# Mini Project On

# Whatsapp Chat Analyzer

## By

## Pankaj R. Bhoir (2021510006)
## Haresh A. Gayakhe (2021510015)

Under the guidance of
**Internal Supervisor**

## Prof. Dr. Pooja Raundale

Department of Master Of Computer Application
Sardar Patel Institute of Technology
Autonomous Institute Affiliated to Mumbai University
2021-23

## CERTIFICATE OF APPROVAL

This is to certify that the following students

**Pankaj R. Bhoir (2021510006)**
**Haresh A. Gayakhe (2021510015)**

Have satisfactorily carried out work on the project
entitled

## "Whatsapp Chat Analyzer"

Towards the fulfilment of project, as laid down
by
Sardar Patel Institute of Technology
during year
2021-23.


Project Guide:
Prof. Dr. Pooja Raundale

# PROJECT APPROVAL CERTIFICATE

This is to certify that the following students

**Pankaj R. Bhoir (2021510006)**
**Haresh A. Gayakhe (2021510015)**

Have successfully completed the Project report on

## "Whatsapp Chat Analyzer",

which is found to be satisfactory and is approved

at

SARDAR PATEL INSTITUTE OF TECHNOLOGY,
ANDHERI (W), MUMBAI

INTERNAL EXAMINER                    EXTERNAL EXAMINER

HEAD OF DEPARTMENT                    PRINCIPAL

# Contents

**5 Limitations**     **19**

**6 Future Enhancements**     **19**

**7 User Manual**     **19**

**8 CONCLUSION**     **20**

**9 Bibliography**     **20**

Pankaj R. Bhoir (2021510006)
Haresh A. Gayakhe (2021510015)

Whatsapp Chat Analyzer

# Abstract

The most used and efficient method of communication in recent times is an application called WhatsApp. WhatsApp chats consists of various kinds of conversations held among group of people.

WhatsApp chat analyzer is the application which provide analysis of WhatsApp group chats as well as individual chat. There are various methodologies available for analysis but here matplotlib, streamlit, seaborn, re, pandas libraries of python and some concept of NLP is used. This is the combination of machine learning and NLP.

This whatsapp chat analyzer take import whatsapp chat file from user and analyze it and give different visualizations as a result.

# Objectives

The Web based Application "Whatsapp Chat Analyzer" is used

- To provides a platform to the user which enables user to analyze whatsapp chats

- This application allows user to browse whatsapp exported (.txt) file and import it to WhatsApp chat analyzer and get analysis according to that txt file.

- The user can Analyze by clicking Show Analysis button.

- To visualize data simple and attractive formate

# List of Figures

# List of Tables

Pankaj R. Bhoir (2021510006)
Haresh A. Gayakhe (2021510015)

# 1   Introduction

## 1.1   Problem Definition

WhatsApp-Analyzer is a statistical analysis tool for WhatsApp chats. Working on the chat files that can be exported from WhatsApp it generates various plots showing, for example, which another participant a user responds to the most. We propose to employ dataset manipulation techniques to have a better understanding of WhatsApp chat present in our phones..

## 1.2   Objectives and Scope

### 1.2.1   Objectives

The Web based Application "Whatsapp Chat Analyzer" is

- To provide an in-depth exploratory data analysis on various types of WhatsApp chats.

- Use Streamlit to build a dashboard

- To present an analysis of the WhatsApp group data to as certain the level of involvement and participation by members in that group chat.

- To extract WhatsApp chat data

### 1.2.2   Scope

The user can provide his/her chat details for in-depth analysis of whatsapp chats.

In the web application the user must enter his/her personal and educational details in the profile section.

Our System is being made for reducing the information loss and smoothening the communication between Training and Placement Co-ordinator and students so that they both have all the required information in their hand.

## 1.3 Existing System

In olden days' there is no analysis for whatsapp chat. If someone wants to analyze there is no CSV file to analyze. WhatsApp Application provide export txt file which is in raw format. It is very complicated for analysis. So we have to forget that system and switch over to the WhatsApp Chat Analyzer.
Some of the disadvantages of existing system are as follows :

- Analysis are not accurate.

- Raw data.

- Time consuming.

- Difficult to Analyze.

## 1.4 Proposed System

The "WhatsApp Chat Analyzer" provides a platform to the user which enables user to analyze whatsapp chats online on heroku link. This application allows user to browse whatsapp exported (.txt) file and import it to WhatsApp chat analyzer and get analysis according to that txt file. And user can Analyze by clicking Show Analysis button.
Some of the advantages of our system are as follows :

- Runs on all devices.

- Shows based on whatsapp chat file.

- Shows different visualizations.

- Shows Total Messages.

## 1.5    System Requirements

- Hardware Requirements on Client Side

Table 1.5.1: Hardware Requirements on Client Side

| Device | Any operating system which supports python |
|---|---|
| Processor | Dual Core Processor or Above |
| RAM | Minimum 2 GB RAM |
| Storage | Minimum 10 GB Storage Space |

- Software Requirements on Client Side

Table 1.5.2: Software Requirements on Client Side

| Operating System | Windows/Linux/Mac OS |
|---|---|
| Server | Not Required |

# 2 Software Requirement Specification (SRS) and Design

## 2.1 Purpose

The purpose of our project is to develop an web application that can help user to easily analyze their Whatsapp data and to keep all the required information handy.
This can save lots of efforts of user and it will be easier to analyze the data the each and every single and small detail. This web app will keep track of all the information regarding the whatsapp chats.

## 2.2 Definition

To build a Web Application so the user can have an easy to analyze the data.

## 2.3 Overall Description

### 2.3.1 Product Functions

The product function includes:

1. Run App: Users are required to run web app on any system.

2. Upload data: Then user have to upload the text file of the data which is contains whatsapp chat.

3. View analysis: After clicking on view analysis it shows the detailed analysis of whatsapp chat

# Whatsapp Chat Analyzer

### 2.3.2 User Characteristics

There are only one types of user:

- User: user have to upload their whatsapp chat txt file. then after clicking on view analysis, web app will show the detailed analysis of whatsapp chat. where user can check the analysis according to what deta he/she needed.

# 3 Project Analysis and Design

## 3.1 Methodologies Adapted

In Waterfall model, very less customer interaction is involved during the development of the product. Once the product is ready then only it can be demonstrated to the end users.

Once the product is developed and if any failure occurs then the cost of such issues is very high, because we need to update everything from document till the logic.

3.1.1: Waterfall Model

## 3.2   Modules

### 3.2.1   Activity diagram



3.2.1: Activity Diagram

### 3.2.2   PERT Chart



3.2.2: PERT Chart

### 3.2.3   Gantt Chart



3.2.3: Gantt Chart

# 4 Project Implementation and Testing

## 4.1 Blank web app



4.1.1: Web View

## 4.2 Select chat text file



4.2.1: File upload

## 4.3   Click on show analysis



4.3.1:  Show analysis

## 4.4   Activity map



4.4.1:  Activity histogram

## 4.5 Weekly activity



4.5.1: Heat map

## 4.6 Wordcloud

## 4.7 Common words



4.7.1: Most Common words

## 4.8 Emoji analysis

# Whatsapp Chat Analyzer

## 4.9 Code 1



4.9.1: app.py

## 4.10 Code 2



4.10.1: app.py

## 4.11    Code 3



4.11.1: app.py

## 4.12    Code 4



4.12.1: app.py

## 4.13    Code 5



4.13.1: app.py

## 4.14    Code 6



4.14.1: helper.py

## 4.15    Code 7



4.15.1: helper.py

## 4.16    Code 8



4.16.1: helper.py

## 4.17    Code 9



4.17.1:  helper.py

## 4.18    Code 10



4.18.1:  preprocessor.py

## 4.19    Code 11



4.19.1: preprocessor.py

# 5 Limitations

- Time format should be in 24 hrs format.

- It does not have a feature to view media files it only shows number of files shared.

# 6 Future Enhancements

- Upload both 12hrs and 24hrs time format.

- Abusive words detection.

- Sentiment analysis of users.

# 7 User Manual

- Run the code

- Click on browse files and select Whatsapp chat file on which you want to perform analysis.

- Click on show analysis to display output

# 8   CONCLUSION

The major objective that has been decided in the initial phase of the require-
ment analysis is achieved successfully. After the implementation, the system
provides reliable results. The system is totally menu and user friendly, which
makes it easy for the users even with limited knowledge of computer environ-
ment to operate the developed system. The system avoids the drawbacks of the
existing manual system and the validation facility of the system totally elimi-
nates the chances of wrong data entry. It has following features:

- User friendly.

- Time saving.

- Runs on any devices.

- Analyzes any WhatsApp imported file.

- Accuracy.

- Reliability.

- Easy to use.

# 9   Bibliography

## 9.1   Web References

[1.] `https://streamlit.io/`
[2.] `https://seaborn.pydata.org/`
[3.] `https://pandas.pydata.org/docs/`
[4.] `https://www.geeksforgeeks.org/generating-word-cloud-python/`
[5.] `https://www.tutorialspoint.com/python_text_processing/python_extract_`
`url_from_text.htm`