

## Topical clustering is used to analyze tweet data

<sup>1</sup>Ravi kumar kuchipudi, <sup>2</sup>Dr V.V. Jaya Rama Krishnaiah,

<sup>1</sup>research scholar, ANU, Computerscience and Engineering, ravikumary4@gmail.com

<sup>2</sup>Associate professor, ASN Engineering college, Tenali

---

### Abstract

Practical management guidance and the spread of "social networking" have made "networking" more important than ever. As a result of social media, massive amounts of data are generated. It is estimated that over 400 million people use Facebook each month, exchanging over 5 billion pieces of information. Online social networks like Twitter, Facebook, LinkedIn, and Instagram connect people from all over the world. An increasing number of applications use social network analysis, which allows us to gather important information about the people in the network, share information, or make connections. Our approach to analyzing Twitter and Facebook profiles by geography is presented in this paper. The user should be able to choose the locations of these other individuals. According to the analysis, we'll compare Facebook and Twitter profiles based on where they're located, and then we'll extract tweets and comments made by people in that area. As a result, the project's focus is on big data analytics..

**Keywords** - Big data Analytics, Twitter and Facebook mining, social networking Analysis

### 1. Introduction

The term "social network" refers to web-based services that allow users to create a public or semi-public profile within a domain and communicate with other users inside the network. A Facebook account can be used for a variety of things, including connecting with friends and sharing photos. Social media users share their thoughts and perceptions on a wide range of topics. It can cover a wide range of topics, such as politics, religion, technology, and the latest movie releases, as well as other topics that are frequently brought up in daily conversation. For years, the number of people using social networking sites such as Facebook has risen rapidly. More recently, users of social networking sites such as Twitter have been referred to as "friends," although this was not the case when Twitter first launched. In contrast to following someone on Twitter, adding someone as a friend on Facebook creates a much stronger connection. As with Facebook, you can establish meaningful connections on Twitter. According to how Facebook users interact with their family and friends, Twitter following may be interested.

Facebook buddies, but on Twitter, people discuss about those who share their interests. users can choose between a completely public profile and one that is only viewable to acknowledged friends on Facebook, for example. There are two secure options on Twitter. Only the people who the user follows will be able to access their private messages. There are no different privacy messages for individual communications. It's possible to combine Facebook and Twitter. Twitter apps allow automatic posting of tweets to Facebook. We can acquire location-based data from Twitter and Facebook using this technique. On the basis of information gleaned from the tweets and posts made by local users, as well as information drawn from the interested domin. The technique of picking words and phrases from a text document is known as keyword extraction. With the LDA method, we can find out how many users are interested in a particular location, and then use supervised and unsupervised algorithms to build a network based on that information. Following that, we have a notification to your mail. Through that we have link through then graph is generated.

## **2. Related work**

### **2.1 Data Retrieval:**

The API4j is used to get data from Twitter. To begin, get a list of all the tweets from the users in that geographic area by utilizing the API that was provided by the user focused region. We used API4J (JAVA) to get data from Facebook. Retrieve the information depending on the geographical location of the users and the comments or posts that they have made on the relevant domain.

### **2.2 Topic identification:**

When it comes to finding hot subjects, the two most common techniques are LDA and PLSA. When used to identify topics, LDA is a probabilistic generative model that may be used to a variety of problems. In the same way, PLSA is a statistical technique for topic modeling that can be used. Although temporal information is lost in these methods, it is crucial for detecting common subjects and is a feature of social media data. The LDA model is used by Facebook and Twitter to discover trending themes in posts and tweets. Despite this, their study exclusively considers the private interests of consumers rather than global issues.

### **2.3 Keyword Extraction:**

Many unsupervised and supervised approaches have been presented for keyword or informative term extraction. Unsupervised keyword extraction approaches rely exclusively on implicit information discovered in a single text or in a text corpus as a whole. Methods that use previously classified training datasets are known as supervised methods. supervised and unsupervised keyword extraction is used.

different but complementary methods. Key words can be effectively extracted using machine learning techniques such as those found in KEA and GenEx, two established supervised frameworks. Other cutting-edge ways to keyword extraction, such as using neural networks, have been proposed in the last few years. Identify relevant phrases in the news articles using the sources' keywords.

### **2.4 Co-Occurrence Similarity:**

There is a correlation between the frequency with which certain word pairings appear in a given document. Supply keyword identification statistics that can be used in other documents you create. As long as the distribution of co-occurrence probabilities between terms is tilted toward some particular subset of common terms, the researchers concluded, the term in question is likely to be a keyword.

Co-occurrence similarity is a metric used to assess the degree to which search results return snippets that show the relationship between phrases. Double checking for co-occurrence is the term they use for this technique (CODC). To determine how closely two words or entities are related, I proposed an approach that makes use of page counts and text samples from Web searches.

### 3. Implementation

#### 3.1 Extraction of Twitter Profiles:

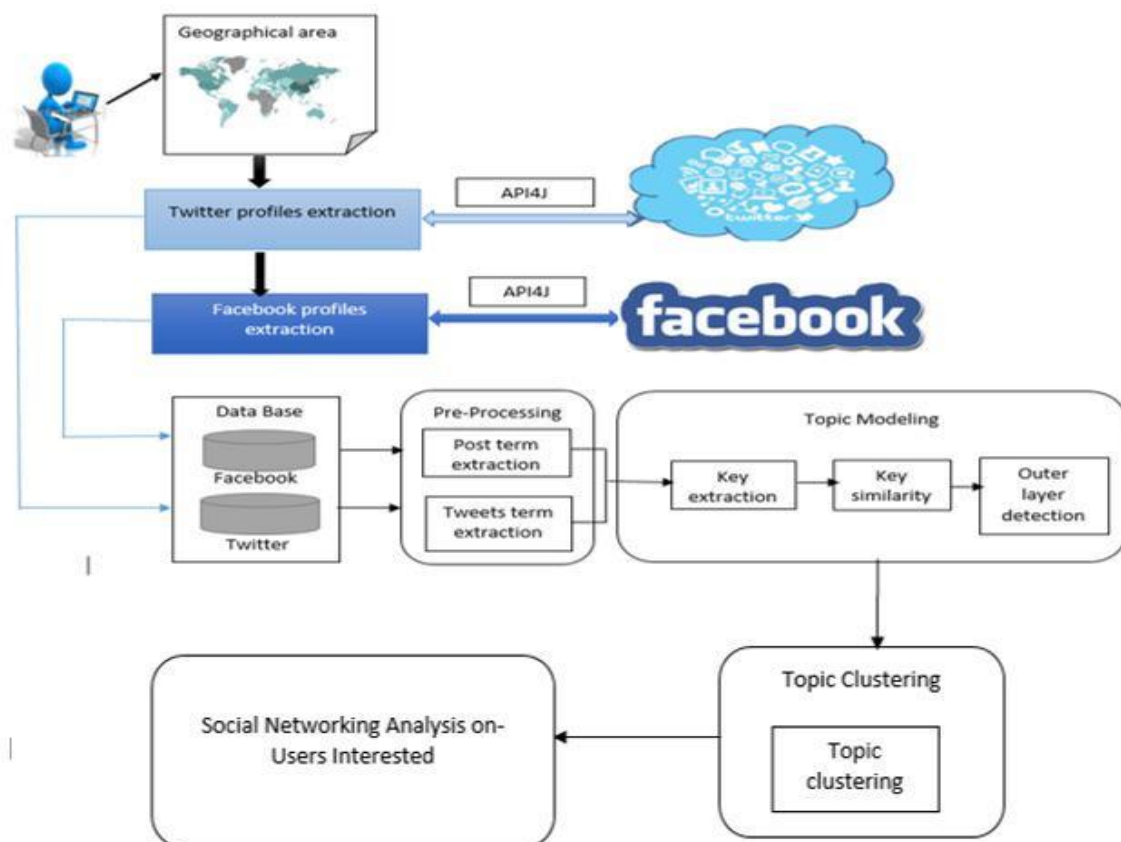
It's at this point that the user decides which user community they wish to focus on. In order to complete this activity, enter a search term depending on your current location. Additionally, the user has the option of choosing from which geographical regions (locations) to get tweets. A csv file contains information about the various geographical areas, such as their names and locations.

Direct access to some Twitter data is possible using the public Twitter API. Due to the fact that Twitter profiles are open to the public, we don't have to ask individuals for this information. To communicate with Twitter, we'll make use of the Twitter4J library. This library is simple to use and can be quickly integrated into our application. It is possible to search for tweets based on the search term and the geographical area supplied by the user with this library's help. Users' profiles (including the search term) are then gathered from the tweets..

#### 3.2 Extraction of Facebook Profile:

Third-party components of Facebook Java APIs are used to extract the data from the social networking site. The training set consists of about 2000 posts from various users. Data from live news streams is used to classify the results using several classifiers..

Figure 3.1: General structure of our approach



### 3.3 Preprocessing:

The following methodology extracts and filters key phrases from news and social data for a specific time period..

- Stop words
- Data filters

### 3.4 Topic Modeling:

It's necessary to analyze and model the data pulled from Twitter and Facebook based on the material that's already there. Using Topic Modeling, which is a text-mining technology, it is possible to find latent semantic structures in a document. In text corpora, it is the finding of "themes" through the grouping of frequently occurring co-occurring terms. In order to uncover the abstract "themes" that appear in a collection of documents, we use a statistical model of this type. When it comes to topic modeling, we employ a number of different strategies. for instance

- a) Latent Dirichlet allocation
- b) Probabilistic latent semantic analysis.

Despite the exclusion of many possibly irrelevant terms, the graph still contains an excessive number of terms (vertices) and co-occurrences (edges). We only want to include term co-occurrences with high QS (Quantitative Similarity) values, so that we may filter out the less important ones. Outliers (irregular co-occurrence values) must be distinguished from regular ones in the graph before meaningful edges can be identified.

#### 3.4.1 Keyword Extraction:

This procedure begins with gathering tweets and posts, which are then analyzed for keywords to extract the most important information. Consider the keyword "movies" as an example. Then, in our studies, we collected all tweets mentioning this key phrase. As a result, the corpus will consist of tweets. The Twitter and Facebook API4J was used to get tweets from the internet.

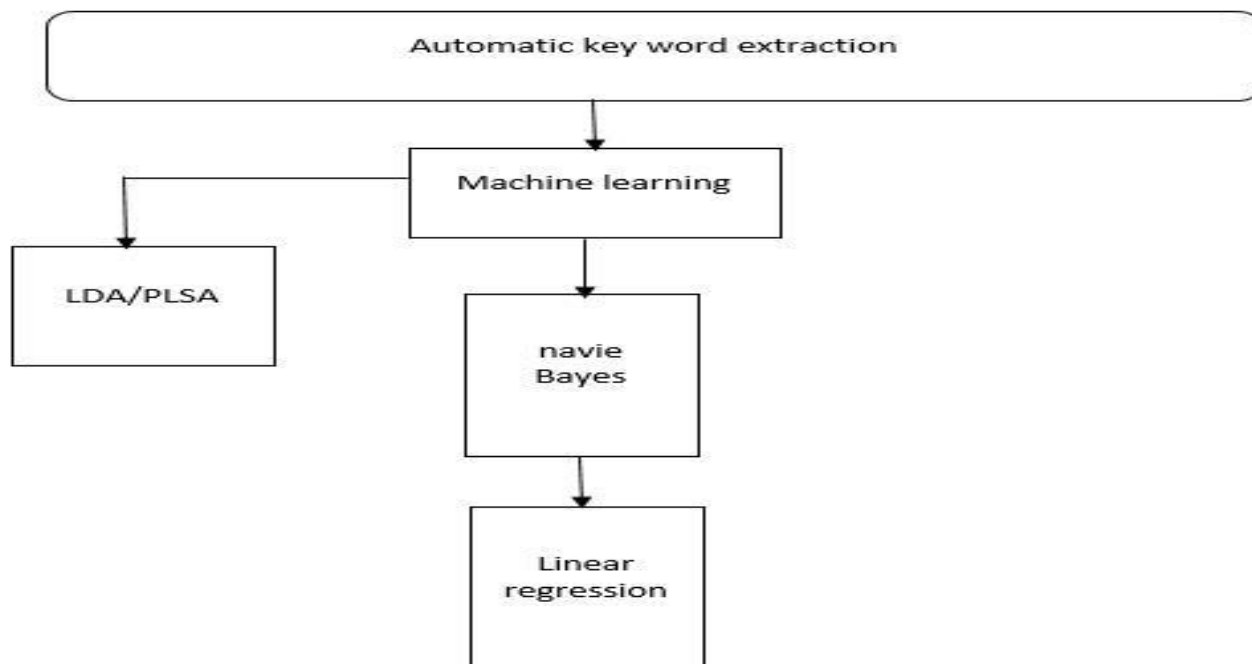


Figure 3.4.1: keyword Extraction

3.4.2 Keyword Similarity:

3.5 Topic Clustering:

Our analysis of Facebook and Twitter profiles uncovered a high degree of term similarity, as well as an understanding of the keyword usage scope among different people. Taken into account as words that can be found in the Oxford English Dictionary. We made use of information from a person's Facebook and Twitter profile's Interests section to do this. As a result, we take into account the positive feedback from users who describe their interests and passions.

3.4.3 Outliers Detection:

Using the user's preferences, outliers are quickly discovered. The parameter sigma can be set by the user to a threshold value for outlier detection. This task's output is a list of outliers and a graph showing the density of the investigated samples.

Traditional clustering algorithms may not be suitable when dynamic content is used. Clustering algorithms should be able to assess and update conclusions incrementally as new data is introduced during the streaming process, according to some researchers. It uses a small amount of primary memory and only processes each piece of data once. To handle streaming data and extract significant keywords and build the model, we propose utilizing a Hybrid technique that combines supervised and unsupervised approaches..

3.6 Social Networking Analysis On-Users Interested:

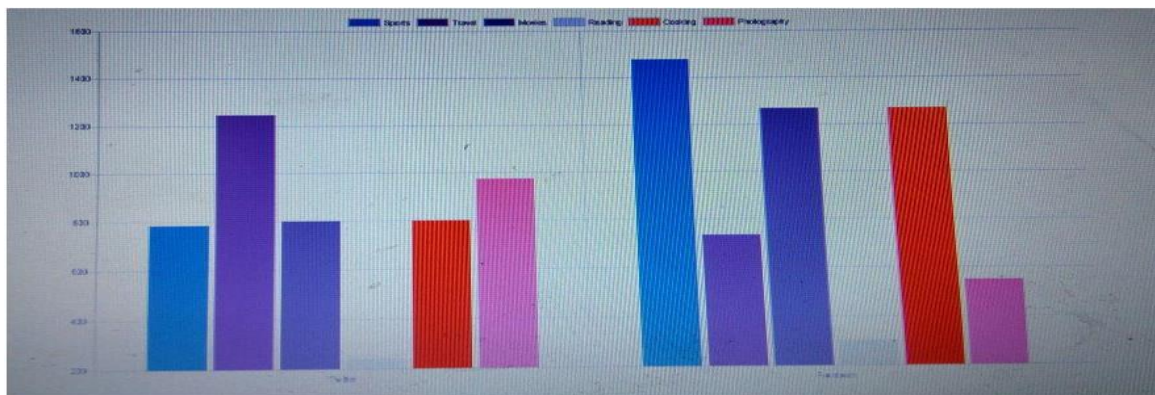
Data clustering based on user interest areas is complete. Through the graph's created URL, an email is sent to a certain mail id. The graph is created when the URL has been generated..

Social Networking Sites	Sports	Travel	Movies	Reading	Cooking	photography
Facebook	1478	734	1267	300	1267	548
Twitter	783	1245	800	240	800	974

Figure 4.45: location wise user interested areas

4. Results

Figure 4: Graph Based On User Interested Area



## 5. Conclusions

A method for analyzing Twitter and Facebook profiles based on their geographical location is presented in this paper. The user should be able to choose the locations of these other individuals. It is proposed to do a comparison between the Facebook and Twitter profiles.

## References

1. P. P. Angelov and X. Zhou, "Evolving fuzzy classifier for novelty detection and landmark recognition by mobile robots," in *Mobile Robots*, 2007, pp. 89–118.
2. T. Jo, M. Lee, and T. M. Gatton, "Keyword extraction from documents using a neural network model," in *Proc. Int. Conf. Hybrid Inf. Technol. (ICHIT)*, vol. 2.2006, pp. 194–197. SociRank: Identifying and Ranking Prevalent News Topics Using Social Media Factors *Davis, Gerardo Figuera, and Yi-Shin Chen* Derek 2016 IEEE. Personal use is permitted
3. Priyanka, Anand.R1 Dr.Rajashekhar M.Patil "comparison of betweenness and closeness centralities using incremental algorithms in dynamically growing networks", *IJACET*, ISSN(PRINT):2394-3408,(ONLINE):2394-3416,VOLUME-3,ISSUE-2,2016.
4. Anand.R, Pushpalatha .M, Dr Rajshekhar M Patil " A parallel algorithm for reading the different variables in social networks using data mining techniques" *Journal of Advanced Computing and Communication Technologies* (ISSN: 2347 - 2804) Volume No.4 Issue No.2, April 2016 .
5. Bharathi M , S. N. Chandra shekara , Anil G.N , Anand R , Muneshwara M.S " K-Anonymity for Real-time Social Network Applications with Network Based anonymization and processing framework" 2012 IACSIT Coimbatore Conferences IPCSIT vol. 28 (2012) © (2012) IACSIT Press, Singapore.