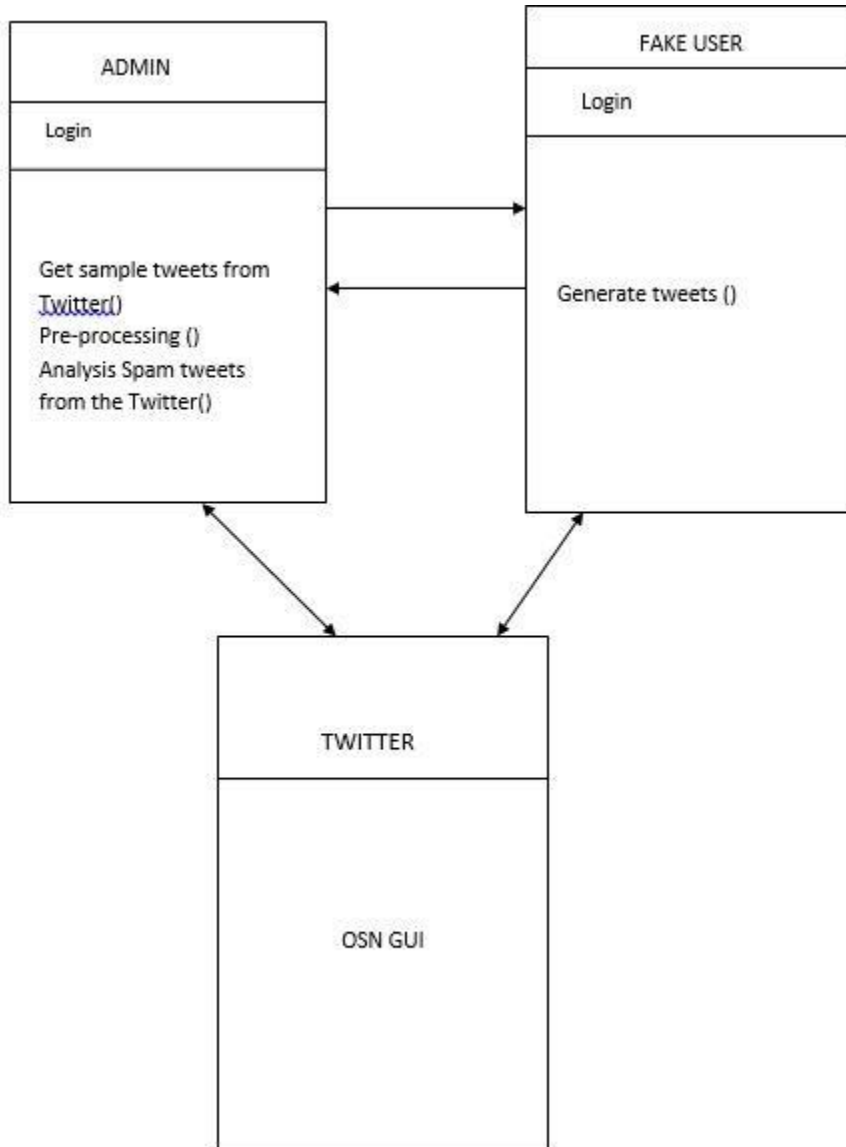
7. Integrate best practices.

**6.1 <u>USE CASE DIAGRAM:</u>**

      A use case diagram in the Unified Modeling Language (UML) is a type of behavioral diagram defined by and created from a Use-case analysis. Its purpose is to present a graphical overview of the functionality provided by a system in terms of actors, their goals (represented as use cases), and any dependencies between those use cases. The main purpose of a use case diagram is to show what system functions are performed for which actor. Roles of the actors in the system can be depicted.

Registration

Login

View following and followers

View Timeline

View profile

Compose tweet

View Trending news

View User Details

View user Tweets

Classification

Fake Content, URL based spam detect and Trending based
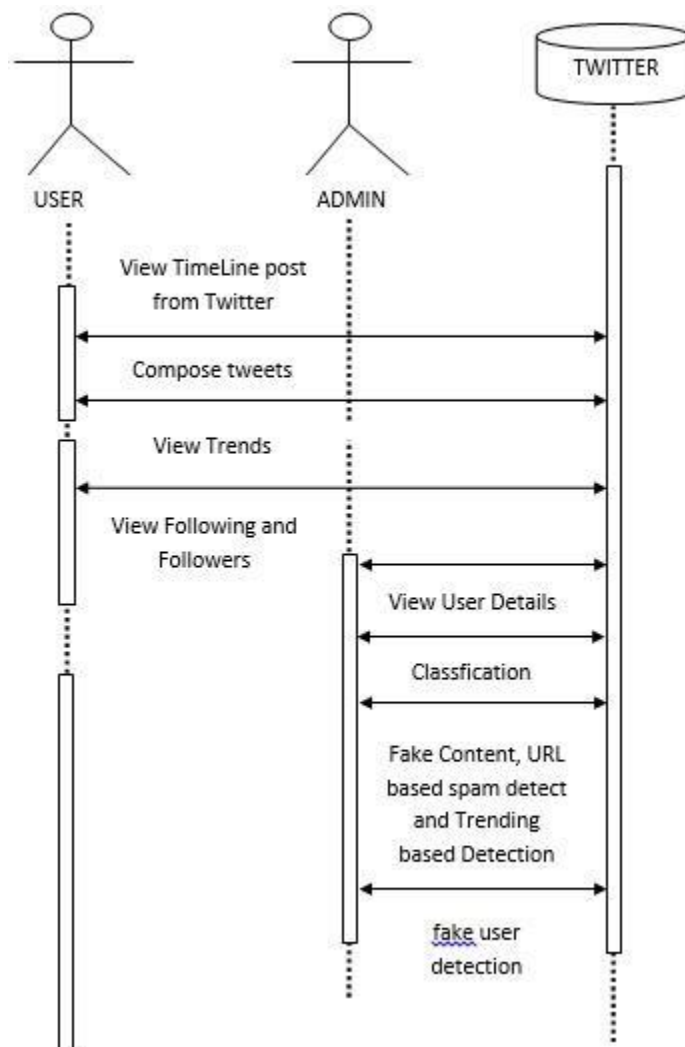
Logout

ADMIN

USER

## 6.2 CLASS DIAGRAM:

In software engineering, a class diagram in the Unified Modeling Language (UML) is a type of static structure diagram that describes the structure of a system by showing the system's classes, their attributes, operations (or methods), and the relationships among the classes. It explains which class contains information.
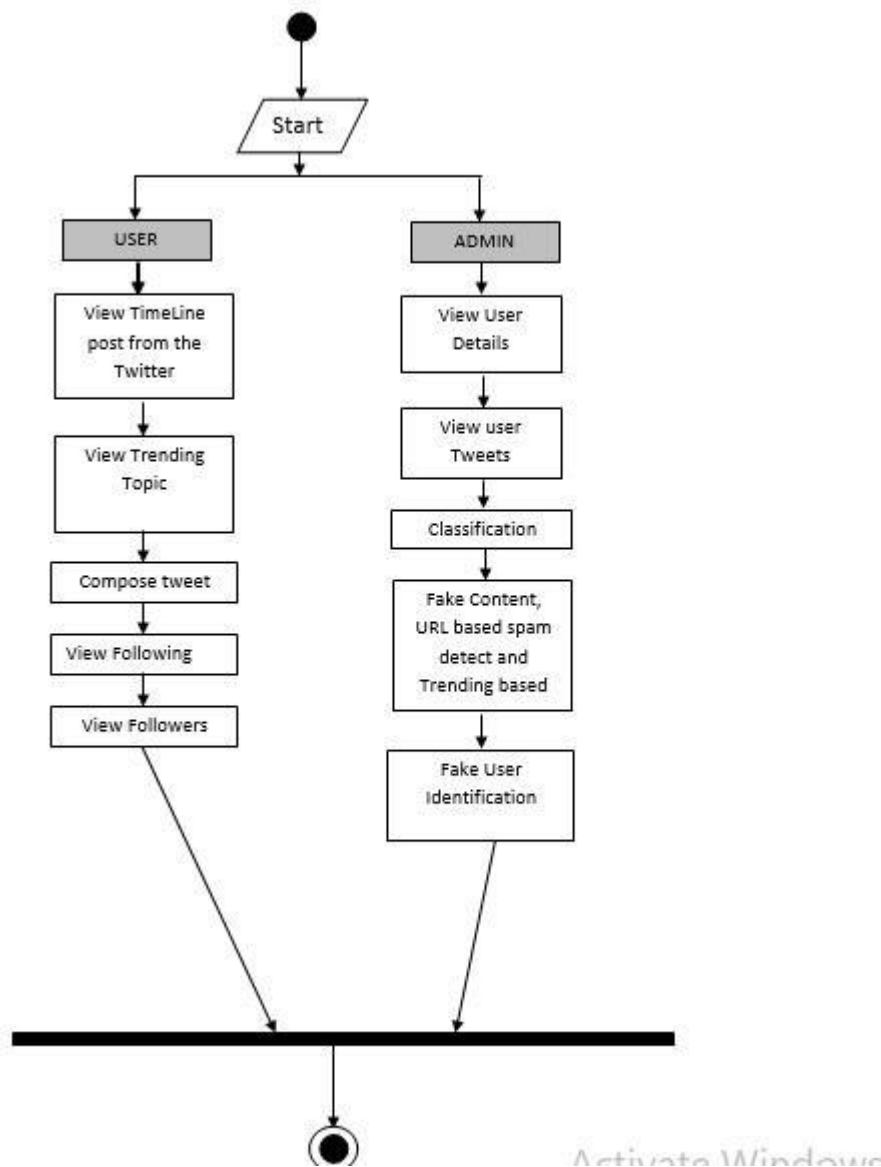
## 6.3 SEQUENCE DIAGRAM:

A sequence diagram in Unified Modeling Language (UML) is a kind of interaction diagram that shows how processes operate with one another and in what order. It is a construct of a Message Sequence Chart. Sequence diagrams are sometimes called event diagrams, event scenarios, and timing diagrams.

## 6.4 ACTIVITY DIAGRAM:

Activity diagrams are graphical representations of workflows of stepwise activities and actions with support for choice, iteration and concurrency. In the Unified Modeling Language, activity diagrams can be used to describe the business and operational step-by-step workflows of components in a system. An activity diagram shows the overall flow of control.

**6.7 SYSTEM TESTING**

The purpose of testing is to discover errors. Testing is the process of trying to discover every conceivable fault or weakness in a work product. It provides a way to check the functionality of components, sub assemblies, assemblies and/or a finished product It is the process of exercising software with the intent of ensuring that the

Software system meets its requirements and user expectations and does not fail in an unacceptable manner. There are various types of test. Each test type addresses a specific testing requirement.

**TYPES OF TESTS**

**UNIT TESTING**

Unit testing involves the design of test cases that validate that the internal program logic is functioning properly, and that program inputs produce valid outputs. All decision branches and internal code flow should be validated. It is the testing of individual software units of the application .it is done after the completion of an individual unit before integration. This is a structural testing, that relies on knowledge of its construction and is invasive. Unit tests perform basic tests at component level and test a specific business process, application, and/or system configuration. Unit tests ensure  that each unique path of a business process performs accurately to the documented specifications and contains clearly defined inputs and expected results.

**Functional test**

 Functional tests provide systematic demonstrations that functions tested are available as specified by the business and technical requirements, system documentation, and user manuals.

Functional testing is centered on the following items:

Valid Input               : identified classes of valid input must be accepted.

Invalid Input          : identified classes of invalid input must be rejected.

Functions              : identified functions must be exercised.

Output                 : identified classes of application outputs must be exercised.

Systems/Procedures: interfacing systems or procedures must be invoked.

Organization and preparation of functional tests is focused on requirements, key functions, or special test cases. In addition, systematic coverage pertaining to identify Business process flows; data fields, predefined processes, and successive processes must be considered for testing. Before functional testing is complete, additional tests are identified and the effective value of current tests is determined.

## System Test

System testing ensures that the entire integrated software system meets requirements. It tests a configuration to ensure known and predictable results. An example of system testing is the configuration oriented system integration test. System testing is based on process descriptions and flows, emphasizing pre-driven process links and integration points.

### 6.7.1 White Box Testing

White Box Testing is a testing in which in which the software tester has knowledge of the inner workings, structure and language of the software, or at least its purpose. It is purpose. It is used to test areas that cannot be reached from a black box level.

### 6.7.2 Black Box Testing

Black Box Testing is testing the software without any knowledge of the inner workings, structure or language of the module being tested. Black box tests, as most other kinds of tests, must be written from a definitive source document, such as specification or requirements document, such as specification or requirements document. It is a testing in which the software

under test is treated, as a black box .you cannot "see" into it. The test provides inputs and responds to outputs without considering how the software works.

# CHAPTER 7

# INTEGRATION TESTING AND TEST RESULTS

## Integration Testing

Software integration testing is the incremental integration testing of two or more integrated software components on a single platform to produce failures caused by interface defects.

The task of the integration test is to check that components or software applications, e.g. components in a software system or – one step up – software applications at the company level – interact without error.

**Test Results:** All the test cases mentioned above passed successfully. No defects encountered.

## Acceptance Testing

User Acceptance Testing is a critical phase of any project and requires significant participation by the end user. It also ensures that the system meets the functional requirements.

**Test Results:** All the test cases mentioned above passed successfully. No defects enc

# CHAPTER 8

# CONCLUSION

In this paper, we performed a review of techniques used for detecting spammers on Twitter. In addition, we also present taxonomy of Twitter spam detection approach es and categorized them as fake content detection, URL based spam detection, spam detection in trending topics, and fake user detection techniques. We also compared the presented techniques based on several features, such as user features, content features, graph features, structure features, and time features. Moreover, the techniques were also compared in terms of their specified goals and datasets used. It is anticipated that the presented review will help researchers find the information on state-of-the-art Twitter spam detection techniques in a consolidated form. Despite the development of efficient and effective approaches for the spam detection and fake user identification on Twitter [34], there are still certain open areas that require considerable attention by the researchers. The issues are briefly highlighted as under :False news identification on social media networks isan issue that needs to be explored because of the serious repercussions of such news at individual as well as collective level [25]. Another associated topic that is worth investigating is the identification of rumor sources on social media. Although a few studies based on statistical methods have already been conducted to detect the sources of rumors, more sophisticated approaches, e.g., social network based approaches, can be applied because of their proven effectiveness.

# REFERENCES

[1] C. Chen, S. Wen, J. Zhang, Y. Xiang, J. Oliver, A. Alelaiwi, and M. M. Hassan, ''Investigating the deceptive information in Twitter spam,'' Future Gener. Comput. Syst., vol. 72, pp. 319–326, Jul. 2017.

[2] I. David, O. S. Siordia, and D. Moctezuma, ''Features combination for the detection of malicious Twitter accounts,'' in Proc. IEEE Int. Autumn Meeting Power, Electron. Comput. (ROPEC), Nov. 2016, pp. 1–6.

[3] M. Babcock, R. A. V. Cox, and S. Kumar, ''Diffusion of pro- and anti-false information tweets: The black panther movie case,'' Comput. Math. Org. Theory, vol. 25, no. 1, pp. 72–84, Mar. 2019.

[4] S. Keretna, A. Hossny, and D. Creighton, ''Recognising user identity in Twitter social networks via text mining,'' in Proc. IEEE Int. Conf. Syst., Man, Cybern., Oct. 2013, pp. 3079–3082.

[5] C. Meda, F. Bisio, P. Gastaldo, and R. Zunino, ''A machine learning approach for Twitter spammers detection,'' in Proc. Int. Carnahan Conf. Secur. Technol. (ICCST), Oct. 2014, pp. 1–6.

[6] W. Chen, C. K. Yeo, C. T. Lau, and B. S. Lee, ''Real-time Twitter content polluter detection based on direct features,'' in Proc. 2nd Int. Conf. Inf. Sci. Secur. (ICISS), Dec. 2015, pp. 1–4.

[7] H. Shen and X. Liu, ''Detecting spammers on Twitter based on content and social interaction,'' in Proc. Int. Conf. Netw. Inf. Syst. Comput., pp. 413–417, Jan. 2015.

[8] G. Jain, M. Sharma, and B. Agarwal, ''Spam detection in social media using convolutional and long short term memory neural network,'' Ann. Math. Artif. Intell., vol. 85, no. 1, pp. 21–44, Jan. 2019.

[9] M. Washha, A. Qaroush, M. Mezghani, and F. Sedes, ''A topic-based hidden Markov model for real-time spam tweets filtering,'' Procedia Comput. Sci., vol. 112, pp. 833–843, Jan. 2017.

[10] F. Pierri and S. Ceri, ''False news on social media: A data-driven survey,'' 2019, arXiv:1902.07539. [Online]. Available: https://arxiv. org/abs/1902.07539

[11] S. Sadiq, Y. Yan, A. Taylor, M.-L. Shyu, S.-C. Chen, and D. Feaster, ''AAFA: Associative affinity factor analysis for bot detection and stance classification in Twitter,'' in Proc. IEEE Int. Conf. Inf. Reuse Integr. (IRI), Aug. 2017, pp. 356–365.

[12] M. U. S. Khan, M. Ali, A. Abbas, S. U. Khan, and A. Y. Zomaya, ''Segregating spammers and unsolicited bloggers from genuine experts on Twitter,'' IEEE Trans. Dependable Secure Comput., vol. 15, no. 4, pp. 551–560, Jul./Aug. 2018.