

Current Trends in Text Mining for Social Media

Tajinder Singh*, Madhu Kumari, Triveni Lal Pal and Ahsan Chauhan

*Department of Computer Science & Engineering,
National Institute of Technology, Hamirpur, H.P, 177005, India
nith2k14@gmail.com, madhu.jaglan@gmail.com,
trivenipal@gmail.com, ahsan.ch23@gmail.com*

Abstract

Online social media has created new paradigms of information sharing which not only provides appropriate platform for the contributors but also for active information seekers. Numerous forms of social media have gained the widespread attention internet users' almost on explosion level. Availability of data on such behemoth scale mandates regular and critical analysis of this information for various perspectives' plausibility. As text mining plays a significant and crucial role in discovery of these insights therefore its challenges and contribution in social media analysis must be explored extensively. This paper focuses on description of assorted social media text mining methods with their viabilities and summarizes various commercial and open source tools for such data analysis

Keywords: *analysis, clustering data mining, event detection, Twitter*

1. Introduction

Social media sites are playing key role in the modern information world. Text explosion on the web is spreading in variety of forms and this gigantic growth of technology is pumping a huge amount of data in social media. Various schemes and training mechanism are available to share and communicate with others which deal with the repetition of text instead of dealing with text only [26]. Micro-blogging is a common platform for all online users and has become very frequent in past few years [27]. As numerous social media sites are increasing the users are communicating, participating in events, sharing views, and receiving breaking news anytime from anywhere [49]. Therefore, social media sites are providing a conducive milieu for all to create and broadcast information. In the past epoch the traffic has turn out to be almost double on internet [5]. The intensification of usage of social media is rapidly growing. Table 2, show the statistics of social media's popularity among various age groups and certainly confirms its significance in internet era [54, 55].

Table 1. Fact Sheet of Social Media Users

Total users of social networking sites	74%
Overall %age of Men	72
Overall %age of Women	76

* Corresponding Author

Table 2. Variation in Usage According to Users

18-29	89%
30-49	82%
50-64	65%
65+	49%

Social media allows the users to deploy the services and at the same time offers the platform to contribute (Twitter, Facebook, Tumbler, Flickr, LinkedIn *etc.*) [6]. The flow of information was unidirectional, in long established media *i.e.*, B2C. But in the current scenario are absolutely changed, social media allows users to carve up the informational needs speedily. Among all the social media sites, Twitter is gaining momentum [47] as it can be accessed through mobile, website interface, SMS, *etc.*, 80% users are accessing through mobiles and it has 320M monthly lively users [60]. It allows users to post, interpret, share and deliver 140 words' called tweets [5].

Table 3. Variation of Users in different Social Medias

Facebook	71
Twitter	23
Instagram	26
Pinterest	28
LinkedIn	28

From all the internet users, surveys say that trend for accessing social sites is rising quickly and twitter is most famous among all social sites [55]. With the advancements variety of services came along which deal with mixture of data formats such as text, images, videos *etc.* Thus text mining is essential to extract the information from the gamut of social media websites. Information composed and published by citizens, journalists, social site users and consumed by thousands of individuals, who give spontaneous feedback to all who are generating text on the social media. Social media enables us to be connected and interact with each other, anywhere and anytime – allowing us to observe human behavior on an unprecedented scale with a new lens as the interaction trend of users' are increasing day by day on internet [9]. This social media novel world provides us the golden opportunities to understand individuals at scale and to mine human behavioral patterns otherwise impossible.

Unfortunately, social media data is drastically different from the traditional facts that we are recognizing within data mining. Apart from enormous size, the mainly user-generated data is noisy and with abundant social relations such as friendships and followers-followees [49]. This new type of data mandates new computational data analysis approaches that can combine social theories with statistical and data mining methods. The pressing demand for new techniques ushers in and entails a new interdisciplinary field social media mining [8]. With text analytics narrative information is represented by various text mining techniques which cannot be analyzed generally. Text mining is designed to help insights from text ended questions, CRM, reviews and sentiments analysis, event detection and prediction, trend analysis, information diffusion *etc.* [56, 57]. Various mining software of texts are used which transpose words from unstructured data into numerical values, pattern identification algorithm which helps to analyze the interested pattern of data. Text mining plays a significant role in summarizing the documents, extracting concepts from the text and indexing it for use in predictive analytics. Thus it is possible to extract the meaning from text in the social media and cluster documents of similar types. With this practice we can easily handle and process the contents of social media data.

Thus text mining is a “Combination of Machine learning (ML) and the Statistical techniques which are helpful to model and structure the data contents of textual resources for research, e-commerce data analysis and for further investigations or processing” [18]. The process of text mining can be visualized as two step progression such as Text refining and knowledge distillation. The former transforms the free-form of text documents into intermediate form which is a graph representation [11] or relational data representation and the later deduces the pattern from the intermediate pattern. Information Extraction (IE) helps to extract entities and keywords which are required for further processing and it helps to mine some valuable, knowledgeable pattern from collected keywords to take decision [34]. Text mining uses several techniques for the information extraction such as bag-of-words (BOW) model for document ranking, matching and clustering *etc.* An overview of the information extraction system is depicted in the Figure 1.

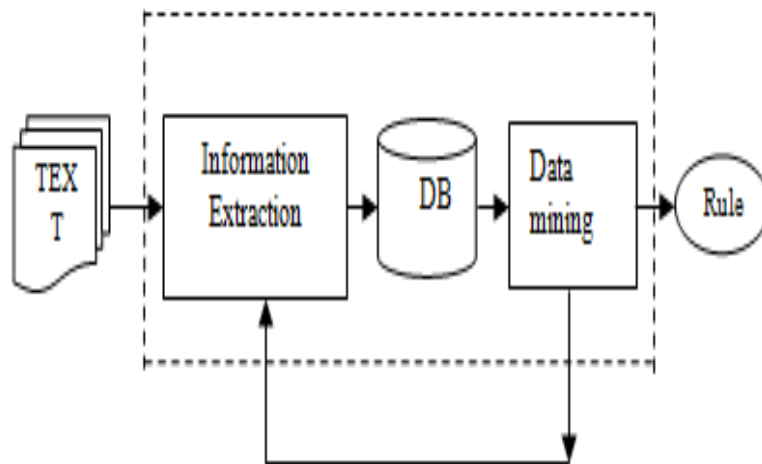


Figure 1. Overview of IE based Text

IE support business to facilitate decision making [45] which demands the normalization of the extracted textual pattern [13] into knowledgeable patterns which helps to take decision for strategic planning in business [45].

As Social media shatters the boundaries between the real world and the virtual world. With this we can integrates social theories and the computational methods together to find out that how individuals interact and how social molecule communities form [29]. The inimitability of social media data calls for novel data mining techniques that can cooperatively knob user generated content with affluent social relations [37]. The study, development and implementation of these innovative techniques are under social media mining, an emerging discipline under the umbrella of data mining.

Social Media mining, introduce necessary concepts and principal algorithms suitable for investigating enormous social media data. It discusses theories and methodologies from different disciplines such as computer science, data mining, machine learning, social network analysis, network science, sociology, ethnography, statistics, optimization, and [37]. Social media mining encompasses the tools to properly represent, measure, model, and mine meaningful patterns from large-scale social media data. It cultivates an inventive kind of data science which is well versed in social and computational theories, specialized to analyze unruly social media data [46], and skilled to help bridge the gap from what we want to know about the infinite social media world with computational tools. According to [61] the internet traffic of September, 2016 social media websites can be classified as describes in Table 4.

Table 4. Top 10 most Popular Websites in the World September 2016

Rank	Website
1	Google
2	Facebook
3	YouTube
4	Yahoo
5	Baidu
6	Amazon
7	Wikipedia
8	Taobao
9	Twitter
10	Tencent QQ

1.1. Social Media Categories

With the proliferation of users on the web arise several platform on which large amount of information is shared and extracted among them. Figure 2 depicts various social media categories which enhances the opportunities of social media users to share their views in the form of text [62].

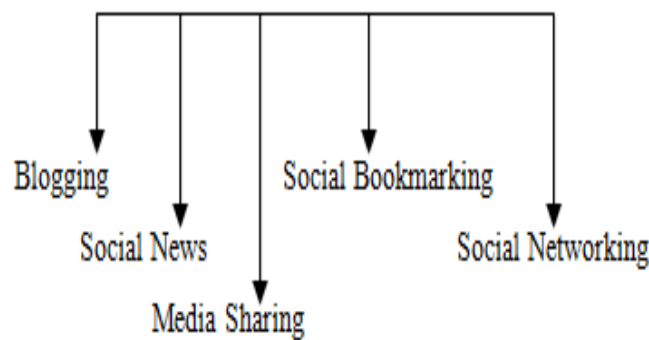


Figure 2. Types of Social Media

The social news platform, allows users to post news links and other information on a web page of a news article. Users vote on the different news items shown on the web page and the highest rating news are outstandingly displayed as headlines [7]. Mostly social media sites not only include text but also picture videos and numerous multimedia forms, these sites are known as multimedia sharing sites which allow users to utilize multiple social features like creating a profile, uploading multimedia data and addition of textual comments. YouTube is the most popular media sharing site in the world where vast number of users are uploading pictures and movies. In the similar way blogging is a social media where the users can publish their information and views on the web [10]. People write information, share links, provide suggestions and can even report news. As the information is pooled on blogs very frequently, this helps to televise information anytime which is proven to be a fast mode of communication broadcasting. Blog information can be utilized to cater the information seeking behavior of any user through analysis of blogs by opinion mining. Social bookmarking is to get ideas about web pages by chipping in each other through tags, URL's links etc. The bookmarking search can be done through different bookmarks which are based on link organization, Query/Click through log data, user generated data and page contents.

With bookmarking users can store, categorize, manage, depict and divide links to those websites, blogs, videos, text etc. which are popular and interesting [14]. The example of the popularity of social media sites can be predicted from the monthly visit of users on

Facebook and twitter (900,000,000 and 310,000,000, estimated unique monthly visitors on Facebook and on twitter) [59].

Thus these all social media helps to advertise the business, news, articles, events and information according to trend all over the world for marketing business perspective Sentiment analysis (SA), Sentiment Polarity Disambiguation (SPD), Event Detection and Tracking, Trend Analysis and prediction, Spammed Profile Detection, Information Diffusion etc. are challenging and promising research areas in text analytics which are major thrust of this article. This paper is organized as, Section 2 expounds on challenges in text analysis in social media, preprocessing of such text along with normalization methods are explained in Section 3, representational techniques for this data for effective utilization is presented in Section 4, Section 5 describes applications areas of social media text mining and Section 6 summarizes the tools available to carry out this sort of text mining.

2. Challenges in Text Analysis in Social Media

There are three major steps in a general framework of text analysis which include pre-processing of text, representation of text and discovery of knowledge [31].

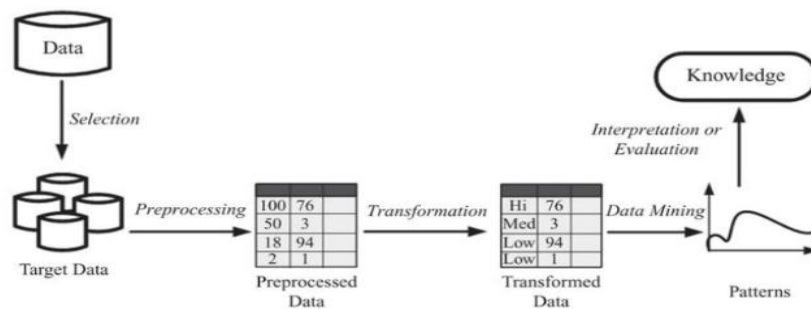


Figure 3. Traditional Framework of Text Analysis

The figure given above (Figure 3) [37] illustrates the traditional text analysis framework.

This framework allows us to extract the information from the text corpus of Facebook, Twitter, and LinkedIn etc. Since it is a very fundamental framework of text analysis therefore it does not ensure the quality of text. As standards of quality for social text is quite different from the traditional text quality thus it is challenging to find parameters of quality to quantify and fit them in the traditional analytical framework.

2.1. Text Acquisition

To acquire text from the social media, text acquisition can be retrieved from existing text corpus or from online text streams. Acquired text is further analyzed, mined, and summarized. NLP methods and techniques provide a strong support to text mining tools to get closer to the semantics of text source [40]. Entities, attributes facts and events are necessary to acquire from the social text.

- **Information Extraction (IE)**

It is a restricted form of NLP (Natural Language Processing), this extraction refers to get useful patterns of text with functional annotation or explanation. In IE it is usually known in advance what information is needed and which part of text is extracted to fill in slots of a predefined template. In extraction process named entity recognition (NER), relation extraction (RE), and co-reference resolution [52] SEO and shallow NLPs are playing a key role. IE methods have become more complicated due to advancement in

NLP and machine learning. All these methods and technologies are used to extract required patterns from text. The extracted data usually contains a large chunk of noise and this noise is required to be removed from the extracted patterns. For this purpose IE includes various tasks which are Named Entity Recognition (NER), Co-reference Resolution (CO), Relation Extraction (RE), Event Extraction (EE) [21] Trend Analysis and Prediction [28].

3. Normalization of Social Text

3.1 Text Quality and Process of Normalization

Text normalization is essential part in social media text analytics to extract meaningful pattern from unrefined noisy text [52] and then refined text helps in sentiment analysis, sentiment polarity disambiguation, event analysis & tracking, trend analysis & prediction *etc.* Various metrics affect the quality of social media text [15]. In order to have better findings from text, data should be normalized. If we consider a twitter community, then every user have own words to write or post something on twitter. Thus posted words are combination of non-standard words, slang or Folksonomies [49]. Thus it is necessary to handle the following properties of Twitter text to improve the quality of text for further processing; a brief account of this aspect is summarized in Table 5.

Table 5. Text Quality Mterics

Challenge	Description
Stop List	Common words frequency of occurrence
Lemmatization	Similarity detection of text/words
Text Cleaning	Removal of unwanted from the data
Clarity of Words	To clear the meaning in text
Tagging	Predicting data annotation and its characteristics
Syntax/Grammar	Scope of ambiguity, data dependency
Tokenization	Various methods to tokenize words or phrases
Representation of Text	Various methods and techniques to represent text
Automated Learning	Similarity measures and use of characterization
Emoticons	Various symbols to express happiness, anger, sad, etc
Folksonomies	Various user's created taxonomies

3.2. Cleaning of Data

Unwanted symbols, abbreviations, spelling mistakes, emoticons, tags (@ tags or # tags) [50] *etc.*, spoils the quality of text. Such kind of words and symbols are usually found on every popular social media forum like Twitter, Facebook, news group blogs, SMS *etc.* Informal semantics and slag words are very famous and handy to use in all media [47]. Without care of grammar, semantics, words representation people post informal text freely but it poses a dire demands of normalization for the analysis to convert informal language into meaningful form [38]. A variety of unnecessary belongings which hassle to be removed from the text are escaping HTML *characters*, Stop-Word-Removal, Punctuation, Removal of URL, @ and # symbols, Special characters, spell checking and correction, Slangs Lookup *etc.* Various assets are available on web which helps to achieve the target of clean text.

Therefore in every social media form the main task it to reduce the impact of noisy text and its issue related to sparseness of textual features. In sentiment analysis (SA), Sentiment polarity disambiguation (SPD), event identification and tracking (EIT), trend analysis and prediction (TAP) normalized text is demanded for better performance.

4. Representation of Text

Text in social media is represented in various forms. Most common and easy form of representing text is bag of word (BOW) model which capture word frequency. On the other hand graph based text representation is related to mathematical terms [32]. Many relational and structural information can be gathered from graphs. In graphs a textual keyword can be represented using vertex as a feature term. Graph based text representation helps in information retrieval applications because it provides easy computations related to various operations like term weight ranking *etc.* [49]. Vector Space Model (VSM) is also used to represent social text but due to few disadvantages like structure of text and meaning cannot be expressed in a fine manner makes this representation less useful. On the other hand if the document has same words then it is not easy to compute the value of textual words and relation among the text cannot be expressed quantitatively. Text on the Web has remarkable discrepancy, covering a broad variety of topics interests and viewpoints. It includes news pages, blogs, and documents on the Web etc. These web pages could be in any language, which complicates an already challenging text mining problem. An early approach for dealing with documents in an information retrieval (IR) setting was the vector space model (VSM). This is the most common method to model documents and web pages to sparse numeric vector and then deal with them with linear algebraic operations. Consider an example of BOW (Bag of Word) which includes two simple text documents as:

***Leena likes to watch movies. Honey likes movies too.
 Leena also likes to watch football games.***

From these two sentences dictionary can be constructed as:

```
{
  "Leena":1,
  "likes":2,
  "to":3,
  "watch":4,
  "movies":5,
  "also":6,
  "football":7,
  "games":8,
  "Honey":9,
  "too":10
}
```

This includes the ten distinct words and by using the indexes of the dictionary each document is represented by a ten entry vectors defined as:

```
[1,2,1,1,2,0,0,0,1,1]
[1,1,1,1,0,1,1,1,0,0]
```

In BOW model, a word is represented as a separate variable having numeric weight of varying importance. The most popular weighting schema is Term Frequency / Inverse Document Frequency (*TF-IDF*). It gives the numeric statistic which defines the importance of word is mostly used in text mining for generating a novel information.

Term Frequency / Inverse Document (*tf-idf*) is the product of two statistics *i.e.*, term frequency and inverse document frequency. In term frequency $tf(t,d)$, the simplest choice is to use the raw frequency of a term in a document, *i.e.*, the number of times that term t occurs in document d . If we denote the raw frequency of t by $f(t,d)$, then the simple tf scheme is $tf(t,d) = f(t,d)$ [58].

Thus most useable schema of Term Frequency / Inverse Document (*tf-idf*) is defined as:

$$tfidf(w) = tf * \log \frac{N}{df(w)}$$

4.1. Similarity Functions in Text Mining

In in order to efficiently capture deep insights of underlying patterns from data it is of paramount importance to analyze relatedness and correlation of phrases and sentences, therefore similarity functions are essential to represent text social media which includes:

- **Jacard Similarity Index** which defines the size of intersection divided by size of union of the sample set.

$$j(X, Y) = \frac{|X \cap Y|}{|X \cup Y|} = \frac{|X \cap Y|}{|X| + |Y| - |X \cap Y|}$$

If 'X' & 'Y' both are empty set then,

$$j(X, Y) = 1 \quad 0 \leq j(X, Y) \leq 1$$

Jaccard distance which measure dissimilarity between data sets is represented as:

$$d_j(X, Y) = 1 - J(X, Y) = \frac{|X \cup Y| - |X \cap Y|}{|X \cup Y|}$$

- **Cosine Similarity**

It is a similarity between two non -zero vectors which represent text. It is computed by using Euclidean dot product as:

$$X \cdot Y = \|X\| \|Y\| \cos \theta$$

$$similarity = \cos \theta = \frac{X \cdot Y}{\|X\| \|Y\|} = \frac{\sum_{i=1}^n X_i Y_i}{\sqrt{\sum_{i=1}^n X_i^2} \sqrt{\sum_{i=1}^n Y_i^2}}$$

Where 'X_i' & 'Y_i' are component of vector X & Y respectively.

4.2. Graphical Representation of Data

Graph contains set of objects called nodes and the connections between nodes are known as edges [11]. In general terms' graph G is represented as: $G(V, E)$, where, V is set of nodes & E represents set of edges. In graphical representation of text, terms are represented by *vertices* (V) and relations between terms by *edges* (E) [42], each vertex is defined as, $V = \{F, S, P, D, C\}$, where F= Feature Term, S= Sentence, P= Paragraph, D= Document, C= Concept [48]. A weighted generally graph G contains 4 tuples as: $G = (V, E, \eta, \kappa)$ where V defines the vertices and E defines edges. $E \subseteq V \times V$ is combining set of edges which connects V, $\eta: V \rightarrow X_v$ is a vertices labeling function whereas $\kappa: V \times V \rightarrow X_e$ is edge labeling function (with X_v and X_e are the sets of labels that can appear on the vertices and edges) . Thus a text document T is a combination of

$T = \{g_1, g_2, \dots, g_n\}$ where $\{g_1, g_2, \dots, g_n\}$ is defining number of graphs in T and unified graph G is defined as : $G = \cup\{g_1, g_2, \dots, g_n\}$.

Data mandates novel computational data analysis approaches that can combine social theories with statistical and data mining methods. Unfortunately, social media data is drastically different from the traditional computational methods.

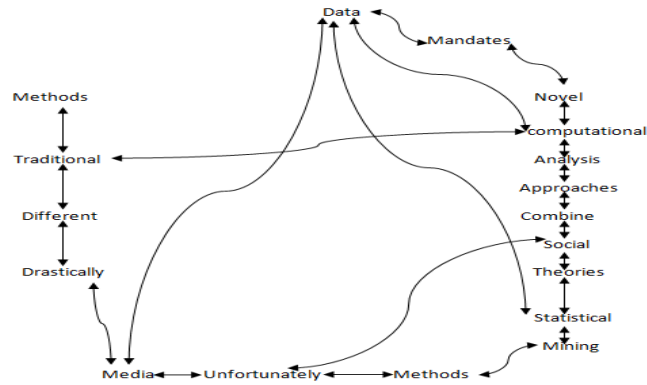


Figure 4. Representation of Text as a Graph

Directed, undirected, weighed or un-weighted graphs are used to represent text. In the undirected graph no information about the flow (quantitative information about nodes' connectedness) is described. Thus edges between two vertices X and Y are identical to the edges between vertices Y and X. Unidirectional are represented as $X \longrightarrow Y$, whereas bidirectional representation is $X \longleftrightarrow Y$.

Figure 4. shows the graphical representation of data where it is represented in the form of graph as depicted in figure given below:

5. Application Areas of Text Mining

5.1. Sentiment Analysis

Sentiment analysis extracts a massive amount of intriguing information from the social media world with the intent to identify the sentiment of users (users information sharing behavior) [44] emotion, attitude or opinion about event, area of interest [22]. This helps to understand and predict future market values and trends of a particular commodity under consideration. Twitter developer's account analysis provides a good opportunity to extract data stream using twitter API (Application Programming Interface) which helps in sentiment analysis profoundly. Research on sentiment analysis in social media text can be further classified in broad categories as machine learning and lexicon based sentiment analysis.

- **Machine Learning**

Sentiment analysis in this domain can further be divides as supervised and unsupervised learning, where supervised techniques include decision tree, linear, rule based and probabilistic classifiers. Polarity scores of sentiments can be predicted using machine learning approaches. In supervised learning labeled class of text documents are needed. Thus normalization or cleaning of text is required to extract informative and decision making patterns from large corpus of text or stream of text. But unsupervised learning can be operated on unlabeled dataset. These methods analyze the text their

features and then arrange/group the text according to their similarity. Clustering of similar documents is performed to collect similar keywords from large text collection.

- **Lexicon based**

Lexicon based approach of sentiment analysis exploits hierarchical information reoccurring and prominent dictionary definitions of lexicon present in corpus. It can be further divided into two broad categories as statistical and semantics based approaches. Lexicon based methods are quite fast but they are not suitable for analysis of text streams as they cannot predict changes over time during streaming. Linguistic methods are also used for sentiment analysis and sentiment polarity disambiguation. In sentiment analysis (SA) of twitter, sentiment can be formally predicted from the corpus or from the stream of text, which can be defined as function $F: T \times C \rightarrow \{0, 1\}$, where T is the set of all possible tweets and C is the set of predefined categories. The value of $F(t, c)$ is 1 if the tweet t belongs to the category c and 0 otherwise. The approximating function $M: T \times C \rightarrow \{0, 1\}$ is called a classifier. Thus in sentiment analysis classifier helps to produce results as "close" as possible to the true category assignment function F . In sentiment analysis (SA) and sentiment polarity disambiguation (SPD) of non-standard, slang and local words are converted into actual word format. Text streams are used for extracting trendy events on query based method e.g., Twitter TAP, provides online streaming of text with Python API of Twitter. Mathematically event can be defined as a function $F: E \times Q \rightarrow \{0, 1\}$, where E belongs to set of all possible event extracted from twitter streaming and Q is a set of predefined query. The value of $F(e, q)$ is 1 if the event e belongs to the query q and 0 otherwise. The function $S: E \times Q \rightarrow \{0, 1\}$, is classifier. In trend analysis prediction influence of particular event is analyzed for Text Analysis and Prediction. This depends upon the probability of influence (p) of a trend (T) in time interval $[t_1, t_2]$ which is a closed interval of t_1, t_2 of twitter streaming. Probability function P is used to calculate the influence of trend T at time t_1 to t_2 . This p will analyze the probability of T within the text streaming and predict the influence of the particular trend T which will define for how long T will be in trend.

- **Social Text Streams**

In data from social media can be processed and analyzed in two ways i.e. offline processing and online stream processing. In offline processing existing data sets are used for various tasks like normalization, sentiment analysis, sentiment polarity disambiguation, event detection, tracking etc. On the other hand if we consider the online processing fast clustering, trend analysis and prediction, rapid text summarization are the major tasks which are performed frequently.

There are two major aspects social text streams analysis, streams extraction and filtering these streams to group clusters for relevant tasks. Twitter API is used for tweet stream extraction which is based on user query matching. Generally stream s is extraction of tweets T in closed time interval of time, $[time_1, time_2]$, which contains sequence of text called as $T = (t_1, t_2, \dots, t_i, \dots, t_r, \dots, t_n)$. A stream sequence S is a combination of all the streams collected during this interval, $S = \{s_1, s_2, \dots, s_i, \dots, s_r, \dots, s_n\}$. There may be some part of S such that s_i contains a tweets t_i , which may be or may not be linked to other entities in the social network such tweets should be scraped or ignored from the s_i . For this filtering generally fast clustering methods are employed.

As clustering of text stream is necessary to categorize the similar contents in the summarized form of clusters. For this purpose text streams $S = (s_1, s_2, \dots, s_i, \dots, s_r, \dots, s_n)$, is distributed among clusters (C) which is defined as $C = (C_1, C_2, \dots, C_j)$ where, $C_1, C_2, \dots, C_j \in S$, such that $S = \cup(c_1, c_2, \dots, c_j)$. Each of the object s_i , belongs to the cluster c_j where $c_j \subset C$. Clustering is based on similarity function $sim_f(s_f)$ which

holds the information about content changed within the clusters and change in stream of text in closed time interval $[time_1, time_2]$ as well.

5.2. Detecting Popular Events and Trend Prediction

Social media platform provides ample opportunity for people to discover, report, write and communicate with other people on single platform about famous events and trends [20]. As social media belongs to dynamic networks class, structure of network of changes gradually not only in terms of nodes but also in contents present in these nodes especially in terms of text. So there are chances of nodes' change due to occurrence of events over a network in time T_1 to T_2 . Most of the events are reproduced by tweets which are related to time and space or at least one of the two dimensions [51]. The twitter event E is real world occurrence at time Γ_E and T_E are related tweets to E posted during Γ_E where $\Gamma_E \in (\Gamma_{E1}, \Gamma_{E2}, \dots, \Gamma_n)$ (occurrence of event in social Twitter stream). There are two types of events occurs in social media which falls into two broad categories as planned and unplanned events. Twitter's PE (Planned Event) consist of predefined topic X or hash tags (#) and Γ (time) at which PE is planned to occur. The event detection methods for UPE (Unplanned Events) generally employ a filter mechanism to select important features signifying characteristics of event among individual or additional characters of consequent tweets T_{UPE} which exhibit burst pattern during event time period Γ_{UPE} .

Trend A belongs to A_q (queries) in S (social stream) is a collection of tweets posts T_p which are associated with A in closed time interval $[t_1, t_2]$. A_q is analyzed from the α (peek time) for A by selecting a cluster with large number of tweets related to A_q . In general A may or may not be popular at various times interval as sometimes it may persists for several hours and sometimes for several days. Thus it is necessary to find the peak time α for analyzing A . Trend can also be classified into two categories [20] exogenous and endogenous. Exogenous trends are related to broadcasting media such as news. whereas endogenous trends are related to re-tweet, group or communities activities etc.

5.3. Social Computing

This computing area of research involves a profound mix of areas like computer science, social interaction theories behavioral science and economics etc. With the information systems' growth and use of social computing is quintessential wherever users' milieu is engaged [19]. The social computing in the context of social media is concerned with the intersection of social behavior and computational systems [17]. A variety of applications in social computing are famous like blogs, email, IM (Instant Messaging), social networking like Facebook, Twitter, LinkedIn etc. Ubicon is an advanced software tool for social computing [35] which helps to process data, store data and serve pipelining [30]. Thus the use of such social applications largely influences the lifestyle of people where the people are interacting socially. These all emerging trends of social media push the social computing to grow in modern world. In social media, text plays an important role [23] and at the same time the behavior and lifestyle of people is also changing because of it. Social media and social networks are not made for single person [41] but serves richly diverse communities including business and corporations. Therefore social computing models for online social media must incorporate not only the

mass's perception for certain trends but also the cause and effect relationships of these trends.

5.4. Opinion Review and Spam/Fake Opinion Detection

Every user who uses the services of social media and other online services shares their experiences and opinions about the services, produces or occurrence of events. 75,000 novel blogs and 1.2 million fresh posts producing opinion on products and services are generated each day [24]. Opinion holder also defines the advantages of events occurring in the media and their perception in opinioned text. There is a vast contribution of opinions from various social media sites as opinions play a significant role in e-commerce and e-business. With this contribution, fake reviewers are also increasing which imparts false value of goods, services and brands in reader's mind. A fake opinion misleads the reviewers and wrong decision can be taken. These fake reviews act as opinion spam in social media. Thus reviewers play an important role in the social media which help online users to make appropriate decisions about the variety of tasks including the hotel services, mobile network services, flight services, restaurants, sale and purchase etc. In order to deal with subjectivity and opinions present in social media an opinion is quantified as a quintuple

$$(o_j, f_{jk}, so_{ijkl}, h_i, t_i), \text{ where}$$

o_j is a target object.

f_{jk} is a feature of the object o_j .

so_{ijkl} is the sentiment value of the opinion of the opinion holder h_i on feature f_{jk} of object o_j at time t_i . so_{ijkl} is positive, negative or a more granular rating.

h_i is an opinion holder.

t_i is the time when the opinion is expressed.

[53] defined and studied the problem of automatically determining review quality using social information contents. The method which was proposed by them is quite generalizable and applicable for quality estimation of other types of user generated contents.

5.5. Spams in Profile

The major task of spammers is to distribute spam messages and promoting personal blogs, advertisements, pages etc. In social media, 45% of users click on the links which are shared by their social friends [25].

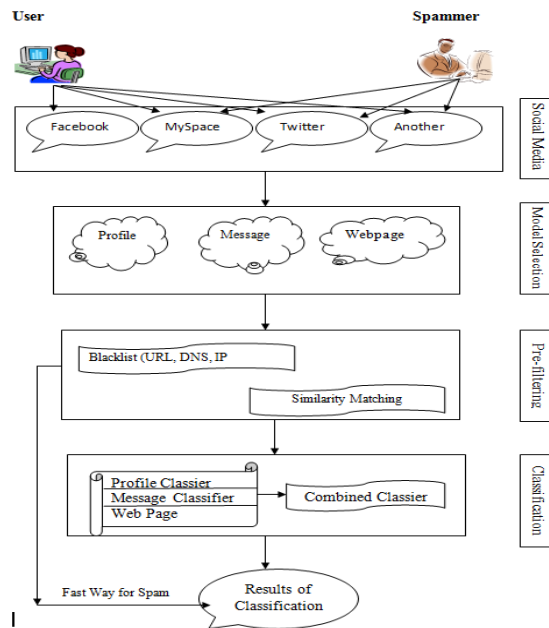


Figure 5. Framework for Spam Detection

This way of distribution of information on social media also attracts other parties which include social spammers; therefore spammers' role in social media cannot be ignored. Various strategies are applied by spammers for reaching into users' network of trust. In the spam messages adult contents, phishing attempts, and other such kind of contents are distributed which influences the users' trust on online activities. Spamming in social media is mainly attributed to spammed profiles to generate a feeling of genuine profile.

- **Social spam profile in social media environment**

With an identity users interact with social friends and their presence is represented by their social profile. This identity allows user to participate in various events and activities on social media [12]. Spammers analyze the behavior of social users to distribute spam using social networks. For this purpose spammer need profile and with this spammer can send friend request, recommend various groups and communities for collecting needful information. Few social media sites employ collaborative filtering for identification of social spam profiles.

- **Spam detection in social media profiles**

Pre-filtering is a process where the similar and repeated spam in social media can be classified and filtered out from spam receiving networks. Various techniques are involved in this process which includes *Blacklisting* of various entities such as URL, IP address or DNS [12]. Another profile spam detection method uses similarity matching which match the contents of a given profile with the previously known spammed profile. Hashtags are profoundly used for this purpose [12] which will perform this action in a quick manner. Classification can be used to filter the spam messages. Various standard algorithms such as Naïve Bayes, SVM (Support Vector Machine) and LogitBoost can be used for this classification purpose.

5.6. Information Diffusion

Generally, OSN are represented by graphs in which nodes act as users and edges describes relationships [1]. In OSN exchange of information is growing rapidly and

various social media like Facebook, Twitter, and LinkedIn etc. are serving to disseminate information of billion users [16]. People communicate ideas, opinions, videos, messages, pictures etc. within their friend circle and all over the world. This helps to examine social behavior of users to offer an opportunity to study the method of human communication with quantitative approach [2-3]. Process of Information Diffusion includes a *sender (S)* or a set of senders which forward information and a *receiver (R)* or a set of receiver which receive the sender's information. In general number of receivers is larger than sender; it also includes a *Medium (M)* through which information passed.

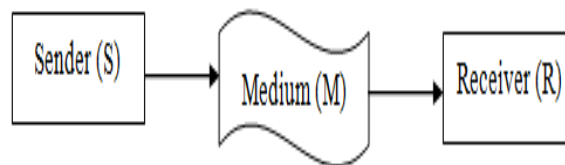


Figure 6. Information Diffusion Process

- **Various Forms of Information Diffusion**

On the basis of behavior, information diffusion can be divided into various forms (see Figure 7). As diffusibility of information decides its fate how far and wide it will travel, there are two major classes to quantify information diffusion models in social media as content centric and network centric. Huge works has been done on network centric paradigms through network connections' properties of graph and profile influence modeling, though content based information diffusion has not been explored so thoroughly. Areas like sentiment based information diffusion promises a great deal for content based information diffusion

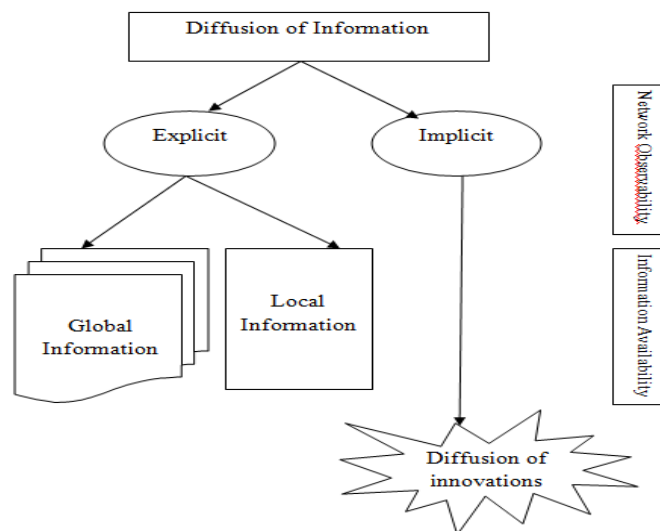


Figure 7. Various forms of Information Diffusion process

7. Tools

This section briefly describes the tools to perform and explore above explained tasks and application on social media texts, some of them are open source software some are commercially available, Table 6 and 7 list them respectively.

Table 6. Features of Commercial & Open Source Data Mining Tools

S. No	Product	Preprocess	Associate	Cluster	Summarize	Categorize	API
Commercial							
1	Clearforest	✓	✓	✓	✓		
2	dtSearch	✓	✓		✓		
3	Insightful Infact	✓	✓	✓	✓	✓	✓
4	Inxight	✓	✓	✓	✓	✓	✓
5	SPSS Clementine	✓	✓	✓	✓	✓	
6	SAS Text Miner	✓	✓	✓	✓	✓	
7	Word Stat	✓	✓	✓	✓	✓	
8	Discover Text	✓		✓	✓	✓	✓
Open Source							
1	R	✓	✓	✓	✓	✓	✓
2	Rapid Miner	✓	✓	✓	✓	✓	✓
3	Weka	✓	✓	✓	✓	✓	✓
4	KNIME	✓		✓	✓		✓
5	GATE	✓	✓	✓	✓	✓	✓
6	KH Coder	✓	✓	✓	✓	✓	✓

Table 7. Open Data Mining Tools Description

S. No	Tool	Developer	Programming Language	License	Purpose	Operating System	GUI/Comm and Line
1	R	Worldwide development	C, Fortran, R	free software, GNU GPL 2+	sci. computation and statistics	Cross Platform	both
2	Rapid Miner	RapidMiner, Germany	Java	open s. (v.5 or lower); closed s., free Starter ed. (v.6)	general data mining	Cross Platform	GUI
3	Weka	Univ. of Waikato, New Zealand	Java	open source, GNU GPL 3	general data mining	Cross Platform	both
4	KNIME	KNIME.com AG, Switzerland	Java	open source, GNU GPL 3	general data mining	LINUX OS X, Windows	GUI
5	Carrot2	Carrot Search	Java	BSD Licence	Clustering and Data Mining	Cross Platform	GUI
6	KH Coder	Koichi HIGUCHI	Perl	GPL Ver2	Text Mining	Cross Platform	GUI
7	Orange	Univ. of Ljubljana, Slovenia	C++, Python	open source, GNU GPL 3	general data mining	Cross Platform	both

8. Conclusion

This paper is an attempt to elucidate text analytic bent of social media with recent updates and text mining technique. As social media has evolved in unprecedented way, it led to many interesting research direction especially in the field of text mining. Thus, the main contribution of this article is to expound and conceptualize the domains of social media which are accessible on an extraordinary range. Rich patterns of online social text can be exploited to extract the relevant information effectively. Looking at the interestingness of the social media text analytics the complexity and veracity of techniques can't be ignored, therefore it encompasses variety of tasks such as semantic hashing for large information representation for event detection and classification, topic modeling and for fast similarity search and sophisticated language modeling techniques for deeper semantic understanding. This paper discusses these tasks along with different application areas. As a finding, text mining in social media includes fecund research areas such as:

- Semantic analysis
- Sentiment analysis
- Personalization
- Event detection and Trend analysis
- Human behavior modeling

Among these research areas the most prevalent and tough problem is to evolve some generic language models which are capable to handle multilingual idiosyncrasies, text's representational discrepancies and contextual resourcefulness of online social text in almost unsupervised manner.

References

- [1] Adrien Guille, Hakim Hacid, Cécile Favre, Djamel A. Zighed, Information Diffusion in Online Social Networks: A Survey, SIGMOD Record, June 2013.
- [2] A. Cho. Ourselves and our interactions: the ultimate physics problem? Science, 325-406, 2009.
- [3] A. Vespignani. Predicting the behavior of techno-social systems. Science, 325(5939):425–428, 2009.
- [4] A. Vespignani. Modelling dynamical processes in complex socio-technical systems. Nature Physics, 8(1):32–39, 2011.
- [5] Aston, N., Munson, T., Liddle, J., Hartshaw, G., Livingston, D. and Hu, W. (2014) Sentiment Analysis on the Social Networks Using Stream Algorithms. Journal of Data Analysis and Information Processing, 2, 60-66.
- [6] Ayushi Dalmia, Manwiter Thish Gupta, Vasudeva Varma, Sentiment Analysis: The Good, the bad, and the neutral, SemEval 2015.
- [7] Bogdan Batrinca, Philip C. Treleaven, Social media analytics: a survey of techniques, tools and platforms, Springer, 2014.
- [8] Charu C. Aggarwal, Karthik Subbian, Event Detection in Social Streams, Army Research Laboratory, 2012.
- [9] C. Aggarwal, Social Network Data Analytics, Springer, (2011).
- [10] Chao Shen, "Text analytics of social media: Sentiment analysis, event detection and summarization" ProQuest ETD Collection for FIU. Paper AAI3705112, January 1, 2014.
- [11] David F. Nettleton, Data mining of social networks represented as graphs, Elsevier, 2012.
- [12] De Wang, DaneshIrani, and CaltonPu, **A Social-Spam Detection Framework**, CEAS 2011 –Eighth annual Collaboration, Electronic messaging, Anti Abuse and Spam Conference Sept. 1-2, 2011.
- [13] Emma Haddi, Xiaohui Liu, Yong Shi, The Role of Text Pre-processing in Sentiment Analysis, First International Conference on Information Technology and Quantitative Management, Elsevier, 2013.
- [14] Enrique Estellés, Esther del Moral, Fernando González, Interdisciplinary Journal of E-Learning and Learning Objects, 2010.
- [15] Eleanor Clarka and Kenji Arakia, Text normalization in social media: progress, problems and applications for a pre-processing system of casual English, Procedia, Elsevier, 2011.
- [16] Eytan Bakshy, Information Diffusion and Social Influence in Online Networks, Doctor of Philosophy (Information) in The University of Michigan, 2011.

- [17] Fei-Yue Wang, *Toward a Paradigm Shift in Social Computing: The ACP Approach*, IEEE Society, 2007.
- [18] Feldman, R., Aumann, Y., Fresko, M., Liphstat, O., Rosenfeld, B., Schler, Y., *Text Mining via Information Extraction*, J.M. Zytlow and J. Rauch (Eds.): PKDD'99, LNAI 1704, pp. 165–173, Springer Berlin, Heidelberg, Germany, 1999.
- [19] F.-Y. Wang, K. Carley, D. Zeng and W. Mao, “Social computing: From social informatics to social intelligence”, *Intelligent Systems*, 22(2), pp. 79–83, 2007.
- [20] Hila Becker, *Identification and Characterization of Events in Social Media*, Thesis Doctor of Philosophy, 2011.
- [21] Jakub Piskorski and Roman Yangarber, *Information Extraction: Past, Present and Future*, Springer, 2013.
- [22] Jasmina Smailovic, *Sentiment Analysis in Streams of Microblogging Posts*, Doctoral Dissertation, November 2014.
- [23] Jinjun Chen, Jianxun Liu, *INTRODUCTION: SOCIAL COMPUTING AND SOCIAL NETWORKS*, *Journal of Organizational Computing and Electronic Commerce*, 24: 119–121, 2014.
- [24] Kim, P.: *The Forrester Wave: Brand Monitoring, Q3 2006*, Forrester Wave (white paper), 2006.
- [25] L. Bilge, T. Strufe, D. Balzarotti, and E. Kirda, “All your contacts are belong to us: automated identity theft attacks on social networks,” in *Proceedings of the 18th International Conference on World Wide Web (WWW'09)*. ACM, 2009.
- [26] Lifna C.S, Dr.Vijaya lakshmi M, *Identifying Concept-Drift in Twitter Streams*, ICACTA- 2015, Elsevier, 2015.
- [27] Lowri Williams, Christian Bannister, Michael Arribas-Ayllon, AlunPreece, Irena Spasic, *The role of idioms in sentiment analysis, Expert Systems with Applications*, Elsevier, 2015.
- [28] Luca Maria Aiello, Georgios Petkos, Carlos Martin, David Corney, Symeon Papadopoulos, Ryan Skraba, Ayse Goker, Ioannis Kompatsiaris, Alejandro Jaimes,” *Sensing Trending Topics in Twitter*, IEEE Transaction, 2013.
- [29] Mansoureh Takaffoli, Farzad Sangi, Justin Fagnan, Osmar R. Zaiane, *Community Evolution Mining in Dynamic Social Networks*, *Procedia*, Elsevier, 2011.
- [30] Martin Atzmueller, Martin Becker, Mark Kibanov, Christoph Scholz, Stephan Doerfel, Andreas Hotho, Bjoern-Elmar Macek, Folke Mitzlaff, Juergen Mueller and Gerd Stumum, Taylor & Francis, 2014.
- [31] Mena B.Habib, Maurice van Keulen, “Information Extraction for Social Media”, *Proceeding of third Workshop on Semantic Web and Information Extraction*, 2014.
- [32] Michael Gamon, *Graph-Based Text Representation for Novelty Detection*, *Workshop on Text Graphs*, 2006.
- [33] Michael W. Berry, Jacob Kogan, *Text Mining Applications and Theory*, Wiley, 2010.
- [34] Novita Hanafiah, Christoph Quix, *Entity Recognition in Information Extraction*, Springer, 2014.
- [35] N. Marz and J. Warren, *Big Data: Principles and Best Practices of Scalable Realtime Data Systems*, Greenwich, CT: Manning Publisher, 2013.
- [36] Parameswaran and A.B. Whinston, “Social computing: An overview”, *Communications of the Association for Information Systems*, 19, 2007.
- [37] Reza Zafarani, Mohammad Ali, Abbasi Huan Liu. *Social Media Mining*, Draft Version: April 20, 2014.
- [38] Sara Rosenthal, Alan Ritter, Preslav Nakov, Veselin Stoyanov, *SemEval-2014 Task 9: Sentiment Analysis in Twitter*, *Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014)*, pages 73–80, Dublin, Ireland, August 23-24, 2014.
- [39] Santhi Chinthala, Ramesh Mande, Suneetha Manne, Sindhura Vemuri, *Sentiment Analysis on twitter streaming Data.*, Springer International Publishing Switzzland 2015.
- [40] S. Jusoh and H. M. Alfawareh, “Agent-based knowledge mining architecture,” in *Proceedings of the 2009 International Conference on Computer Engineering and Applications*, IACSIT. Manila, Phillippines: World Academic Union, June 2009, pp. 602–606.
- [41] Syed K. Tanbeer, Carson K. Leung & Juan J. Cameron, *Interactive Mining of Strong Friends from Social Networks and Its Applications in Ecommerce*, *Journal of Organizational Computing and Electronic Commerce*, 24: 157–173, 2014.
- [42] S. S. Sonawane, Dr. P. A. Kulkarni, *Graph based Representation and Analysis of Text Document: A Survey of Techniques*, *International Journal of Computer Applications (0975 8887)*, Volume 96 - No. 19, June 2014.
- [43] Shamanth Kumar, Fred Morstatter, Huan Liu, *Twitter Data Analytics*, Springer, August 19, 2013.
- [44] Stefan Stieglitz & Linh Dang-Xuan, *Emotions and Information Diffusion in Social Media—Sentiment of Microblogs and Sharing Behavior*, Taylor and Francis, 2014.
- [45] Sumali J. Conlon, Jason G. Hale, Susan Lukose & Jody Strong, *Information Extraction Agents for Service-Oriented Architecture Using Web Service Systems: A Framework*, 04 Jan 2016.
- [46] Sukanya Dutta, Tista Saha, Somnath Banerjee, Sudip Kumar Naskar, *Text Normalization in Code-Mixed Social Media Text*, IEEE, 2015.
- [47] Tajinder Singh, Madhu Kumari, “Role of Text Pre-processing in Twitter Sentiment Analysis” *ICIP*, Bangalore, Elsevier Procedia, 2016.
- [48] Thi Ngoc Quynh Do, *A graph model for text analysis and text mining*, master thesis, 2012.

- [49] Xia Hu, Huan Liu, Text Analytics in Social Media, Springer, 2012.
- [50] Xia Hu, Jiliang Tang, Huiji Gao, and Huan. Liu, Unsupervised sentiment analysis with emotional signals., Proceedings of the 22nd international conference on World Wide Web, WWW'13, ACM, 2013.
- [51] Xiaowen Dong, Dimitrios Mavroudis, Francesco Calabrese, Pascal Frossard, Multi-scale Event detection in social media, Springer, 2015.
- [52] Youngmin PARK, Sangwoo KANG, Jungyun, Information Extraction Using Distant Supervision and Semantic Similarities, Advances in Electrical and Computer Engineering, 2016.
- [53] Yue Lu, Panayiotis Tsaparas, Alexandros Ntoulas, Livia Polanyi, Exploiting Social Context for Review Quality Prediction, ACM, 2010.
- [54] <http://www.pewinternet.org/factsheets/social-networking-fact-sheet>
- [55] <http://recruiterstoday.blogspot.in/2014/08/social-media-fact-sheet2014.html.VFmrr1fX7cs>
- [56] <https://textalytics.com/home>
- [57] <http://www.optimizationgroup.com/methods/text-mining/>
- [58] <http://en.wikipedia.org/wiki/Tf%E2%80%93idf>
- [59] <http://www.ebizmba.com/articles/socialnetworking-websites>
- [60] <http://twittercommunity.com>
- [61] <http://www.abcnews.com/top-10-most-visited-websites-in-the-world-2015/>
- [62] <https://blog.hootsuite.com/types-of-social-media/>
- [63] <http://www.iprospect.com/en/ca/blog/10-sentiment-analysis-tools-track-social-marketing-success/>
- [64] <http://marcobonanzani.com/2015/03/02/mining-twitter-data-with-python-part-1/>