

# CHAPTER -1

## PROJECT INTRODUCTION, PURPOSE AND SCOPE

---

### INTRODUCTION

This project is based on the data science as we know that data science is the combination of various fields like machine learning, Big data, data visualization, data mining, deep learning etc. This project is based on Big data and data mining. We talk about data structure so in this project we used semi structure data so it regard to big data problems also we used large amount of data. And another point in the project is the data mining as mentioned above it contain so many outliers in the dataset so we refine that outlier and select only that dataset field that is useful for us.

In this project we used dummy dataset that take place from Github repository.

For processing and storage we used hadoop ecosystem it totally complete on hadoop framework. We provide storage on hadoop distributed file system and for processing purpose we used Pig Latin languages. this languages implement those developer that is not good in writing a map reduce code. It just like SQL like language so it is very easy to learn.

We store raw data on hadoop distributed file system. After processing or finalise the result we store that data again in 3 distributed file.

The data set that used in this project is tweets based so we process tweets data. In this project we find out the sentiment of tweets that comes in form of tweets.

We check polarity in this project. Polarity is of three type positive polarity mean positive sentiment in form tweets, negative polarity, and neutral polarity or sometimes is called no polarity in term of sentiment means tweets is form of felling no positive view or negative sentiment in tweets. All mining of data is done with help of Pig Latin language that was developed in 2006 when hadoop comes in market.

For check the polarity of tweets we give rating to the word that part of tweets. Every word has unique rating. Some words have a positive rating or more that 0 numerical value and some have negative rating means less than 0. And some words have 0 rating.

After taking useful data with the help of data mining. We select the tweets attribute and serial number and with the help of Pig Latin language we apply matching or joining operation and calculate the polarity of whole line that comes in tweets format. From here we get three type of numerical values first is positive ( $>0$ ), negative ( $<0$ ) or equal to zero ( $=0$ ). What means numerical values we discuss above?

Finally output result stored on hadoop in three distributed file system one for positive sentiment, second for negative sentiment , another one is 0 sentiment means no sentiment tweets in form of felling type.

## **1.2 PURPOSE OF PROJECT**

The main motive of this project is to find the sentiment and opinion of the tweets and also we can say that we check the quality of tweets. So this type of project is necessary because today is the era of socialization anyone want to connect through social media for taking a more advantage. So there are various problems in existing system so every organisation need such type of person or technology that tell the organisation of their brands what the people thinks about any organisation so that the take decision of their problems so that organisation and enterprise take more and more benefits.

The main purpose of this project is as follows:-

- Determining marketing strategy
- Improve campaign success
- Improve product messaging
- Improve customer services
- Generate leads

### 1.3 SCOPE OF PROJECT

Sentiment analysis is a uniquely powerful tool for businesses that are looking to measure attitudes, feelings and emotions regarding their brand.

The majority of sentiment analysis projects have been conducted almost exclusively by companies and brands through the use of social media data, survey responses and other hubs of user generated content

.

- The future of sentiment analysis is going to continue to dig deeper, far past the surface of the number of likes, comments ,shares and reach, and truly understand significance of social media interaction and what they tell us about the brand.
- Data pre-processing using more parameters to get best sentiment.
- Updating dictionary for new synonyms and antonyms already existing words.
- It helps to analyze the intent like views complain suggestions etc.
- This type of project is useful for natural processing language.
- This type of project is help in Consumer voice
- This project based on rating of tweets so this type of project help to tells the brand reputation checking the sentiment of peoples.
- This type of project is playing a wonderful role in online advertising and online commerce

So this project is best projects related to the data science it also so this type so project is help in politic also for their Voting advise application and Clarification of politicians positions.

This type or sentiment analysis is help in Public Actions like real-world events monitoring, legal matters, and intelligent transaction system.

## **CHAPTER – 2**

### **SOFTWARE AND HARDWARE REQUIREMENT STUDY**

---

#### **INTRODUCTION :-**

This chapter is regards study of hardware and software required during project complementation. As you know this project is regard the problems of Big Data so we need high configuration system hardware.

This project based new technology so here we used new software like Hadoop framework and new languages.

So here we mentioned that technologies, hardware and software-

#### **2.1 HARDWARE REQUIREMENT**

It deal what type of hardware is required to complete the project and what is the system configuration is required is mandatory to install is available because Hadoop framework require red very high configuration Linux system.

##### **2.1.1 LINUX AS OPERATING SYSTEM :**

As we know that Hadoop does not support other operating system except Linux operating system. So it is mandatory to used Linux operating system for Hadoop.

##### **2.1.2 RAM REQUIRED:**

For Hadoop minimum 8 GB RAM is required.

##### **2.1.3 HARDDISK :**

Minimum 1TB hard disk is enough because we are not used Hadoop as distributed manner.

## **2.2 SOFTWARE REQUIRMENT**

We used various tools and languages as per project requirement. So we used here Hadoop framework, Pig, and Bitwise software for data ingestion.

### **2.2.1 HADOOP**

Hadoop is the solution of Big Data problems so we used Hadoop framework. With the help of hadoop we can process, and store the large amount of data.

Hadoop ecosystem contains so many tools to handle this type of problems that discuss above.

For storage Hadoop provide storage for large amount of data in distributed manner through which we handle such type of problems very easily. It provide fast processing system than the other commodity system. It is developed by Dug Cutting and Mike Cafarella in 2005 it the project of yahoo that is Nutch search engine. We used HDFS for storage in ecosystem for raw data storage and also final result or data.

### **2.2.2 PIG**

Pig is the tool or platform that used Pig Latin languages for process the large amount of data. This is developed in 2006 when hadoop comes in market because most of the programmer faces so many problems to written Map Reduce code than Pig Latin language is developed by Yahoo because it is the part of Yahoo project.

### **2.2.3 BITWISE**

This software is used for data ingestion process with help of this software we transfer data from one system to another by taking remotely the both system.

## **CHAPTER – 3**

### **PURPOSED METHODOLOGY**

---

#### **INTRODUCTION**

This chapter is regards with the purposed methodology used during the project starting till the complementation. This project divided into the various phase like data ingestion, data wrangling, refining of data and last phase of this project is data analysis or result analysis. So understand one by one.

#### **3.1 Data Ingestion**

This is the first stage of this project in this phase we take data from Github repository and import that data into hadoop distributed file system.

Data ingestion deals with data import and export process from one database to another database. We stored raw data in to distributed manner through which we can process that data easily and fast as we can see figure no. 2.0

Figure 2.0 show data stored in hadoop distributed file system in that figure one is .text file that is the AFFIN word rating dictionary and another is the excel file that contain demonetization of tweets.

#### **3.2 Data wrangling**

This the second phase of project. This stage comes after the data ingestion in this stage we take data in standard format through which we can handle or process easily when data comes un this format every rows become observation and every column become variable so we can handle data easily.

In this stage we load data from hadoop distributed file system to Pig engine because we use Pig Latin languages to process the data and pig also a platform that stored data temporary for processing purpose and after processing we can stored data into hadoop distributed file system.

### **3.3 Data refining**

This is the third stage of project and this is very important phase of project or we can say that in data science project this is the very complex and important stage because in this stage we deal with the outlier of the data.

Outlier may something that affect our data. In this tweets section there are various type of outlier like in comment most tweets contain braces, some type of special symbol. They type object gives error during the process of data

So in this stage we select tweets part of data and clean it for further usage and this is done by the help of Pig Latin languages.

### **3.4 Data Mining**

This is main phase of the project in the phase we deal with the data mining stage as we know that in Data mining we select that data that is useful for us. Means that is the form of information.

In the phase we select two variables one is text means the column that contain demonetization tweets and process it further as per requirement.

### **3.5 Analysis Methodology**

The is the final stage of data set or this project in the phase we process data as per our requirement like first we find out the average rating of the demonitization tweets by the data that obtained from mining process.

In the divided data or information into three category one for average rating, second is positive rating, and last one is negative rating. Positive rating means sentiment of tweets is positive and negative means negative sentiment and if average rating is equal to the zero means there is no sentiment and we can say that tweets is the polar form means this tweets is the felling type.

Finally we stored process information into three directory in the hadoop.

## CHAPTER – 4

### CODING AND SCREEN SHOTS

---

#### INTRODUCTION

This chapter include all coding part of project because project have several phase like data ingestion, data wrangling, data mining and last one is data analysis of data that comes after the process the raw data or tweets data.

Here we mentioned all code that used for refining of data, mining, and analysis the result with their sentiment. Basically all code written in PIG Latin languages and also code of Hadoop that how to write a data on Hadoop in block format.

All snapshots of project is mentioned in this chapter that defined how HDFS stored data in block format. What is final output, and how namenode stored metadata.

#### 3.1 CODING

/\*Loading the csv file using pig\*/

```
load_tweets = LOAD '/bharat_project_data/demonetization-tweets.csv' USING  
PigStorage(',');
```

/\* Extracting the needed data.... \*/

/\* \$0 represents first field and \$1 represents the second field.... \*/

/\* "id" is the alias name of \$0 and "text" is the alias name of \$1 \*/

```
extract_details = FOREACH load_tweets GENERATE $0 as id,$1 as text;
```

/\* The following characters are considered to be word separators: space, double quote("), coma(,) parenthesis(()), star(\*).The following characters are considered to be word separators: space, double quote("), coma(,) parenthesis(()), star(\*). \*/

/\* TOKENIZE will split the records based on above seperators and give the bag of words.... \*/

/\* FLATTEN will remove the parenthesis like () and {} \*/



/\* Below code of FLATTEN(TOKENIZE(text)) will remove "{}" from the bag of words and tokens i.e. words from each record will get associate with each record as 3rd field in a iterative manner until the tokens will be finished \*/

tokens = foreach extract\_details generate id,text, FLATTEN(TOKENIZE(text))  
As word;

/\* For checking the affect of step 'tokens', use the below commands \*/

tokens\_limit = LIMIT tokens 8;

dump tokens\_limit;

/\* Loading the dictionary which have rating for words \*/

dictionary = LOAD '/input\_files/AFINN.txt' using PigStorage('\t')  
As(word:chararray,rating:int);

/\* 'replicated' join is the special type of join in which second relation is small enough to fit into the main memory which will help in efficient join.... \*/

word\_rating = join tokens by word left outer, dictionary by word using  
'replicated';

/\* This will give sample data i.e. first 3 rows of word\_rating \*/

word\_rating\_limit = limit word\_rating 3;

dump word\_rating\_limit;

/\* We are iteratively selecting the data from the relation "word\_rating" and selecting fields from respective relations \*/

/\* To select the field which is the part of a certain relation we need to use double colon "::" as below in case of some operation has

already been applied between two relations like in previous steps \*/

rating = foreach word\_rating generate tokens::id as id,tokens::text as text,  
dictionary::rating as rate;

/\* This will give sample data i.e. first 3 rows of rating \*/

rating\_limit = limit rating 3;

/\* This below command will group relations "rating" on the basis of id and text combinedly \*/

word\_group = group rating by (id,text);

/\* First 2 rows of word\_group \*/

word\_group\_limit = limit word\_group 2;

/\* This line is crucial in deciding the rating for any tweet on twitter data as it will sum up all the word's rating in a tweet \*/

/\* Below line will generate group which is from relation word\_group and as all the data will get group on the basis of (id,text) we can perform average of those tokens also as per their ratings \*/

avg\_rate = foreach word\_group generate group, AVG(rating.rate) as tweet\_rating;

/\* First 100 rows of avg\_rate \*/

avg\_rate\_limit = limit avg\_rate 100;

dump avg\_rate;

/\* Filter the positive tweets\*/

positive\_tweets = filter avg\_rate by tweet\_rating >= 0;

/\* Filter the negative tweets\*/

negative\_tweets = filter avg\_rate by tweet\_rating < 0;

/\* Storing the positive and negative tweets output in output\_files folder\*/

store positive\_tweets into '/project\_files/positive\_tweets\_output1';

store negative\_tweets into '/project\_files/negative\_tweets\_output1';

NOTE : The final output data stored in HDFS in Hadoop framework . as we know that Hadoop stored data in block format so final output data is in two blocks but after filtering data it stored in 3 different file or folder format as given below.

1<sup>st</sup> - for positive tweets

2<sup>nd</sup> – for neutral tweets

3<sup>rd</sup> – for negative tweets

So we make three Excel file or sentiment of tweets one for Negative tweets, neutral or feeling and other Excel file contained positive sentiment of tweets data.

## 3.2 SCREEN SHOTS

This contain all the screen shots of project . Fig1.0 is contain raw data that used for project

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T
	text	favorited	favoriteCount	replyToScreenName	created	truncated	replyToId	replyToUi	statusSource	screenName	retweetCount	isRetweeted	retweeted						
1	RT @rsshur	FALSE	0	NA	#####	FALSE	NA	8.01E+17	NA	<a href="https://twitter.com/HASHTAGI	331	TRUE	FALSE						
2	RT @Hemant_80	FALSE	0	NA	#####	FALSE	NA	8.01E+17	NA	<a href="https://twitter.com/PRAMODK	66	TRUE	FALSE						
3	RT	FALSE	0	NA	#####	FALSE	NA	8.01E+17	NA	<a href="https://twitter.com/rahulja130	12	TRUE	FALSE						
4	RT @ANI	FALSE	0	NA	#####	FALSE	NA	8.01E+17	NA	<a href="https://twitter.com/deeptiyv	338	TRUE	FALSE						
5	RT @satis	FALSE	0	NA	#####	FALSE	NA	8.01E+17	NA	<a href="https://twitter.com/CPIMBadli	120	TRUE	FALSE						
6	@DerekScis	FALSE	0	DerekScis	#####	FALSE	NA	8.01E+17	2.59E+09	<a href="https://twitter.com/ambazaari	0	FALSE	FALSE						
7	RT @gaauri	FALSE	0	NA	#####	FALSE	NA	8.01E+17	NA	<a href="https://twitter.com/bhodia1	637	TRUE	FALSE						
8	RT	FALSE	0	NA	#####	FALSE	NA	8.01E+17	NA	<a href="https://twitter.com/KARUNAS	112	TRUE	FALSE						
9	RT	FALSE	0	NA	#####	FALSE	NA	8.01E+17	NA	<a href="https://twitter.com/sumitbhat	1	TRUE	FALSE						
10	National r	FALSE	0	NA	#####	TRUE	NA	8.01E+17	NA	<a href="https://twitter.com/HelpIndia	0	FALSE	FALSE						
11	Many	FALSE	1	NA	#####	FALSE	NA	8.01E+17	NA	<a href="https://twitter.com/sumitbhat	1	FALSE	FALSE						
12	RT @Joydi	FALSE	0	NA	#####	FALSE	NA	8.01E+17	NA	<a href="https://twitter.com/MonishGa	120	TRUE	FALSE						
13	@Jaggesh	FALSE	0	Jaggesh2	#####	FALSE	8.01E+17	8.01E+17	1.23E+09	<a href="https://twitter.com/yuvaraj_k	0	FALSE	FALSE						
14	RT	FALSE	0	NA	#####	FALSE	NA	8.01E+17	NA	<a href="https://twitter.com/PMKejri	45	TRUE	FALSE						
15	RT @sona	FALSE	0	NA	#####	FALSE	NA	8.01E+17	NA	<a href="https://twitter.com/hkgupta16	50	TRUE	FALSE						
16	RT @Dipar	FALSE	0	NA	#####	FALSE	NA	8.01E+17	NA	<a href="https://twitter.com/aazaadpar	45	TRUE	FALSE						
17	RT	FALSE	0	NA	#####	FALSE	NA	8.01E+17	NA	<a href="https://twitter.com/darkdestir	12	TRUE	FALSE						
18	RT	FALSE	0	NA	#####	FALSE	NA	8.01E+17	NA	<a href="https://twitter.com/snooveme	95	TRUE	FALSE						
19	RT @pGur	FALSE	0	NA	#####	FALSE	NA	8.01E+17	NA	<a href="https://twitter.com/Vishwaam	76	TRUE	FALSE						
20	RT	FALSE	0	NA	#####	FALSE	NA	8.01E+17	NA	<a href="https://twitter.com/PoliticalC	12	TRUE	FALSE						
21	RT @Hemant_80	FALSE	0	NA	#####	FALSE	NA	8.01E+17	NA	<a href="https://twitter.com/MdShuaib	66	TRUE	FALSE						
22	RT	FALSE	0	NA	#####	FALSE	NA	8.01E+17	NA	<a href="https://twitter.com/BharatPar	12	TRUE	FALSE						

**Fig.1.0 (Contain raw data)**

This file contain various type attribute in data set. we take only 2 attribute dataset. Here first column is serial number and another filed is of text type that contain tweets of dataset

This figure1.2 contains the Dictionary that used for find out the sentiment form data. This figure contains words with their rating. Every word have own rating some word have negative rating, some have positive rating and some have 0 rating.

abandon	-2abandoned	-2abandons	-2abducted	-2abduction	-2abductions
-2abhor	-3abhorred	-3abhorrent	-3abhors	-3abilities2ability	2aboard
1absentee	-1absolve	2absolved2absolves	2absolving	2absorbed	1abuse
-3abused	-3abuses	-3abusive	-3accept	1accepted	1accepting
1accepts	1accident	-2accidental	-2	accidentally	-2accidents
-2accomplish	2accomplished	2accomplishes	2	accusation	-2accusations
-2accuse	-2accused	-2accuses	-2accusing	-2ache	-2
achievable	1aching	-2acquit	2acquits	2acquitted	2acquitting
2acrimonious	-3active	1adequate	1admire	3admired	3admires
3admiring	3admit	-1admits	-1	admitted	-1admonish
-2admonished	-2adopt	1adopts	1adorable	3adore	3adored
3	adores	3advanced	1advantage	2advantages	2adventure
2adventures	2adventurous	2affected	-1affection	3affectionate	3afflicted
-1affronted	-1	afraid	-2aggravate	-2aggravated	-2aggravates
-2aggravating	-2aggression	-2aggressions	-2aggressive	-2aghost	-2agog
2agonise	-3agonised	-3agonises	-3agonising	-3agree	1
agreeable	2agreed	1agreement	1agrees	1alarm	-2alarmed
-2alarmist	-2	alarmists	-2alas	-1alert	-1alienation
-2alive	1allergic	-2allow	1alone	-2amaze	2
amazed	2amazes	2amazing	4ambitious	2ambivalent	-1amuse
3amused	3amusement	3amusements	3anger	-3angers	-3angry
-3anguish	-3anguished	-3animosity	-2	annoy	-2annoyance
-2annoyed	-2annoying	-2annoys	-2antagonistic	-2	anti
-1anticipation	1anxiety	-2anxious	-2apathetic	-3apathy	-3apeshit
-3apocalyptic	-2apologise	-1apologised	-1apologises	-1apologising	-1apologize
-1	apologized	-1apologizes	-1apologizing	-1apology	-1appalled
-2appalling	-2appease	2appeased	2appeases	2appeasing	2applaud
2applauded					

Fig 2(Words dictionary)

The word have rating 0 is the type of felling so thier is no sentiment in the that word, the word have positive rating means good sentiment regards the tweets and negative rating means negative sentiment towards tweets.

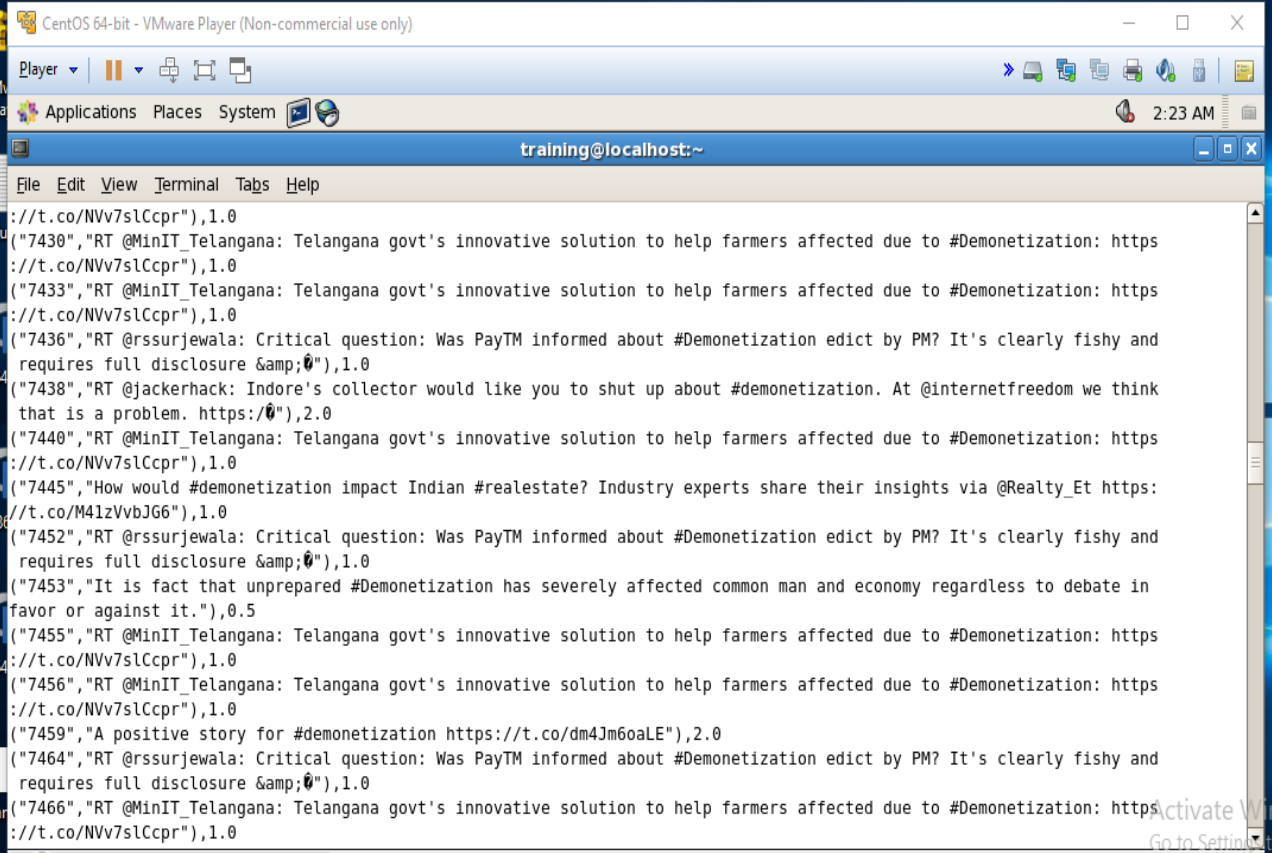
```

[training@localhost ~]$ hadoop fs -ls /project
Found 2 items
-rw-r--r-- 1 training supergroup 28093 2018-09-28 22:00 /project/AFINN.txt
-rw-r--r-- 1 training supergroup 2670787 2018-09-28 10:04 /project/demonetization-tweets.csv
[training@localhost ~]$

```

Fig 3.0 (raw data)

This figure contain raw data of the project that stored on the hadoop distributed file system. First file contain words dictionary and second file contain raw data or demonetization of tweets dataset.



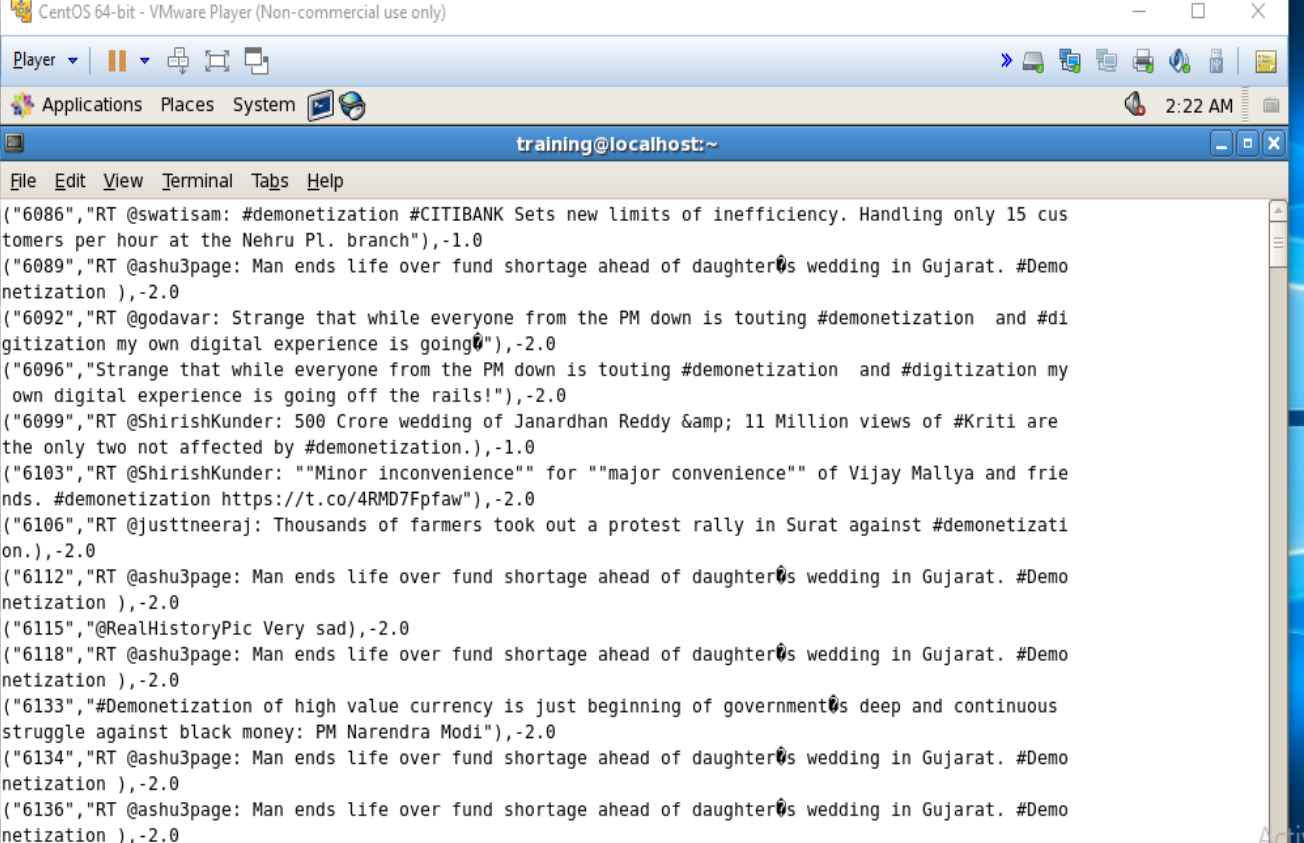
The image shows a terminal window titled "CentOS 64-bit - VMware Player (Non-commercial use only)". The terminal displays a list of tweets with their corresponding sentiment ratings. The tweets are formatted as JSON-like strings: ("tweet\_id", "tweet\_text", "rating"). The ratings are numerical values ranging from 0.0 to 2.0. The tweets discuss the demonetization of Indian currency and its impact on farmers and the economy. The terminal window has a menu bar with "File", "Edit", "View", "Terminal", "Tabs", and "Help". The status bar at the bottom shows "training@localhost:~".

```
://t.co/NVv7slCcpr"),1.0
("7430","RT @MinIT_Telangana: Telangana govt's innovative solution to help farmers affected due to #Demonetization: https://t.co/NVv7slCcpr"),1.0
("7433","RT @MinIT_Telangana: Telangana govt's innovative solution to help farmers affected due to #Demonetization: https://t.co/NVv7slCcpr"),1.0
("7436","RT @rssurjewala: Critical question: Was PayTM informed about #Demonetization edict by PM? It's clearly fishy and requires full disclosure &#x201c;"),1.0
("7438","RT @jackerhack: Indore's collector would like you to shut up about #demonetization. At @internetfreedom we think that is a problem. https://"),2.0
("7440","RT @MinIT_Telangana: Telangana govt's innovative solution to help farmers affected due to #Demonetization: https://t.co/NVv7slCcpr"),1.0
("7445","How would #demonetization impact Indian #realestate? Industry experts share their insights via @Realty_Et https://t.co/M41zVvbJG6"),1.0
("7452","RT @rssurjewala: Critical question: Was PayTM informed about #Demonetization edict by PM? It's clearly fishy and requires full disclosure &#x201c;"),1.0
("7453","It is fact that unprepared #Demonetization has severely affected common man and economy regardless to debate in favor or against it."),0.5
("7455","RT @MinIT_Telangana: Telangana govt's innovative solution to help farmers affected due to #Demonetization: https://t.co/NVv7slCcpr"),1.0
("7456","RT @MinIT_Telangana: Telangana govt's innovative solution to help farmers affected due to #Demonetization: https://t.co/NVv7slCcpr"),1.0
("7459","A positive story for #demonetization https://t.co/dm4Jm6oaLE"),2.0
("7464","RT @rssurjewala: Critical question: Was PayTM informed about #Demonetization edict by PM? It's clearly fishy and requires full disclosure &#x201c;"),1.0
("7466","RT @MinIT_Telangana: Telangana govt's innovative solution to help farmers affected due to #Demonetization: https://t.co/NVv7slCcpr"),1.0
```

**Fig 4.0(Positive rating)**

This image or figure2.1 contain positive rating dataset that comes after processing the dataset.

This figure 2.2 contain the negative rating of tweets or we can say negative sentiment of tweets



The screenshot shows a terminal window titled 'training@localhost:~' within a VMware Player. The terminal displays a list of tweets and their corresponding negative sentiment ratings. The ratings are either -1.0 or -2.0. The tweets are related to the demonetization of Indian currency and the wedding of Janardhan Reddy.

```

("6086","RT @swatisam: #demonetization #CITIBANK Sets new limits of inefficiency. Handling only 15 customers per hour at the Nehru Pl. branch"),-1.0
("6089","RT @ashu3page: Man ends life over fund shortage ahead of daughter's wedding in Gujarat. #Demonetization ),-2.0
("6092","RT @godavar: Strange that while everyone from the PM down is touting #demonetization and #digitization my own digital experience is going"),-2.0
("6096","Strange that while everyone from the PM down is touting #demonetization and #digitization my own digital experience is going off the rails!"),-2.0
("6099","RT @ShirishKunder: 500 Crore wedding of Janardhan Reddy & 11 Million views of #Kriti are the only two not affected by #demonetization.),-1.0
("6103","RT @ShirishKunder: "Minor inconvenience" for "major convenience" of Vijay Mallya and friends. #demonetization https://t.co/4RMD7Fpfaw"),-2.0
("6106","RT @justtneeraj: Thousands of farmers took out a protest rally in Surat against #demonetization.),-2.0
("6112","RT @ashu3page: Man ends life over fund shortage ahead of daughter's wedding in Gujarat. #Demonetization ),-2.0
("6115","@RealHistoryPic Very sad),-2.0
("6118","RT @ashu3page: Man ends life over fund shortage ahead of daughter's wedding in Gujarat. #Demonetization ),-2.0
("6133","#Demonetization of high value currency is just beginning of government's deep and continuous struggle against black money: PM Narendra Modi"),-2.0
("6134","RT @ashu3page: Man ends life over fund shortage ahead of daughter's wedding in Gujarat. #Demonetization ),-2.0
("6136","RT @ashu3page: Man ends life over fund shortage ahead of daughter's wedding in Gujarat. #Demonetization ),-2.0

```

**Fig 5.0(Negative Rating)**

These all are the images of raw dataset and final result of dataset that obtained after processing.

## CHAPTER – 5

### CONCLUSION AND REFERENCES

---

#### 5.1 CONCLUSION

This project is totally based on the Big data problems that Is the sub part of Data Science so this project is very good for any type of organization to help the understand the sentiment of their customers and users.

For doing this project there is another technologies that current comes in market is mandatory to learn because there is required much or deep knowledge of this area to complete such type of project. So for this project we need good knowledge of Hadoop framework.

For this project we are used Pig Latin Language for both processing the data.

We provide storage on Hadoop distributed file system.and analyze the result with Pig Latin. We separate the tweets on the based on the their polarity.

Polarity in the sense their tweets quality of tweets whether tweet is negative, positive and and polarity is neutral(feeling type).

So this project is very help to understand the quality of tweets so it help to understand the sentiment of the people. It is the project of Big data problem so we are lots of think we learn.

In this project we take data from Github repository. Data is dummy because this project is learning purpose we stored that data on hadoop distributed file system to handle the data because data is semi structure so hadoop provide best solution for this type of problems.

Second phase of this project is the data cleaning or refining of data in the phase we clean the data as per our requirement and take off only that data are useful for us.

Last phase of this project is data analysis. In this stage we further store data(or clean data) on the hadoop distributed file system for further.

In this phase we checks the polarity of the sentiment of the comment. We define polarity term in chapter 2. Polarity are three type negative, positive and neutral(feelings) .and final stored result of hadoop distributed file system.



## 5.2 BIBLIOGRAPHY

Chuck Lam, Hadoop in Action,2010, [www.manning.com/HadoopinAction](http://www.manning.com/HadoopinAction).

Hadley Wicham , R for Data Science,2016, [www.r4ds.had.co.nz](http://www.r4ds.had.co.nz)

Allen B. Downey, Think Stats,2011, [http://bit.ly/think\\_stats\\_2e](http://bit.ly/think_stats_2e).

Wikipedia.