

SUMMARIZATION AND TIMELINE GENERATION OVER EVOLUTIONARY TWEET STREAMS

¹Mr. Pravin V. Thakare, ²Prof. Harsha. R. Vyawahare, ³Prof. Manjusha M. Patil

Computer Engineering, SIPNA College of Engineering and Technology, Amravati, Maharashtra, India¹
Assistant Professor, Computer Science and Engineering, SIPNA College of Engineering and Technology,
Amravati, Maharashtra, India², Senior Training and Placement Officer, G. H. Raisoni Institute of Engineering
and Technology Wagholi, Pune, Maharashtra, India³

ABSTRACT

In recent years, number of users are being interested in the Social Networking site as well as micro blogging websites for example Twitter, Facebook etc. In a single day Twitter counts Tweets over a 500 million. The very complex part in this system is to control the real time applications sharing, keeping as well as managing like large data. Because of huge amount of data generated by the user, it goes from different concentrating problems like noisy as well as frequent information. By the researchers view, Querying as well as retrieval of like large information, it is very essential and one of the critical problems. Previous system tends to only work on the static as well as the limited information. Number of previous systems were tried to solve this problem and additionally given different solution over the problem. Summarization is a procedure consists of a text document in like a manner that short summary created through implementing the essential keywords of the original document. In this paper, we have tendency to propose the new method that builds the appropriate content-based summary in limited period of time. Our proposed system is also time efficient as compared with previous systems. We proposed a TCV-Rank summarization technique for generating online summaries and historical summaries of arbitrary time durations. Which monitors summary based/volume-based variations to produce timelines automatically from tweet streams; we design an effective topic evolution detection method. Our experiments on large-scale real tweets demonstrate the efficiency and effectiveness of our framework

Keyword:-*Continuous summarization, tweet clustering, summary, Tweet stream, TCV*

1. INTRODUCTION

In social networking, twitter also is a highly used and an online social networking service that permits every user to send as well as read short 140-character messages. Social networks as well as micro blogging services like Twitter, Facebook etc. are the reasons behind to the increment in case of generated huge quantity of short-text messages. Only Twitter receives 400 million tweets on a single day. There is not need to registration to read the Tweets, unregistered user can also able to read the Tweet Registered users post their tweets and unregistered users also read them. Twitter users have accessibility via the website interface, SMS or mobile device app. On the basis of analysis, now days, Twitter is a one of the ten most-visited websites as well as Twitter have more than 500 million users and out of 500 million more than 302 million are active users. The main benefit behind the Twitter is that tweets are implemented for data communication as well as sharing concepts and also individual's point of view. Previous systems are designed to work over the static as well as limited information. Actually, the real time applications tasks such as sharing as well as managing large data are very complex task. This huge amount of information produced by the social networking sites are goes from the complex problems like noisy and frequent information.

On the famous topic, Twitter may generate millions of tweets and extends over weeks. So, a user has to face the not required to visit the millions of tweets that are not even related with the user and it is not probable each time. Thus, to prevent this problem a method given called as filtering. If filtering is permitted then reaping for essential contents within like huge amount of tweets is also very critical task and this has to face due to immense amount of unrelated tweets. We have one another probable solution called summarization for data

overload issue. Summarization is a procedure consists with a text document in an order that short summary created through implementing the essential keywords over the original document. There is some of the search engines are implemented summarization methods for example Twitter, Facebook and Google etc. Summarization has other categories such as document summarization, image collection summarization and video summarization. The major concept of the summarization is to search a characteristic and common subset of the information that presents unique data of the whole set.

2. RELATED WORK

Tweet summarization method has gone through two stages. First step needs tweet information clustering and second step performs summarization.

Various authors are also determines the algorithm for stream data clustering in literature. BIRCH is an algorithm that manages iterative lessening and clustering by implementing hierarchy's algorithm. This algorithm is an unsupervised data mining algorithm [7]. These algorithms are specially implemented to execute continuous grouping on huge data sets. Implementation of BIRCH algorithm has a benefit that is, it has ability to create cluster in enhanced and fundamentally. It crates bunch of approaching and multi-dimensional metric data points for develop the better clustering for a provided set of resources, such as, memory as well as time limitations. CluStream is a popular between the most keep running of the ordinary stream clustering techniques. It has online little scale clustering portion besides detached from the net full scale clustering portion. The pyramidal time period was also proposed by authors to survey chronicled littler scale gropus for different time ranges [6].

A twitter post is as long as 140 characters in length and here we consider English posts. The twitter posts are casual, non-standard spelling and repeatedly do not have any accentuation. The hybrid TF-IDF based algorithm utilized for multi-post summaries of twitter post. Here few file abstract techniques are illustrated. Irregular Summarizer is a technique which self-assertively takes k posts or each subject for rundown. This framework was significant with a particular main objective to provide most cynical situation execution moreover set the lower bound of execution. Most recent Summarizer methodology takes the most recent k posts as summary from the determination pool. It may take the initial segment of a news article as summary. This system is executed in such a manner the brilliant summarizers can't perform better than anything fundamental summarizer. This summarizer only utilizes the initial segment of the report as outline [9]. Sum basic methodology uses clear possibilities of word with an upgrade ability to process the best k posts. This approach is useful in case of it completely depends upon the repetition of words in the initial content. It is sensibly great fundamental. Sum Basic framework was made by Nenkova and Vanderwende in 2005. This framework generates flat multi-record designs. Its design is impelled by the observation that, the words that are from time to time occurring within the chronicle group with higher probability [1].

Real-time event summarization [4] provides data about event at whatever point any sub-events happen. This strategy is a two-stage procedure to reporting sub-events happen. Initial step is to distinguish sub-events as of late happen and in second step tweets about sub-occasions are chosen. Later consolidating these two stages we get summary of game from set of tweets. LexRank summarizer uses a chart based framework. It recognizes pair wise equivalence between two sentences or between two posts. It creates the similarity score that is the largeness of the edge between the two sentences. The recent score of posts is designed in perspective of the weights of the edges that are connected with each other. This summarizer is helpful to give summary in perspective of graph instead of direct repeat layout. Regardless it relies on repeat; this framework uses the associations between sentences to incorporate more information. This is more difficult algorithm than recurrence based algorithm [2].

Text Rank summarizer [3] is one another graph based strategy. This technique uses the Page Rank algorithm. This provides another graph based summarizer that melds possibly a bigger number of information than LexRank. This occurs in the way that it repeatedly modifies the weights of posts. The recent score of each post is liable to how it is recognized with immediately related posts and the way in which presents and are associated over various posts. Text Rank joins the whole complexity of the graph rather than only pair wise similarities.

Twitter is a decent stage for individuals to express their point of view. The tweets are available in immense volume; it need push to comprehend what occur inside events. Here new strategies for summarizing events that discloses great correspondents and produces live sport upgrades from Twitter posts on events. Great correspondents chose logical tweets from dominant part of non-informative tweets [8]. ETS (Evolutionary Timeline Summarization) [10] is a web mining service that produces timelines for huge scale of information. ETS provides developmental directions over specific dates. ETS provides summery as indicated by score of timeline attributes. The benefit is that it encourages quick news browsing as well as information appreciation. ETS assignment has an adjusted optimization issue through iterative substitution.

Zhenhua Wang et al. introduce a summary framework called Sumblr. Sumbler is the predictable summary by stream clustering. This is the initial that determined constant tweet stream summarization. This framework involves three essential sections, to be particular the Tweet Stream Clustering module, the High-level Summarization module and the Timeline Generation module. Sumblr is profitable to work over dynamic, quick arriving and huge-scale tweet streams.

TCVs are considered as potential sub-topic delegates and maintained dynamically in memory during stream processing. The second structure is the pyramidal time frame (PTF), which is used to store and organize cluster snapshots at different moments, thus allowing historical tweet data to be retrieved by any arbitrary time durations.

3. IMPLEMENTATION

3.1 System Overview

Fig.1 demonstrates that proposed system contains three modules such as the tweet stream clustering module, the high-level summarization module and the timeline generation module. The tweet stream clustering module maintains the online statistical information. The topic-based tweet stream is given due to it has capability to effectively cluster the tweets as well as keep condensed cluster data. The high-level summarization separated within two kinds such as online and historical summaries. An online summary illustrates the trends between the public. So, the input for creating online summaries is accessed directly from the ongoing clusters kept in memory and a historical summary useful for people to determine the major communication within a particular period of time, so we have to remove the effects of tweet contents from the outside of that period of time. Thus, to produce historical summaries is more difficult due to retrieval of the needed data for producing historical summaries.

In the tweet stream clustering module, we design an efficient tweet stream clustering algorithm, an online algorithm allowing for effective clustering of tweets with only one pass over the data. This algorithm employs two data structures to keep important tweet information in clusters. The first one is a novel compressed structure called the tweet cluster vector (TCV). TCVs are considered as potential sub-topic delegates and protect dynamically in memory during stream processing. The second structure is the pyramidal time frame (PTF) [1], which is used to store and organize cluster snapshots at different moments, thus allowing historical tweet data to be retrieved by any arbitrary time durations. The high-level summarization module supports generation of two

kinds of summaries: online and historical summaries. (1) To generate online summaries, we propose a TCV-Rank summarization algorithm by referring to the current clusters maintained in memory. This algorithm first computes centrality scores for tweets kept in TCVs, and selects the top-ranked ones in terms of content coverage and novelty. (2) To compute a historical summary where the user specifies an arbitrary time duration, we first retrieve two historical cluster snapshots from the PTF with respect to the two endpoints (the beginning and ending points) of the duration. Then, based on the difference between the two cluster snapshots, the TCV-Rank summarization algorithm is applied to generate summaries. The core of the timeline generation module is a topic evolution detection algorithm, which consumes online/historical summaries to produce realtime/range timelines. In our design, we consider three different factors respectively in the algorithm. First, we consider variation in the main contents discussed in tweets. To quantify the summary based variation (SUM), we use the Jensen-Shannon divergence (JSD) to measure the distance between two word distributions in two successive summaries. Second, we monitor the volume-based variation (VOL) which reflects the significance of sub-topic changes, to discover rapid increases (or “spikes”) in the volume of tweets over time. Third, we define the sumvol variation (SV) by combining both effects of summary content and significance, and detect topic evolution whenever there is a burst in the unified variation.

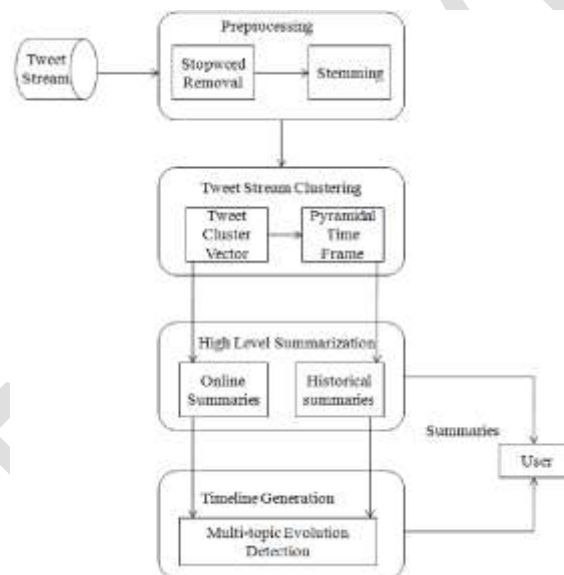


Fig. 1: System Architecture

3.2 Mathematical Model

Formulas:

Term Frequency $tf(d)$ of term t in document d

The number of times that t occurs in d .

Inverse Document Frequency estimate the rarity of a term in the whole document collection

$$idf_i = \log \frac{|D|}{|\{j : t_i \in d_j\}|}$$

Where $|D|$ = Total no: of documents

j = no: of documents containing the term t_i

$$\text{Cosine Coefficient} = \frac{|X \cap Y|}{\sqrt{|X|} \sqrt{|Y|}}$$

System S is represented as $S = \{T, D, TSC, V, P, S, G\}$

3.2.1 Process:-

(a) Input

1. Input Tweet Stream $T = \{t_1, t_2, t_3, \dots, t_n\}$

Where, T is the set of tweet streams and $t_1, t_2, t_3 \dots t_n$ are the number of streams.

2. Historical Tweeter datasets $D = \{d_1, d_2, d_3, \dots, d_n\}$

Where, D is representing as a set of Tweeter datasets.

3. Tweet Stream Clustering $TSC = \{V, P\}$

Where, TSC is represent as a set of Tweet Stream Clustering and $V = \{v_1, v_2, v_3, \dots, v_n\}$ Where, V is represent as a set of tweet cluster vector and $v_1, v_2, v_3, \dots, v_n$ number of vectors and is the sum of weighted textual vectors.

$P = \{p_1, p_2, p_3 \dots p_n\}$

Where, P is representing as a set of pyramidal time frame and $p_1, p_2, p_3 \dots p_n$ number of frames.

$C = \{c_1, c_2 \dots c_n\}$

Where, C is set of clusters generated using K-means algorithm.

4. High Level Summarization using TCV rank summarization algorithm.

$S = \{O, H\}$

Where, S is represent as a set of High Level Summarization, O= online summaries and H= Historical summaries

5. Timeline Generation using Topic detection evolution algorithm.

G= Topic Evolution Generation.

(b) Output

1. Multi topic summarization

3.3 Algorithms

3.3.1 Algorithm 1: Multi-topic version of Sumblr

Input: Multiple tweets (online or dataset), Number of cluster k;

Output: summary of Multiple Topic

Process:

1. While!topic.end () do
2. Topic t= topic.next ();
3. Study continuous tweet stream summarization.
4. For Tweet StreamClusteringmodule run Algorithm 1
5. Input the clusters CL generate using Algorithm 1 to Algorithm 2
6. ForHigh-level Summarization module run Algorithm 2
7. ForTimeline Generation modules run Algorithm 3.
8. Get the output of algorithm as Summary of multiple Topics.
9. END

3.3.2 Algorithm 2: Tweet stream clustering

Input: a cluster set C_set

Output:

1. While! stream.end () do

2. Tweet $t = \text{stream.next}()$;
3. Choose C_p in C_set whose centroid is the closest to t ;
4. If $\text{MaxSim}(t) < \text{MBS}$ then
5. Create new Cluster $C_{new} = \{t\}$
6. $C_set.add(C_{new})$
7. Else
8. update C_p with t
9. If $\text{TS}_{current} \% (\alpha_i) == 0$ then
10. Store C -set into PTF.

Algorithm 3.3.3: TCV-Rank Summarization

Input: a cluster set $D(c)$

Output: a summary set S

1. $S = \emptyset, T = \text{All tweets}$
2. Build a similarity graph on T ;
3. Compute LexRank scores LR ;
4. $T_c = \text{tweets with the highest LR in each cluster}$;
5. While $|S| < L$ do
6. For each tweet t_i in $T_c - S$ do
7. Calculate v_i ;
8. Select t_{max} with the highest v_i ;
9. $S.add(t_{max})$;
10. While $|S| < L$ do
11. for each tweet t_i in $T - S$ do
12. Calculate $v_0 i$;
13. Select $t_0 max$ with the highest $v_0 i$;
14. $S.add(t_0 max)$;
15. Return S ;

Algorithm 3.3.4: Topic Evolution Detection

Input: tweet stream binned by time units

Output: timeline node set TN

1. $TN = \emptyset$;
2. While $\text{stream.end}()$ do
3. Bin $C_i = \text{stream.next}()$;
4. If $\text{hasLargeVariation}()$ then
5. $TN.add(i)$;
6. Return TN ;

4. RESULTS AND DISCUSSION

Figure 2 graph demonstrates time comparison between previous system and proposed system. Previous system need more time to create single topic summery but in the similar amount of time proposed system creates multi-topic summaries.

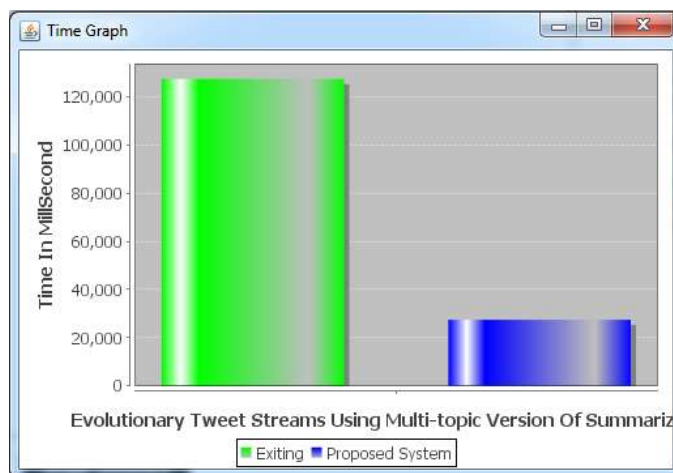


Fig. 2: Time Comparison graph.

Figure 3 graph demonstrates memory comparison graph. Memory graph demonstrates the memory comparison between the systems. Previous system needs more CPU consumption as compared with proposed system.

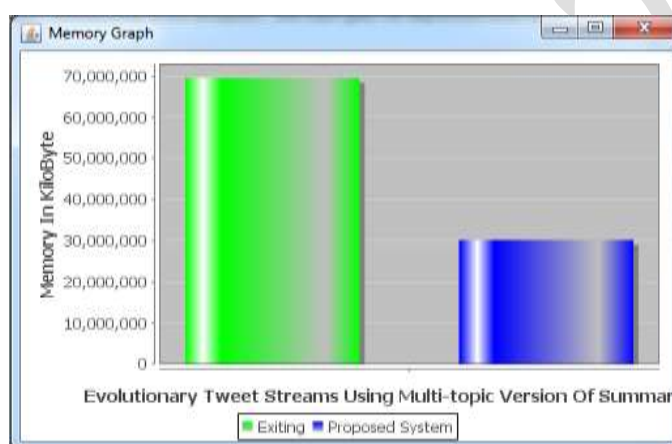


Fig. 3: Memory Comparison graph.

5. CONCLUSIONS and FUTURE SCOPE

In recent years, social networking sites as well as micro blog sites are the interested area for the user. Because of incremented utilization of social networking sites as well as micro blogging sites like Facebook, Twitter etc. it has generated huge quantity of short-text messages. In real-world to control the real time applications sharing, storing as well as managing huge amount of data is very complex. Due to huge amount of data it goes from number of critical problems like noisy as well as frequent information. In this paper, we analyzed different methodologies for document summarization like filtering as well as tweet summarization. These methodologies are implemented to control the large amount of tweets. But these methodologies have different problems like noise as well as repetitive information. So, to avoid these problems, there is requirement to build a dynamic methodology to summarize information generated by Twitter feeds. In proposed system, multi-topic summarization over online dataset is that needs minimum period of time as contrast with other previous systems.

REFERENCES

1. L. Vanderwende, H. Suzuki, C. Brockett, and A. Nenkova, Beyond SumBasic: Task-focused summarization with sentence simplification and lexical expansion, *Information Processing Management*, vol. 43, no. 6, pp. 16061618, 2007.
2. G. Erkan and D. Radev, Lexrank: graph-based centrality as salience in text summarization, *Journal of Artificial Intelligence Research*, vol. 22, pp. 457480, 2004.

3. R. Mihalcea and P. Tarau, TextRank: Bringing order into texts, in EMNLP. Barcelona: ACL, 2004, pp. 404411.
4. ArkaitzZubiaga, DamianoSpina, Enrique Amigó and Julio Gonzalo, Towards Real-Time Summarization of Scheduled Events from Twitter Streams , in Proc. 23rd ACM Conf. Hypertext Social Media,2012.
5. C. Shen, F. Liu, F. Weng, and T. Li, A participant-based approach for event summarization using twitter streams, in Proc. Human Lang. Technol. Annu. Conf. North Amer. Chapter Assoc. Comput. Linguistics, 2013, pp. 1152–1162.
6. C. C. Aggarwal, J. Han, J. Wang, and P. S. Yu, A framework for clustering evolving data streams, in Proc. 29th Int. Conf. Very Large Data Bases, 2003, pp. 8192.
7. T. Zhang, R. Ramakrishnan, and M. Livny, BIRCH: An efficient data clustering method for very large databases, in Proc. ACM SIGMOD Int. Conf. Manage. Data, 1996, pp. 103114.
8. Mitsumasa Kubo, RyoheiSasano, HiroyaTakamura, and Manabu Okumura, Generating Live Sports Updates from Twitter by Finding Good Reporters, in IEEE,2013.
9. David Inouye, Jugal K. Kalita, Comparing Twitter Summarization Algorithms for Multiple Post Summaries, IEEE Trans. Knowl. Data Eng., 23(8):12001214, 2011
10. Zhenhua Wang, LidanShou, Ke Chen, On Summarization and Timeline Generation for Evolutionary Tweet Streams, iee transactions on knowledge and data engineering, vol. 27, no. 5, may 2015.