# Mining the Hierarchical Dimensions of Twitter Data for OLAP

[1]Chitlaa Harshithha, [2]N.Naveen Kumar

[1]Student, [2]Associate professor
Software Engineering,
School of Information Technology, Hyderabad, India.

***Abstract:*** Social media platforms like, Twitter, Facebook disclose lot of information regarding the tastes of the people. In order to help in promotion of products and investigation of sentiments, there are many researches that focus on the content analysis. Also, OLAP (online analytical processing) has been demonstrated to be exceptionally powerful to examine multidimensional structured data. Text OLAP, the sole target behind applying OLAP to the text messages is to mine and develop the dimension which is hierarchical depending upon the data content which is unstructured. Usually text OLAP handles plain texts, but in the contrary the social platforms' information/content incorporates an abundance of information regarding social relationships which can be utilized to dig out an effective dimensional hierarchy. So here we come up with a topic model which can be called THLDA, Twitter Hierarchical Latent Dirichlet allocation. The target of THLDA model is to naturally mine the hierarchical dimension of tweets' topics, which can be additionally utilized for text OLAP on the tweets. Moreover, THLDA utilizes a procedure word2vec to break down the semantic connections of words in tweets to get an increasingly viable measurement. We lead broad examinations on tremendous amounts of Twitter information and assess the viability of THLDA. The results of the experiment show that the method beats other current subject models in mining and building the various leveled measurement of tweeters' points.

*Index Terms* **- Topic Modeling, Text Modeling, Twitter data,Latent Dirchlet Allocation.**

## I. INTRODUCTION

During the previous scarcely any years, Twitter has gotten progressively famous as a developing social stage for informing and correspondence among people. The enormous amounts of Twitter information aggregated so far is making it conceivable in order to find conveyance, float of mass taste and suppositions, which incredibly aids target marketing, product recommendation, etc. Then again, Online Analytical processing, empowers investigation to intelligently see information from various perspectives in layered granularities, which has just been demonstrated particularly helpful for the business insight. Lamentably, OLAP strategies are effective in managing shape information which are organized and formalized, however face troubles in handling literary substance, for example, Twitter information. To effectively apply the OLAP methods to Twitter, it is essential to mine the concealed delegate measurements from its broad substance.

Latent Dirichlet Allocation(LDA) model, which is typical unsupervised model, is proficient at factually breaking down literary information for the hidden themes. We have put forward MS-LDA, which is an LDA-based model to recognize the concealed layered interests from the Twitter information. MS-LDA, as the expansion of LDA, incorporated tweets and social connections between people who use twitter. All things considered, the crude LDA model can just mine single layer/monolayer themes, as opposed to the various leveled ones which OLAP requires. Then again being unsupervised hierarchical topic model,hLDA can acquire the sibling-sibling relationships connections among themes and also can arrange the points into a progressive tree consequently. Truth be told, Twitter information contain plentiful social conduct data about tweeters, for example, following, mentioning, retweeting. Moreover, some semantic connections also exist among the words in the tweets, which can influence the viability of the demonstrating procedure. At the end of the day, to viably find the concealed layers of subjects from the Twitter information in order to develop the progressive measurement for OLAP, we have to put forward another topic model, thatwhich can use the qualities of Twitter in the demonstrating procedure.

## EXISTING SYSTEM

During the previous scarcely any years, Twitter has gotten progressively famous as a developing social stage for informing and correspondence among people. The enormous amounts of Twitter information aggregated so far is making it conceivable to find the conveyance and float of mass tastes and suppositions, which incredibly aids target marketing, product recommendation etc. Then again, online Analytical processing or OLAP, empowers investigation to intelligently see information from all perspectives in layered granularities, which has just been demonstrated particularly helpful for business insight. Lamentably, OLAP strategies are effective in managing shape information which are organized and formalized, however face troubles in handling literary substance, for example, Twitter information. To effectively apply OLAP methods to Twitter, it is basic to mine the concealed delegate measurements from its broad

substance.

**Disadvantages**

- Facing difficulties in processing data/textual content.
- It is basic to mine the concealed agent measurements from its broad substance

## PROPOSED SYSTEM

We thought of a model, which is based on LDA, MS-LDA ,called Multilayered Semantic LDA, to discover the concealed layered interests from the information got from Twitter, a social networking platform. As the expansion of LDA, MS-LDA incorporated the social connections among people who tweet and tweets. By and by, the crude LDA model can just mine single layer subjects, instead of the hierarchical ones which OLAP actually wants. Then again, as an unaided various leveled topic model, hLDA can get the kin connections among points and can sort out the subjects into a hierarchical tree consequently. Truth be told, Twitter information contain bountiful social conduct data about tweeters, for example, retweeting, mentioning and following. Some semantic connections among the words in tweets also exist, which may influence the viability of the displaying procedure. As such, to adequately find the concealed layers of subjects from Twitter information for building the progressive measurement for OLAP, we need to put forward another topic model which use the attributes of Twitter in its demonstrating procedure.

**Advantages**

- Detection of the layered interests which are hidden, from the data obtained i.e., Twitter data.
- HLDA can attain the sibling-sibling relationships between topics and can put forth the topics into a hierarchy based tree consequently.

### I. SYSTEM REQUIREMENTS

#### HARDWARE REQUIREMENTS

- Hard Disk : 20GB
- RAM     : 256 MB
- Processor : Core – i5

#### SOFTWARE REQUIREMENTS

- Operating System : Windows 8
- Coding Language : Java
- IDE        : Eclipse

### II. RELATED WORK

**Mining Hidden Interests from Twitter Based on Word Similarity and Social Relationship for OLAP [1]**
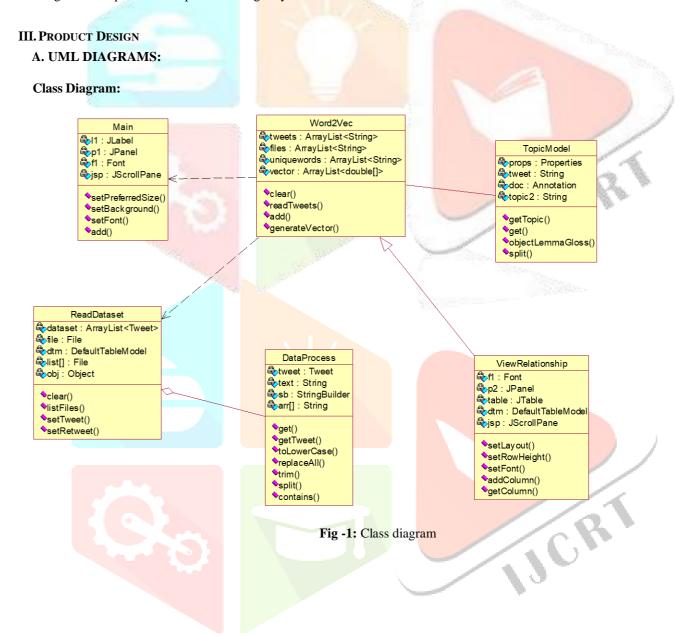
Online Analytical Processing, is a way to deal with answering multidimensional analytical questions in a natural way. Be that as it may, the conventional OLAP approaches can just do arrangement with organized information, however not with unstructured textual information like tweets. To work on this issue, we put forward LDA (Latent Dirichlet Allocation)- based model, which can be identified as, MS-LDA , Multilayered Semantic LDA ,which can distinguish the concealed layered premiums from information obtained from Twitter dependent on Latent Dirichlet Allocation. The layered component of interests can also be used additionally to apply Online Analytical Processing ,OLAP procedures to information obtained from Twitter. Besides, Multilayered Semantic LDA utilizes the semantic comparability between tweets' expressions dependent on word2vec, and furthermore the social relationship among people who tweet ,so that adequacy can be improved. The broad examinations exhibit that Multilayered Semantic LDA can viably remove the measurement pecking order of tweeters' inclinations towards Online Analytical Processing.

Now, we put forward an improvised subject model, for example MS-LDA, Multilayered Semantic LDA ,which can be utilized to extricate the measurement progressive systems of tweeters' inclinations, regularly covered up in the huge measure of unstructured Twitter information. We directed broad investigations on Twitter information to assess the viability of MS-LDA. The outcomes show that Multilayered Semantic LDA in fact has the great acknowledgment impact. The word2vec model utilized in this project is prepared utilizing news given by Google. In any case, the introduction of news when all is said in done is to some degree thorough, while tweets are increasingly casual.

**A Method for Online Analytical Processing of Text Data [2]**

There are progressively noticeable requests for organized/unstructured data reconciliation and progressed investigation. Be that as it may, ordinary database innovation has not had the option to introduce a vigorous and down to earth execution of a really coordinated engineering for such purposes. Subsequent to taking a shot at a few mechanical applications (specifically, in the life sciences and medicinal services region), we have recognized essential problems and specialized ways to deal with the issues.Here, we come up with information portrayals and logarithmic activities for coordinating semantic data (Ex: ontologies) into OLAP frameworks, which will permit to break down a tremendous arrangement of literary archives with their fundamental semantic data. The presentation of the model execution has been assessed utilizing genuine world datasets, and the high adaptability and adaptability of our methodology have been affirmed regarding the calculation time.

In this paper, we put forward an information portrayal and its polynomial math activities to coordinate ontologies with OLAP frameworks to break down an immense arrangement of literary reports. By using our strategy, two kinds of data (structured and unstructured data) can commonly improve data disclosure and investigation extent. Utilizing preorder and post order in a hierarchy, the proposed strategy was executed with a tireless store. The proficiency of our methodology has been affirmed as for the calculation time. Our strategy is so effective and strong that it empowers an expert to intelligently.

### III. PRODUCT DESIGN

#### A. UML DIAGRAMS:

**Class Diagram:**



**Fig -1:** Class diagram

**Use Case Diagram:**



**Fig -2:** Use Case diagram
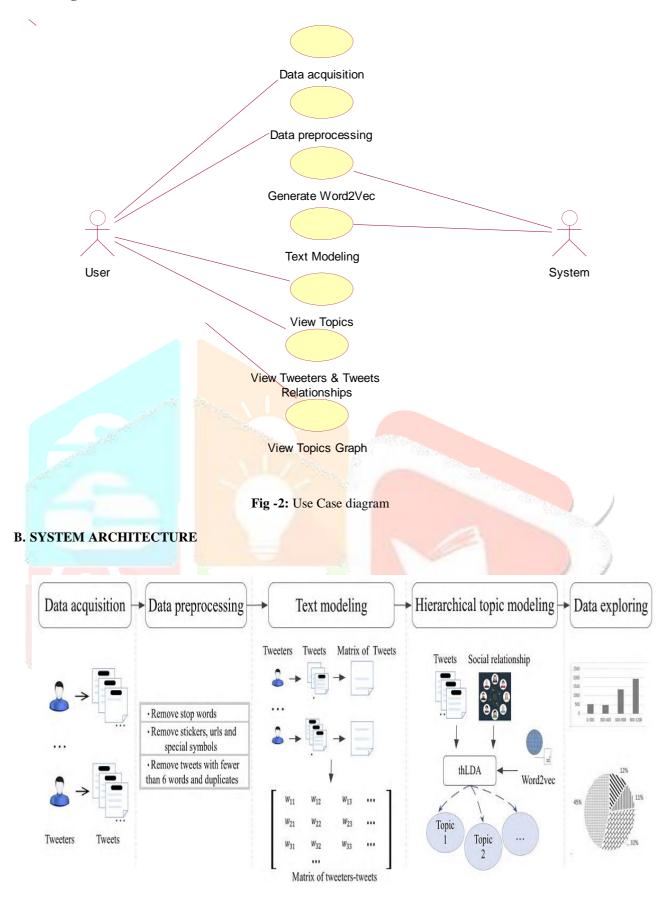
## B. SYSTEM ARCHITECTURE



**Fig -9:** System Architecture

## IV. IMPLEMENTATION

The general procedure of investigating Twitter information dependent on the OLAP method can be portrayed as follows:

• **Data acquisition:** Using the REST APIs which are provided by the twitter, tweets, social relationship and tweeters' profiles, can be obtained.

• **Data preprocessing:** Short words i.e., the most common words like the, is, at etc., should be removed along with the web links. Then parts of speech analysis should be done to remain with verbs and nouns in the unstructured tweets.

• **Text modeling:** The Relationship between tweeters' and tweets can be recognized based on the text modeling.

 • **Hierarchical topic modeling:** Interests(or topics) from the data obtained from Twitter can be extracted, and the hierarchical topic dimension can be constructed based on the probability distribution of various topics and subtopics.

• **Data exploring:** Using OLAP, analysis of tweeters from multiple dimensions can be done.

### Methodologies

In this paper the concept to extract topics from tweets and can extract relationships between tweeters and tweets. Relationship can be extracted by analyzing two person's tweets and look for semantic similarity between their tweets and if both are talking on same matter then similarity will be higher and both tweets users will have similarity and relationship can be bond between them.

In this paper to extract relationship and topic author is asking to generate WORD2VEC (word to vector) and also called as BAG of WORDS (BOG). After forming vector we can easily extract relationship between two vectors and can also extract topics.

WORD2VEC conversion means converting tweets into vector array and each array will be consider as one tweet and in each tweet each occurrence or count of each word will put inside that array. If two tweets have common words then that array column will have value > 0 and similarity will be found.

### IV. TESTING AND RESULTS

The procedure or strategy for discovering errors or defects in an application or software program with the goal that the application functions as indicated by the end client's prerequisite is called testing. The test case is defined as a set of conditions or factors under which a tester will decide if a system or application under test satisfies requirements and works properly.

Following are the testing strategies followed during the test period:

**TEST CASES:**

**Table:** Test Cases and Results

| Test Case Id | Test Case Name | Test Case Desc. | Test Steps | | | Test Case Status | Test Priority |
|---|---|---|---|---|---|---|---|
| | | | **Step** | **Expected** | **Actual** | | |
| 01 | Data Acquisition | Verify dataset is available or not | If it is available | We can load the dataset | Tweets dataset loaded | High | High |
| 02 | Data preprocessing | Verify dataset loaded or not | If it's loaded | We can clean tweets | All special symbols and URL's remove from tweet | low | High |

| 03 | Generate Word2Vec | Verify the tweets are cleaned or not | If it is cleaned | We can generate vector | We get the average occurrence of words in tweets | Medium | High |
| 04 | Text modeling | Verify vector generated or not | If it's generated | We can apply the OLAP & LDA | We will get processing of each tweet to extract topic | High | High |
| 05 | View topics | Verify topics extracted or not | If it's extracted | We can view the topics | We can view the all topics from tweets | High | High |
| 06 | View tweeters & tweets relationships | Verify topics are viewed or not | If it's viewed | We can view the relationships | We get the cosine similarity between tweeters & tweets | High | High |
| 07 | View topics graph | Verify all topics viewed or not | If it's viewed | We can get the graph | We get the topics graph | High | High |

**RESULTS:**

The dataset which I have used is downloaded from Kaggle website. Then I have made each tweet into a separate file.



On running the batch file, we will get this screen. Here we can see the buttons that are created and in the empty field the headers are properly arranged. On clicking the buttons the functional changes can be seen in th e empty area of the screen.

On clicking Data acquisition button on the screen we can selects the folder in which the tweet individual documents are stored as input to the application. This is uploading the data obtained from Twitter. After choosing the required folder we can click on open. Then the changes can be seen on the screen.

OLAP

HIERARCHICAL TOPIC MODELING OF TWITTER DATA FOR ONLINE
ANALYTICAL PROCESSING

| Tweet File | Tweet Text | Retweet Count |
|---|---|---|
| 1.txt | RT @rssurjewala: Critical question: Was PayTM informed about #Demonetization edict by PM? It's clearly fishy and re... | 75 |
| 100.txt | RT @Gadgets360: After the #demonetizatio | 96 |
| 11.txt | RT @sumitbhati2002: Many opposition leaders are with @narendramodi on the #Demonetizatio | 70 |
| 13.txt | National reform now destroyed even the essence of sagan. Such instances urge giving #demonetization a second though... | 43 |
| 14.txt | Many opposition leaders are with @narendramodi on the #Demonetizatio | 15 |
| 16.txt | RT @Joydas: Question in Narendra Modi App where PM is taking feedback if people support his #DeMonetization strat... | 52 |
| 17.txt | @Jaggesh2 Bharat band on 28??<ed><U+00A0><U+00BD><ed><U+00B8><U+0082>Those who are protesting #demonet... | 89 |
| 18.txt | RT @Atheist_Krishna: The effect of #Demonetization ! | 32 |
| 2.txt | RT @Hemant_80: Did you vote on #Demonetization on Modi survey app? | 46 |
| 20.txt | RT @sona2905: When I explained #Demonetization to myself and tried to put it down in my words which are not laced ... | 95 |
| 21.txt | RT @Dipankar_cpiml: The Modi app on #DeMonetization proves once again that the govt is totally indifferent to the m... | 41 |
| 22.txt | RT @roshankar: Former FinSe | 90 |
| 24.txt | RT @Atheist_Krishna: BEFORE and AFTER Gandhi ji heard they are standing there against #Demonetizatio | |
| 26.txt | RT @pGurus1: #Demonetization The co-operative banking sector in Kerala is as good as a tax haven. Is Keral | |
| 27.txt | RT @roshankar: Former FinSe | |
| 29.txt | RT @Hemant_80: Did you vote on #Demonetization on Modi survey app? | |
| 3.txt | RT @roshankar: Former FinSe | 89 |
| 30.txt | RT @roshankar: Former FinSe | 64 |
| 32.txt | RT @Atheist_Krishna: BEFORE and AFTER Gandhi ji heard they are standing there against #Demonetizatio | 50 |
| 34.txt | RT @MahikaInfra: @narendramod | 82 |
| 36.txt | RT @Hemant_80: Did you vote on #Demonetization on Modi survey app? | 72 |
| 37.txt | RT @roshankar: Former FinSe | 13 |
| 39.txt | RT @kapil_kausik: #Doltiwal I mean #JaiChandKejriwal is ""hurt"" by #Demonetization as the same has rendered US... | 43 |
| 40.txt | RT @roshankar: Former FinSe | 58 |
| 42.txt | RT @kapil_kausik: #Doltiwal I mean #JaiChandKejriwal is ""hurt"" by #Demonetization as the same has rendered US... | 44 |
| 43.txt | RT @AAPVind: #Demonetization Is Disaster! @naam_pk | 80 |
| 44.txt | RT @Hemant_80: Did you vote on #Demonetization on Modi survey app? | 96 |

Message

Dataset Loaded

OK

| Data Acquisition | Data Preprocessing | Generate Word2Vec | Text modeling | View Topics | View Tweeters & Tweets Relationships |
|---|---|---|---|---|---|
| | | View Topics Graph | Exit | | |

After the data gets loaded the alert is displayed.

OLAP

HIERARCHICAL TOPIC MODELING OF TWITTER DATA FOR ONLINE
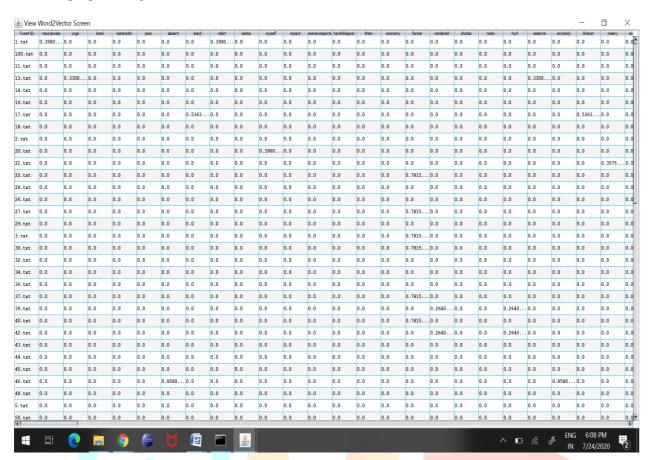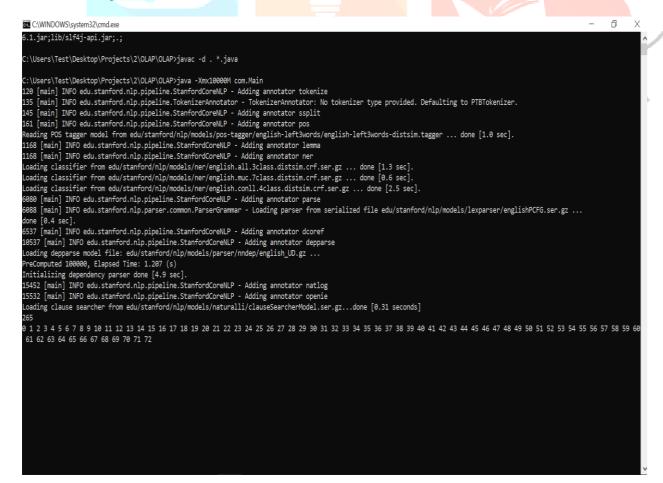ANALYTICAL PROCESSING

| Tweet File | Tweet Text | Retweet Count |
|---|---|---|
| 1.txt | rssurjewala critical question paytm informed demonetization edict clearly fishy requires disclosure | 75 |
| 100.txt | gadgets after demonetizatio | 96 |
| 11.txt | sumitbhati many opposition leaders narendramodi demonetizatio | 70 |
| 13.txt | national reform destroyed even essence sagan such instances urge giving demonetization second eyysireiuq | 43 |
| 14.txt | many opposition leaders narendramodi demonetizatio | 15 |
| 16.txt | joydas question narendra modi taking feedback people support demonetization strategy pygk | 52 |
| 17.txt | jaggesh bharat band protesting demonetization different party leaders | 89 |
| 18.txt | atheist krishna effect demonetization | 32 |
| 2.txt | hemant vote demonetization modi survey | 46 |
| 20.txt | sona when explained demonetization myself tried words which laced heavy technical | 95 |
| 21.txt | dipankar cpiml modi demonetization proves again govt totally indifferent mounting misery hards | 41 |
| 22.txt | roshankar former finse | 90 |
| 24.txt | atheist krishna before after gandhi heard standing demonetizatio | 30 |
| 26.txt | pgurus demonetization operative banking sector kerala good haven kerala black money | 17 |
| 27.txt | roshankar former finse | 21 |
| 29.txt | hemant vote demonetization modi survey | 19 |
| 3.txt | roshankar former finse | 89 |
| 30.txt | roshankar former finse | 64 |
| 32.txt | atheist krishna before after gandhi heard standing demonetizatio | 50 |
| 34.txt | mahikainfra narendramod | 82 |
| 36.txt | hemant vote demonetization modi survey | 72 |
| 37.txt | roshankar former finse | 13 |
| 39.txt | kapil kausik doltiwal jaichandkejriwal hurt demonetization same rendered useless acquired funds | 43 |
| 40.txt | roshankar former finse | 58 |
| 42.txt | kapil kausik doltiwal jaichandkejriwal hurt demonetization same rendered useless acquired funds | 44 |
| 43.txt | aapvind demonetization disaster naam | 80 |
| 44.txt | hemant vote demonetization modi survey | 96 |

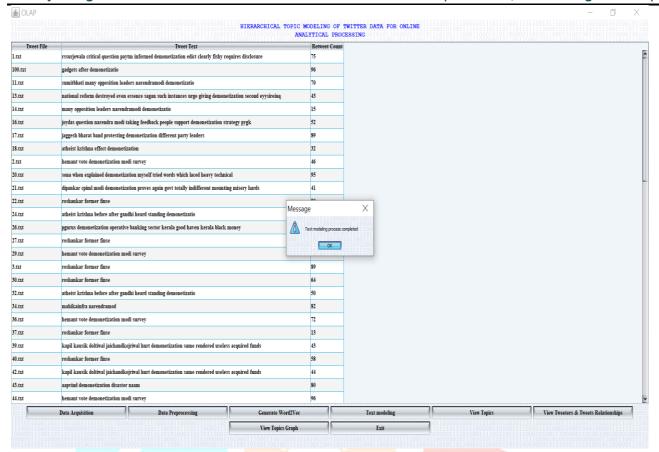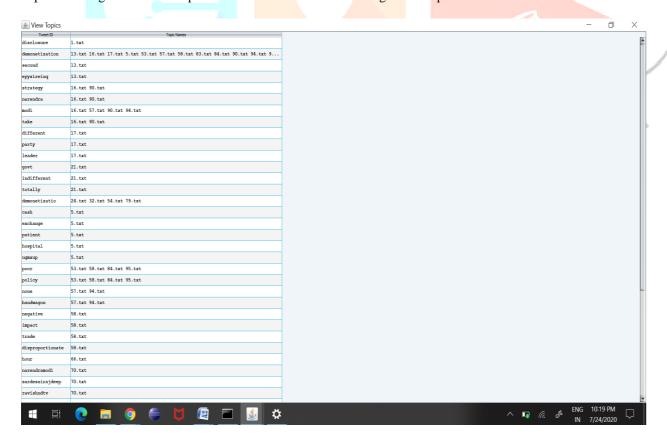| Data Acquisition | Data Preprocessing | Generate Word2Vec | Text modeling | View Topics | View Tweeters & Tweets Relationships |
|---|---|---|---|---|---|
| | | View Topics Graph | Exit | | |

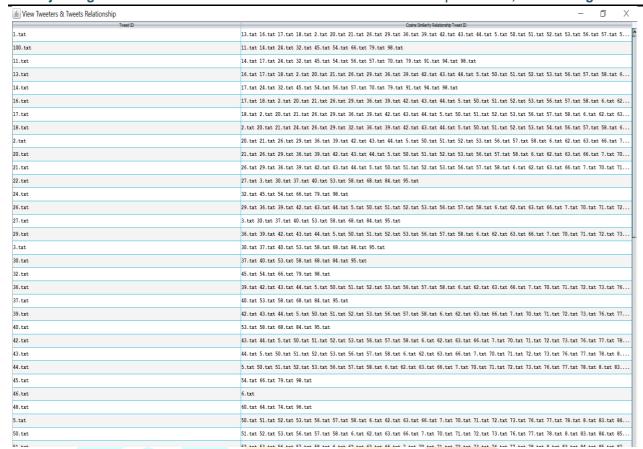Data after pre processing.



Word2Vector is generated.



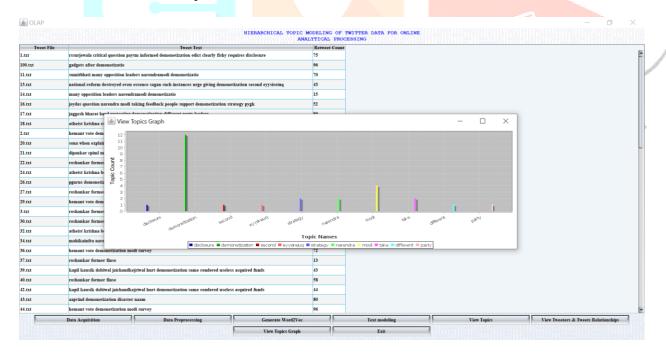Here we can see that each tweet is getting processed on clicking Text Modeling button.

Topic Modeling is done and Topics obtained can be seen on clicking View Topics.



List of topics obtained from given dataset.

Tweeters and Tweets Relationship can also be obtained.



In above graph x-axis showing topic name and y-axis showing number of times that topic appear in all tweets.

## VI CONCLUSION

Here, we have proposed Twitter Hierarchical Latent Dirchlet Allocation method i.e.,thLDA which is a novel hierarchical topic model. This method is applied to the large quantity of unstructured Twitter data in order to mine the dimension hierarchy of tweets' topics. The effectiveness of this method is found out by performing extensive experiments on real data from Twitter. Results of the experiment confirm that this method is more effective than any other models.

## VII FURUTE SCOPE

Indirect social relationships between the tweeters can be analyzed in order to enhance our current model, in the future.

### REFERENCES

[1] D. Yu et al., ''Mining hidden interests from Twitter based on word similarity and social relationship for OLAP,'' Int. J. Softw. Eng. Knowl. Eng., vol. 27, nos. 9–10, pp. 1567–1578, 2017.

[2] A. Inokuchi and K. Takeda, ''A method for online analytical processing of text data,'' in Proc. 16th ACM Conf. Conf. Inf.    Knowl. Manage., 2007, pp. 455–464.