



iJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 9 Issue: VII Month of publication: July 2021

DOI: <https://doi.org/10.22214/ijraset.2021.37193>

www.ijraset.com

Call:  08813907089

E-mail ID: ijraset@gmail.com

Hadoop Based Generic Template for Performing Sentiment Analysis Using Apache PIG

Gadige Vishal Sai¹, Chukka Nikhil², Panyala Harsha Vardhan Reddy³, Sai Harsha Bandurupally⁴

¹Software Development Engineer, NCR Corporation, Hyderabad

²Application Developer, Accenture, Hyderabad

³Associate Software Engineer, CGI, Hyderabad

⁴Full Stack Engineering Analyst, Accenture, Hyderabad

Abstract: Every day over 2.5 quintillion data is generated using various channels like online surveys, transactional data tracking, social media monitoring, etc. Out of these majority of the data is generated using social media platforms. This raw data contains information that can be used for industrial, economic, social and business purposes. To facilitate this, sentiment analysis has become a prospect for various tech-based industry giants to review and analyze their products. Hadoop has been established as one of the best tools for storing, processing, and streaming data in the market. In this paper, we present a generic approach to performing sentiment analysis using Apache PIG which classifies the given data taken from a dataset to either positive or negative to get the people's sentiment over an object or an issue.

Keywords: Big Data Analysis, Hadoop, Apache PIG, Sentiment Analysis, Opinion Mining, Business Analysis

I. INTRODUCTION

Since the dawn of time up until 2005, humans have created nearly 130 exabytes of data. Nevertheless, this value has been significantly increased to 9,100 and later to 40,900 during 2005 to 2015 and 2015 to 2020. This became a massive challenge for the World Wide Web to understand and satisfy the requirements of people. At present, people are comprehensively conveying their thoughts over discussion forms, online blogs and social networking applications. Industrial organizations have understood this and applied various techniques to study the people's sentiment over their products, which help them improve their productivity by predicting their products' success rate. However, to perform this, large scale data storing and processing needs to be done. To answer this, Hadoop has been recommended as one of the best tools for data storing and processing techniques.

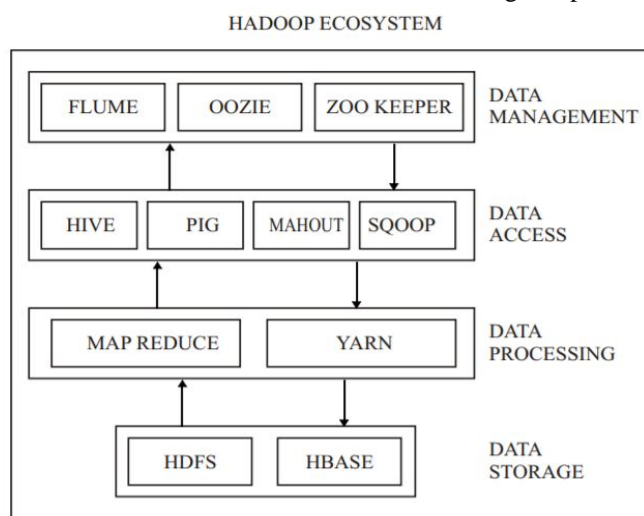


Figure 1. Hadoop Ecosystem

The above figure clearly illustrates the various tools that are available in the Hadoop ecosystem. Hadoop consists of various storing and processing components that can handle all forms of data – structured, semi-structured, quasi structured, and unstructured. In this paper, we use Hadoop's processing component Pig to create a generic template to perform sentiment analysis on a pre-processed data set.

II. RELATED WORK

Big Data analytics is the investigative process that includes recording, storing, analyzing, unsheathing helpful information from the data. It uses different scientific methods, algorithms and techniques to analyze the data effectively and efficiently; data can be either structured or unstructured. Many industries like banking, finance, e-commerce, manufacturing, etc. use data science to manage gigantic data from multiple sources to derive valuable information and make smarter data-driven decisions.

Ajit Noonia and team performed sentimental analysis on Twitter using PIG and HIVE [1], where they collected over 5000 tweets, performing pre-processing over it. They performed sentimental analysis so that they can differentiate between positive and negative tweets. In this paper, researchers have changed the data from unstructured data to structured data to perform pre-processing. After pre-processing, they used tokenization to convert text into tokens before transforming it, and data comparison is made using a data dictionary.

Another related work includes Opinion Mining of Twitter Data using Hadoop and Apache Pig [2]. It proposes a method of using a dictionary-based approach for opinion mining. With the help of pig statements, they have analyzed Twitter data to check the tweets' polarity based on polarity dictionary to separate tweets into positive and negative tweets (opinions). This paper has analyzed the Twitter data using some visualizing techniques to analyze the results effectively.

III.METHODOLOGY

All paragraphs must be indented. All paragraphs must be justified, i.e. both left-justified and right-justified.

In this paper, we propose a generic template which facilitates the user to perform sentiment analysis by implementing the following methodology.

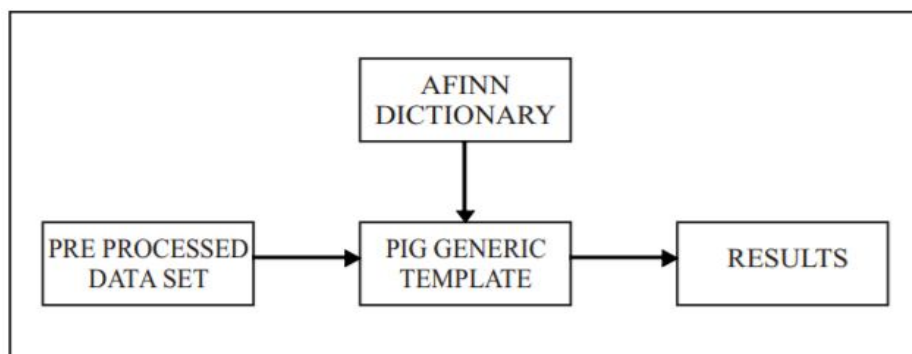


Figure 2. The above diagram illustrates the proposed methodology

A. Preprocessed Dataset

In this proposed experiment, we perform sentiment analysis on a preprocessed dataset. The process of performing preprocessing is outside of the scope of this experiment. It should be done in such a way that all the types of data – structured, semi-structured, quasi structured, and unstructured should be preprocessed to a set of English language statements. This set of sentences should be fed as input to our proposed template, where it helps the user to perform sentiment analysis along with the help of the AFINN dictionary.

B. AFINN Dictionary

AFINN dictionary is an open-source distributed lexicon that contains a list of English words, each associated with a corresponding integer value. Finn Arup Nielsen developed it. The dictionary consists of 2,477 coded words, each having a value score between –5 and +5. A negative score implies that the word is associated negatively, and an upbeat score implies that the word is associated with a positive context. By mapping the word score values to the input statements, we perform sentiment analysis.

absorbed	1
abuse	-3
abused	-3
abuses	-3
abusive	-3
accept	1
accepted	1

Figure 3. Screenshot containing data in AFINN dictionary

C. Proposed PIG Generic Template

The sentiment analysis template consists of a series of PIG commands that uses both the input set of statements and the AFINN dictionary. The sequence of statements executed are shown below using the flowchart:

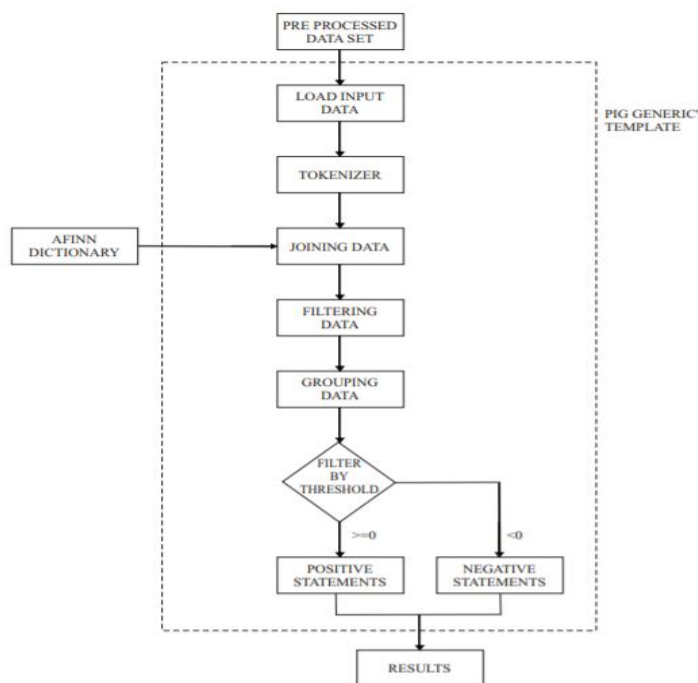


Figure 4. The above figure illustrates a detailed step by step implementation of the proposed generic template

Generally, PIG applications can be executed in two modes – local mode and HDFS (Hadoop Distributed File System) mode. Usually, the local mode is used for software development, whereas the HDFS mode is used for performing software testing. Here, we use the local mode for executing our PIG script. All PIG applications are written in PIG Latin language that can be compiled in the grunt shell.

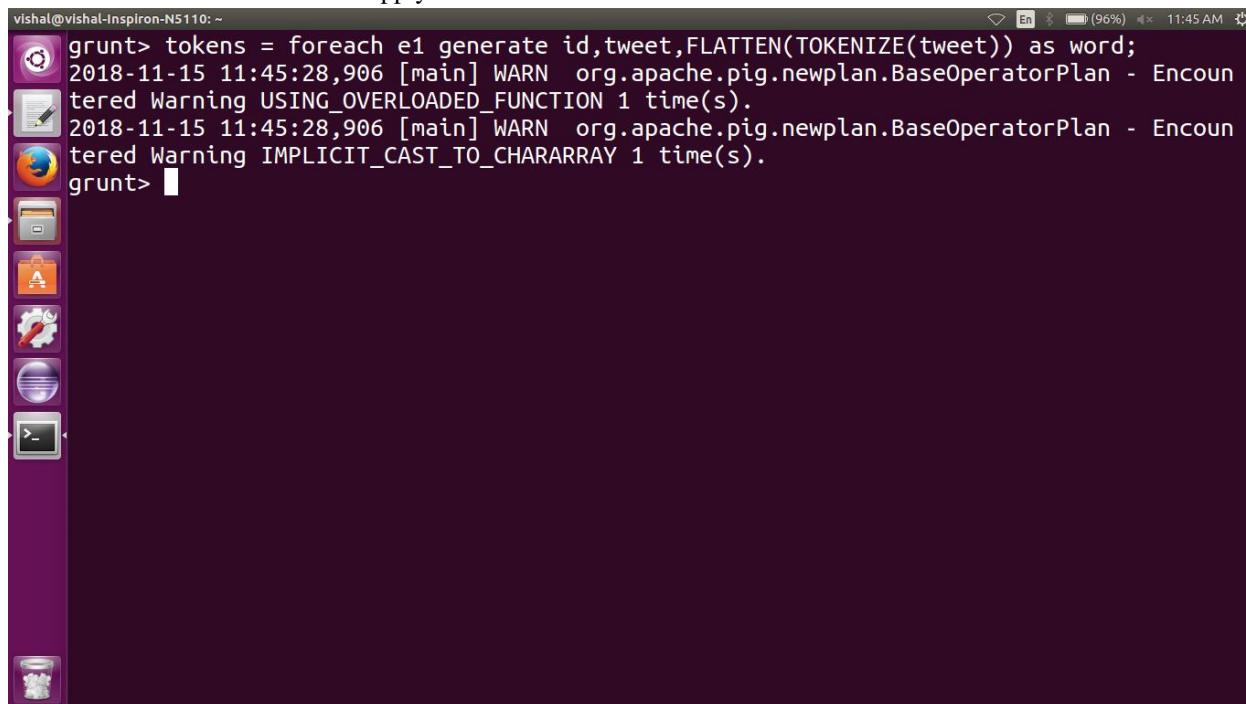
Here, we used the dataset related to the people's opinion on demonetization in India for testing the application.

- 1) **Load Input Data:** Here we use the PIG Latin's load command to load the data onto the PIG local system. Using the PigStorage() function, we perform the load operation from a given local path. We can use the dump command to view the loaded data, which prints the loaded data on the grunt shell. The dump command is like a select query in RDBMS. We can also use the describe command, which helps us view the structure of the requested data.



Figure 5. Screenshot demonstrating the load command

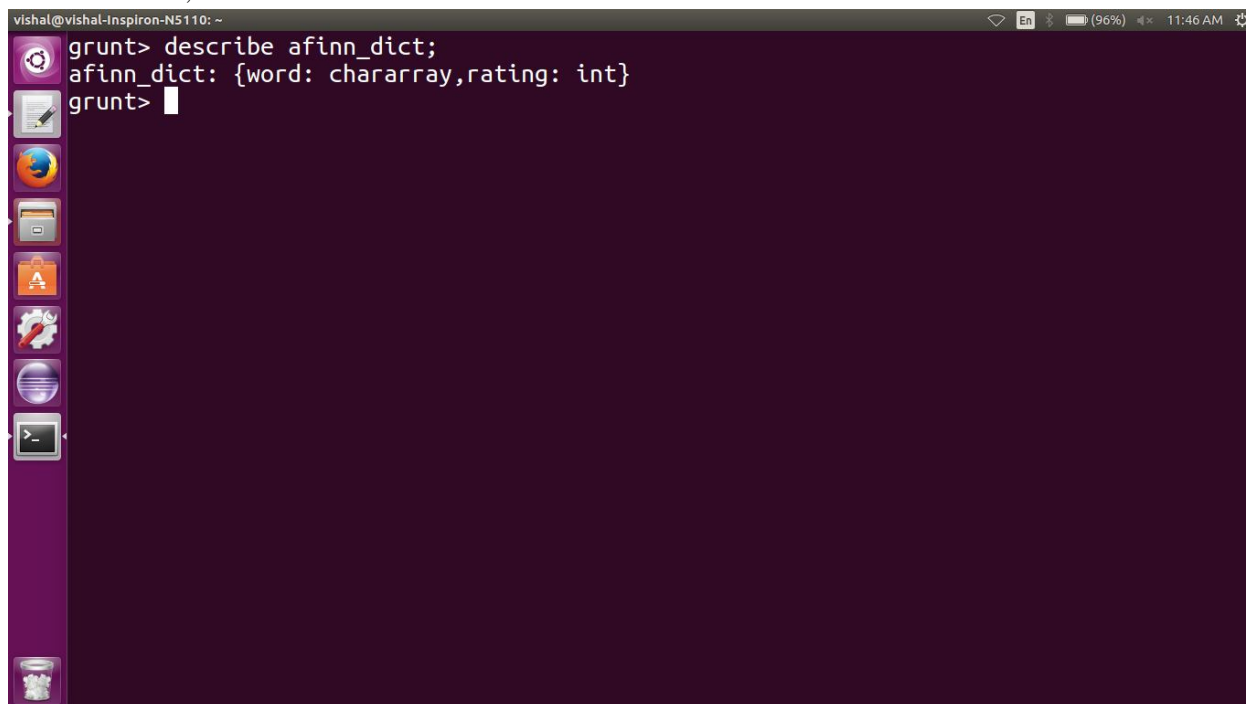
- 2) *Tokenize the Data*: A sentence is a collection of words. The meaning of the words combined represents the entire meaning of the sentence. To know the sentence's sentiment, we break down each sentence into tokens of the format word and word count. Later we use the flatten command to apply it to all the available sentences. The structure of the data looks as follows:



```
vishal@vishal-Inspiron-N5110: ~
grunt> tokens = foreach e1 generate id,tweet,FLATTEN(TOKENIZE(tweet)) as word;
2018-11-15 11:45:28,906 [main] WARN org.apache.pig.newplan.BaseOperatorPlan - Encoun
tered Warning USING_OVERLOADED_FUNCTION 1 time(s).
2018-11-15 11:45:28,906 [main] WARN org.apache.pig.newplan.BaseOperatorPlan - Encoun
tered Warning IMPLICIT_CAST_TO_CHARARRAY 1 time(s).
grunt>
```

Figure 6. Screenshot demonstrating the process of tokenization

- 3) *Load AFINN dictionary*: We use the same load command to load the AFINN dictionary to PIG local system. Also, to view the structure of the data, we can use the describe command.



```
vishal@vishal-Inspiron-N5110: ~
grunt> describe afinn_dict;
afinn_dict: {word: chararray,rating: int}
grunt>
```

Figure 7. Screenshot demonstrating the structure of AFINN dictionary

- 4) *Joining the Data*: To perform sentiment analysis, we need all the data to merge into a single large chunk of data. To achieve this, we use the concept of joins. Specifically, we use the left outer join as we are more concerned about the words present in both the AFINN dictionary and the data set.



```

vishal@vishal-Inspiron-N5110: ~
grunt> join_data_dict = join tokens by word left outer, afinn_dict by word using 'replicated';
2018-11-15 11:47:17,788 [main] WARN org.apache.pig.newplan.BaseOperatorPlan - Encountered Warning USING_OVERLOADED_FUNCTION 1 time(s).
2018-11-15 11:47:17,788 [main] WARN org.apache.pig.newplan.BaseOperatorPlan - Encountered Warning IMPLICIT_CAST_TO_CHARARRAY 1 time(s).
grunt>

```

Figure 8. Screenshot showing the join operation on the data

- 5) *Filtering the Data*: Here we use the filter command, which helps us filter out the data columns essential for the analysis. Here, we can also apply the describe command to view the structure of the filtered output.



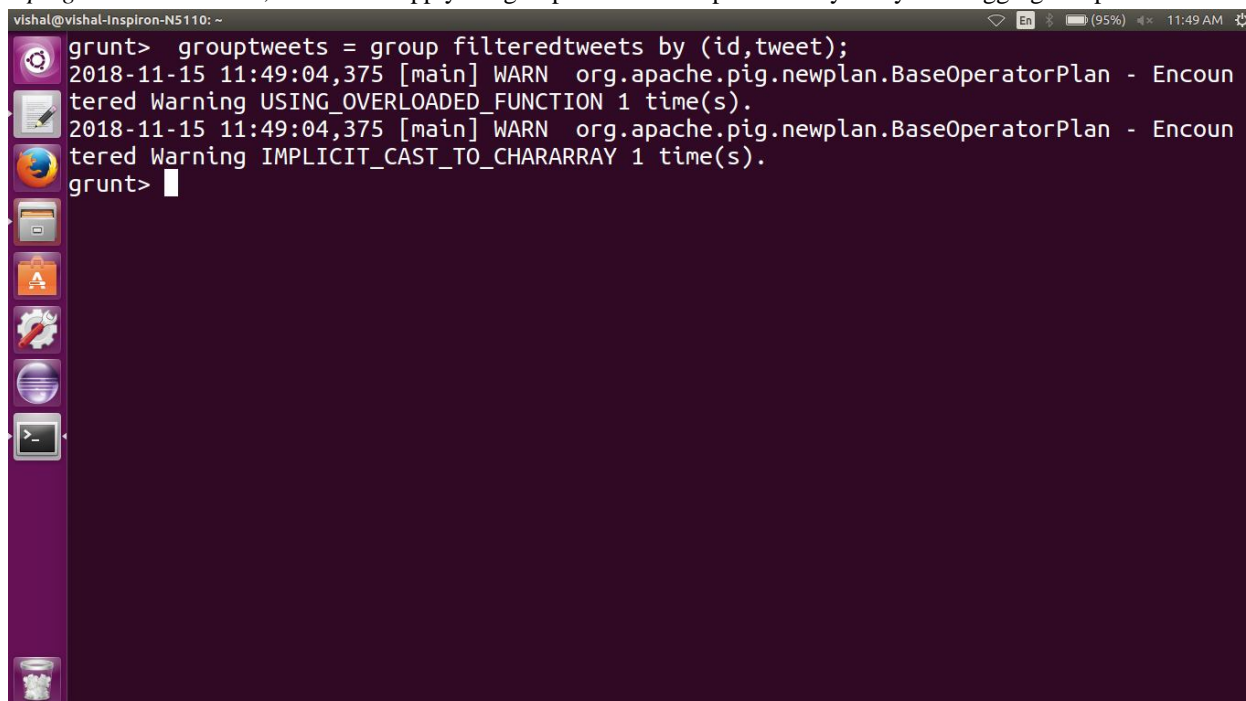
```

vishal@vishal-Inspiron-N5110: ~
grunt> filteredtweets = foreach join_data_dict generate tokens::id as id,tokens::tweet as tweet,afinn_dict::rating as rating;
2018-11-15 11:48:11,118 [main] WARN org.apache.pig.newplan.BaseOperatorPlan - Encountered Warning USING_OVERLOADED_FUNCTION 1 time(s).
2018-11-15 11:48:11,118 [main] WARN org.apache.pig.newplan.BaseOperatorPlan - Encountered Warning IMPLICIT_CAST_TO_CHARARRAY 1 time(s).
grunt>

```

Figure 9. Screenshot showing the filtering process

- 6) *Grouping the Data:* In PIG, we need to apply the group command to perform any analysis or aggregate operation.

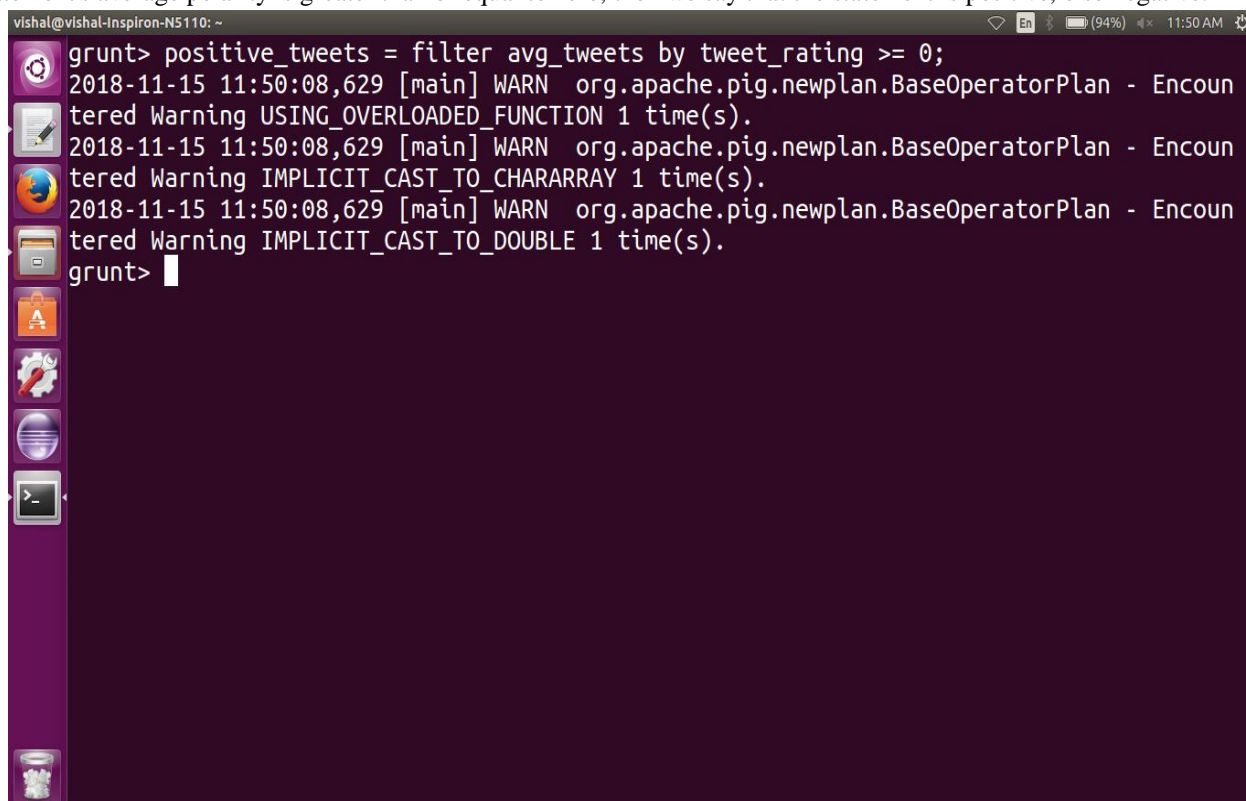


```

vishal@vishal-Inspiron-N5110: ~
grunt> grouptweets = group filteredtweets by (id,tweet);
2018-11-15 11:49:04,375 [main] WARN org.apache.pig.newplan.BaseOperatorPlan - Encoun
tered Warning USING_OVERLOADED_FUNCTION 1 time(s).
2018-11-15 11:49:04,375 [main] WARN org.apache.pig.newplan.BaseOperatorPlan - Encoun
tered Warning IMPLICIT_CAST_TO_CHARARRAY 1 time(s).
grunt>
  
```

Figure 10. Screenshot showing the execution of group command

- 7) *Filter by threshold:* The main advantage of the AFINN dictionary is that all the negative words have a negative polarity, and all the positive words have a positive polarity. Here we use the same decision logic to classify data into positive or negative. If the statement's average polarity is greater than or equal to zero, then we say that the statement is positive, else negative.



```

vishal@vishal-Inspiron-N5110: ~
grunt> positive_tweets = filter avg_tweets by tweet_rating >= 0;
2018-11-15 11:50:08,629 [main] WARN org.apache.pig.newplan.BaseOperatorPlan - Encoun
tered Warning USING_OVERLOADED_FUNCTION 1 time(s).
2018-11-15 11:50:08,629 [main] WARN org.apache.pig.newplan.BaseOperatorPlan - Encoun
tered Warning IMPLICIT_CAST_TO_CHARARRAY 1 time(s).
2018-11-15 11:50:08,629 [main] WARN org.apache.pig.newplan.BaseOperatorPlan - Encoun
tered Warning IMPLICIT_CAST_TO_DOUBLE 1 time(s).
grunt>
  
```

Figure 10. Screenshot showing the logic on how to obtain positive polarity data



```

vishal@vishal-Inspiron-N5110: ~
grunt> negative_tweets = filter avg_tweets by tweet_rating < 0;
2018-11-15 11:51:13,324 [main] WARN org.apache.pig.newplan.BaseOperatorPlan - Encoun
tered Warning USING_OVERLOADED_FUNCTION 1 time(s).
2018-11-15 11:51:13,324 [main] WARN org.apache.pig.newplan.BaseOperatorPlan - Encoun
tered Warning IMPLICIT_CAST_TO_CHARARRAY 1 time(s).
2018-11-15 11:51:13,324 [main] WARN org.apache.pig.newplan.BaseOperatorPlan - Encoun
tered Warning IMPLICIT_CAST_TO_DOUBLE 2 time(s).
grunt>

```

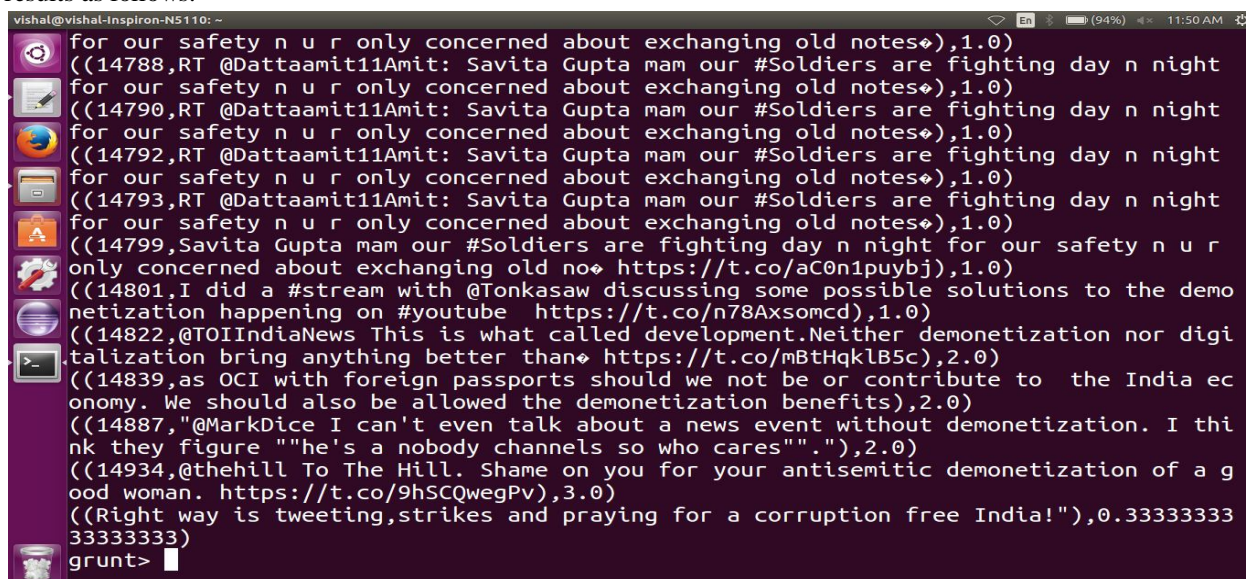
Figure 11. Screenshot showing the logic on how to obtain negative polarity data

We collect all the positive sentences and negative sentences and send them as the output, which can be collected in the subsequent phases. This data can be used in business analytics, data visualization, progress tracking, and other various applications.

IV. RESULTS AND DISCUSSIONS

The generic sentiment analysis template takes the input as processed data statements and produces positive and negatively polarized statements. After fetching the results from the generic sentiment analysis template, one can either use that data for business analysis or data visualization, mainly for predicting the success rate of a product.

Here in this experiment, we considered the Twitter data set containing people's opinion on demonetization in India. The dataset consisted of 14,941 tweets with 16 attributes. However, after performing sentiment analysis using the proposed template, one can see the results as follows:

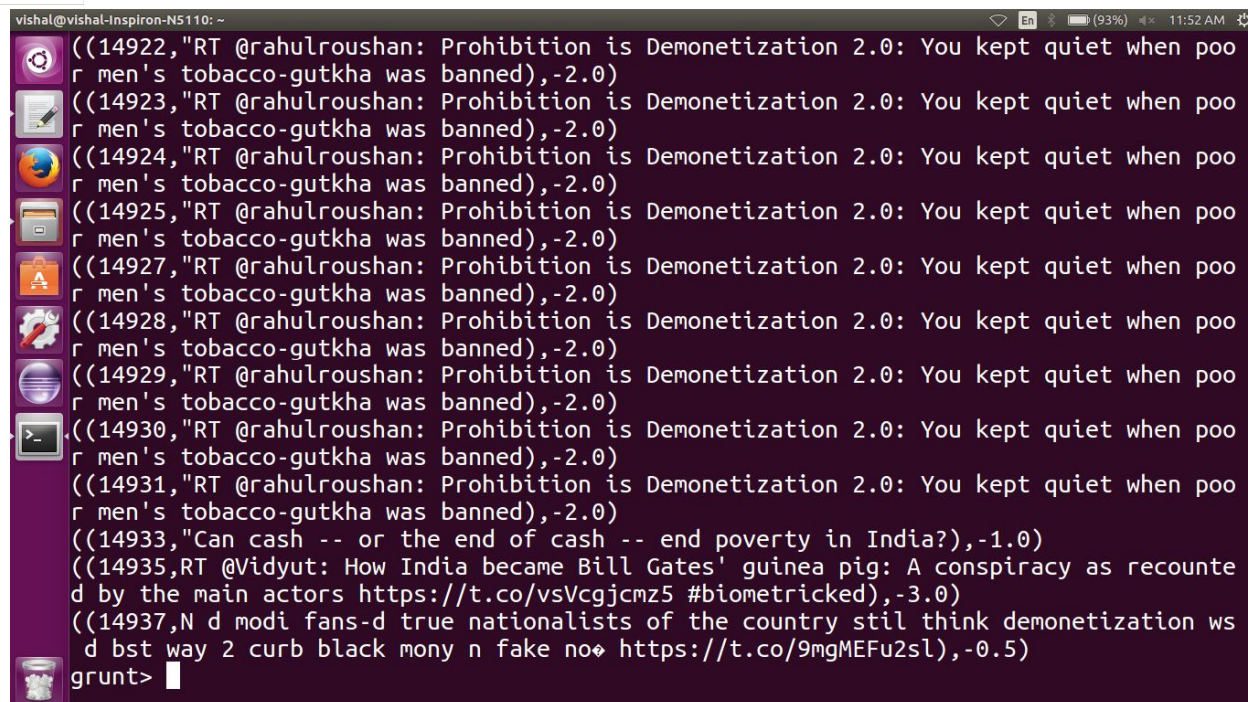


```

vishal@vishal-Inspiron-N5110: ~
for our safety n u r only concerned about exchanging old notes),1.0)
((14788,RT @Dattaamit11Amit: Savita Gupta mam our #Soldiers are fighting day n night
for our safety n u r only concerned about exchanging old notes),1.0)
((14790,RT @Dattaamit11Amit: Savita Gupta mam our #Soldiers are fighting day n night
for our safety n u r only concerned about exchanging old notes),1.0)
((14792,RT @Dattaamit11Amit: Savita Gupta mam our #Soldiers are fighting day n night
for our safety n u r only concerned about exchanging old notes),1.0)
((14793,RT @Dattaamit11Amit: Savita Gupta mam our #Soldiers are fighting day n night
for our safety n u r only concerned about exchanging old notes),1.0)
((14799,Savita Gupta mam our #Soldiers are fighting day n night for our safety n u r
only concerned about exchanging old no https://t.co/aC0n1puybj),1.0)
((14801,I did a #stream with @Tonkasaw discussing some possible solutions to the demo
netization happening on #youtube https://t.co/n78AXsomcd),1.0)
((14822,@TOIIndiaNews This is what called development.Neither demonetization nor digi
talization bring anything better than https://t.co/mBtHqkLB5c),2.0)
((14839,as OCI with foreign passports should we not be or contribute to the India ec
onomy. We should also be allowed the demonetization benefits),2.0)
((14887,"@MarkDice I can't even talk about a news event without demonetization. I thi
nk they figure ""he's a nobody channels so who cares""),2.0)
((14934,@thehill To The Hill. Shame on you for your antisemitic demonetization of a g
ood woman. https://t.co/9hSCQwegPv),3.0)
((Right way is tweeting,strikes and praying for a corruption free India!),0.33333333
33333333)
grunt>

```

Figure 12. Screenshot showing the positive tweets



```
vishal@vishal-Inspiron-N5110: ~
((14922,"RT @rahulroushan: Prohibition is Demonetization 2.0: You kept quiet when poo
r men's tobacco-gutkha was banned),-2.0)
((14923,"RT @rahulroushan: Prohibition is Demonetization 2.0: You kept quiet when poo
r men's tobacco-gutkha was banned),-2.0)
((14924,"RT @rahulroushan: Prohibition is Demonetization 2.0: You kept quiet when poo
r men's tobacco-gutkha was banned),-2.0)
((14925,"RT @rahulroushan: Prohibition is Demonetization 2.0: You kept quiet when poo
r men's tobacco-gutkha was banned),-2.0)
((14927,"RT @rahulroushan: Prohibition is Demonetization 2.0: You kept quiet when poo
r men's tobacco-gutkha was banned),-2.0)
((14928,"RT @rahulroushan: Prohibition is Demonetization 2.0: You kept quiet when poo
r men's tobacco-gutkha was banned),-2.0)
((14929,"RT @rahulroushan: Prohibition is Demonetization 2.0: You kept quiet when poo
r men's tobacco-gutkha was banned),-2.0)
((14930,"RT @rahulroushan: Prohibition is Demonetization 2.0: You kept quiet when poo
r men's tobacco-gutkha was banned),-2.0)
((14931,"RT @rahulroushan: Prohibition is Demonetization 2.0: You kept quiet when poo
r men's tobacco-gutkha was banned),-2.0)
((14933,"Can cash -- or the end of cash -- end poverty in India?,-1.0)
((14935,"RT @Vidyut: How India became Bill Gates' guinea pig: A conspiracy as recounte
d by the main actors https://t.co/vsVcgjcmz5 #biometricked),-3.0)
((14937,"N d modi fans-d true nationalists of the country stil think demonetization ws
d bst way 2 curb black mony n fake no https://t.co/9mgMEFu2sl),-0.5)
grunt>
```

Figure 13. Screenshot showing the negative tweets

V. CONCLUSION AND FUTURE SCOPE

Hadoop has been effectively proved as one of the best tools for supporting processing, storage, and streaming data. With a vast arsenal of open source software and tools, one can use them according to the demands as specified by the application. PIG facilitates processing the data to obtain favourable outcomes, which can be helpful in the future. In this paper, we created a generic template that helps the businesses understand their product reviews, how clients view their brand and considering public opinion to make business decisions more customer-centric. The project's future scope could be performing sentiment analysis on live data, which can be captured by using either a third-party application or APIs provided by the data vendors themselves. The project's scalability could also be improved if the PIG environment is run on HDFS mode rather than the local system mode.

REFERENCES

- [1] Dr. Ajit Noonia (2020), "Sentiment Analysis on Twitter Using Pig and Hive", International Journal of Innovative Technology and Exploring Engineering, ISSN: 2278- 3075, Volume 9 Issue 3, January 2020W.
- [2] Anjali Barskar and Ajay Phulre (2017), "Opinion Mining of Twitter Data using Hadoop and Apache PIG", International Journal of Computer Applications (0975-8887), Volume 158-No 9, January 2017
- [3] Kharde and Sonawane (2016), "Sentimental Analysis of twitter data- A survey of techniques", International Journal of Computer Applications, 6(10), pp. 2321-9637
- [4] Kumar and Sebastian (2015), "Sentiment Analysis on Twitter", IJCSI International Journal of Computer Science.
- [5] Dhawan and Rathee (2013), "Big Data Analytics using Hadoop Components like Pig and Hive", American International Journal of Research in Science, Technology, Engineering & Mathematics, pp. 88-93.
- [6] Shubham Goyal (2016), "Sentimental Analysis of Twitter Data using Text Mining and Hybrid Classification Approach", International Journal of Advance Research, Ideas and Innovations in Technology, 2(5) , pp. 2454-132XJ.
- [7] Rahul Kumar Chawda, Dr. Ghanshyam Thakur, "Big Data and Advanced Analytics Tools", 2016 Symposium on Colossal Data Analysis and Networking (CDAN), IEEE 2016, ISSN: 978-1-5090-0669-4/16.
- [8] Mahalakshmi R, Suseela S , "Big-SoSA: Social Sentiment Analysis and Data Visualization on Big Data", International Journal of Advanced Research in Computer and Communication Engineering, Vol. 4, Issue 4, April 2015 , pp 304-306, ISSN : 2278-1021.
- [9] Manoj Kumar Danthala, "Tweet Analysis: Twitter Data processing Using Apache Hadoop", International Journal Of Core Engineering & Management (IJCEM) Volume 1, Issue 11, February 2015, pp 94-102.



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)