

On Novel summarization and Time line Generation for Evolutionary Tweet Streams

Revathi
Student, MTech, (CSE)
DEPT OF CSE
Svp cet,

P. Sujitha,
Asst. Professor,
Dept. of CSE,
SVP CET,

V. Janardhan Babu,
Professor,
Department of CSE,
SVP CET, Puttur

Abstract: As increasing popularity of social sites like Tweeter, Facebook and Instagram etc. We will get lot of tweets and „N“ no. of short messages being shared at unpredictable rate which is very high. As this data is large enough it will become critical to understand and analyze therefore redundancy and noisy data must be removed. To overcome these drawbacks of existing system we propose sumblr framework, in comparison with other regular approaches of summarization which depends on static data and small datasets where sumblr is dynamic and works on large data set. Firstly we have proposed tweet cluster vector algorithm for maintaining statistical data and compact cluster information to maintain dynamically in memory during stream processing, store and organize cluster snapshots of different moments. Generation of online and historical summaries with arbitrary time durations, we propose TCV rank summarization algorithm. We have proposed an evaluation method which generates timeline, categorization based on topic evaluation.

Keywords: Tweet stream, continuous summarization, summary, timeline.

I. INTRODUCTION

Now a day a socially generated stream has become popular on WWW (World Wide Web). As rapid growth in an internet, use of social media also increases. There are many social sites like Twitter, Facebook, Instagram etc. in which twitter has become one of the most popular social site for users to share information like text, audio, video etc. Short messages are being created and shared at massive rate. Twitter receives thousands of tweets per hour. It is in raw form, the solution for this is summarization of tweets. Summarization represents a set of document which contain summary of related data. We have proposed Tweet Cluster Vector (TCV) algorithm which is used for making cluster of those retrieved tweets among which summarization will take place. Tweet Cluster Vector (TCV) algorithm includes two data structure to keep important tweet information in cluster. These data structures are tweet cluster vector and pyramidal time frame. TCVs are considered as potential sub topic representative and maintained dynamically during stream processing in memory.

The second data structure pyramidal time frame which used for storing and organizing cluster snapshot. So historical and online tweets data extracted by any random time duration which will give more relevant in results. In the summarization we will adding category such as news, politics, entertainment etc. we can summarize the tweets as per category. In the tweet summarization many tweets are repeated so using summarization we can avoid redundancy. The summarization consists of four issues efficiency, topic evolution, performance. Tweet streams or many messages of social site are very large in size so the summarization algorithm is very efficient. Performance of summarization is very effective. We have proposed TCV rank summarization algorithm which is used for generating historical and online summaries. This algorithm selects the top rank tweets from the Tweet Cluster Vector (TCV), to generate historical and online summaries where user specifies random time duration. We retrieve cluster snapshots from the Pyramidal Time Frame (PTF) with respect to beginning and ending of time duration, based on two clusters TCV rank summarization algorithm generates summaries. Also we proposed Topic Evolution Detection algorithm which takes the input of already generated summaries to produce timeline. Also we are working on other social stream which include clustering, timeline generation, Topic evolution etc.

2. LITERATURE SURVEY

We have studied the paper “**A framework for clustering evolving data streams**” (C. C. Aggarwal, J. Han, J. Wang, and P. S. Yu) in which TCVs are considered as potential sub-topic; for stream clustering, Clustream method is used. It includes online and offline micro clustering component. For recalling historical micro cluster, pyramidal time frame also proposed for random time duration. [1]
For using function lexrank in TCV rank algorithm we have studied “**LexRank: Graph based lexical centrality as salience in text summarization**” (G. Erkan and D. R. Radev) in this paper lex ranking is calculated. Depending on the similar data graph is created; Lexrank is used for finding top ranked tweets among large data set. [2] Also we referred, “**Text stream clustering based on adaptive feature selection**” (L. Gong, J. Zeng, and S. Zhang) worked on a various services on the Web such as news filtering, text crawling, etc. It mainly focuses on topic detection and

tracking (TDT). Clustering is used for analyzing text stream. [3]

Again we have studied paper “**Evolutionary timeline summarization A balanced optimization framework via iterative substitution**” (R. Yan, X. Wan, J. Otterbacher, L. Kong, X. Li, and Y.Zhang) evolutionary timeline summarization which consist of time stamped summaries which is used to generate timeline dynamically during the process of continuous summarization (Sumblr).[4]

For summarization we have studied “**Summarizing sporting events using twitter**” (J. Nichols, J. Mahmud, and C. Drews) in which Summarization algorithm creates sentence level summaries of important moments and then concatenated to generate an event summary of paragraphs. [5]

Lastly we have referred “**on summarization and timeline generation for evolutionary tweet stream**” we have referred Tweet Cluster Vector (TCV), TCV Rank algorithm, Topic evolution. In which TCV used for making effective clustering of tweet with the help of pyramidal time frame and tweet cluster vector, TCV rank summarization algorithm is used for generating online and historical summaries by evaluating top ranked function, depending upon top ranked tweets summarization is done. Topic evolution detection generates timeline by considering large variation of sub-topics in stream processing. [6]

3. PROPOSED SYSTEM

The system architecture is for historical and online summarization of social stream. In today’s world, summarization becomes necessity of social stream as millions of information posted on social sites. It is the simplest way to understand exact information using summarization by avoiding redundancy and noisy data.

Fig 3.1 mainly focuses on three module Stream Clustering, Summarization, and Timeline Generation. Here Categorization also done on summary generated.

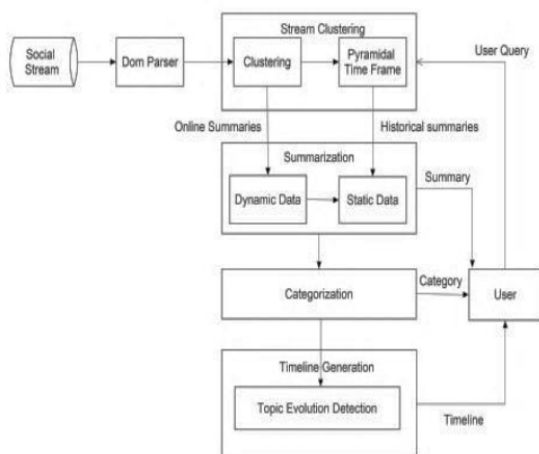


Fig 3.1 System Architecture

1. Stream Clustering

For making the cluster of social stream we use the clustering algorithm, by checking relevant data we are making cluster. [1]

1.1 Pyramidal Time Frame

Pyramidal time frame include time frame. Here user can provide the two points that is starting and ending time.

2. Summarization

After apply clustering on social data set, summarization is done by using TCV- Rank summarization algorithm. Static as well as dynamic data is summarized. [6]

3. Categorization

From generated summary, categorization is done like summary is relates to entertainment, sports, politics etc.

4. Timeline Generation

The base of the timeline generation is topic evolution detection algorithm which uses online summaries and generates timeline. Topic evolution describes changes in subtopics by monitoring variation in stream clustering. [6]

4.RELATED WORKS

4.1 Stream Data Clustering

Stream data clustering has been widely studied in the literature. BIRCH [2] cluster the data based on an in-memory structure called CF-tree instead of the original large data set. Bradley et al. [3] proposed a scalable clustering framework which selectively stores important portions of the data, and compresses or discards other portions. CluStream [1] is one of the most classic stream clustering methods. It consists of an online micro-clustering component and an offline macro-clustering component. The pyramidal time frame was also proposed in [1] to recall historical microclusters for different time durations. A variety of services on the Web such as news filtering, text crawling, and topic detecting etc. have posed requirements for text stream clustering. A few algorithms have been proposed to tackle the problem [1], [3], [4], [7]. Most of these techniques adopt partition-based approaches to enable online clustering of stream data. As a consequence, these techniques fail to provide effective analysis on clusters

formed over different time durations. In [2], the authors extended CluStream to generate duration-based clustering results for text and categorical data streams. However, this algorithm relies on an online phase to generate a large number of “micro-clusters” and an offline phase to re-cluster them. In contrast, our tweet stream clustering algorithm is an online procedure without extra offline clustering. And in the context of tweet summarization, we adapt the online clustering phase by incorporating the new structure TCV, and restricting the number of clusters to guarantee efficiency and the quality of TCVs.

4.2 Document/Microblog Summarization

Document summarization can be categorized as extractive and abstractive. The former selects sentences from the documents, while the latter may generate phrases and

sentences that do not appear in the original documents. In this paper, we focus on extractive summarization. Extractive document summarization has received a lot of recent attention. Most of them assign salient scores to sentences of the documents, and select the top-ranked sentences [9], [10], [11]. Some works try to extract summaries

without such salient scores. Wang et al. [12] used the symmetric non-negative matrix factorization to cluster sentences and choose sentences in each cluster for summarization. Heet al. [13] proposed to summarize documents from the perspective

of data reconstruction, and select sentences that can best reconstruct the original documents. In [14], Xu et al. modeled documents (hotel reviews) as multi-attribute uncertain data and optimized a probabilistic coverage problem of the summary.

While document summarization has been studied for years, microblog summarization is still in its infancy. Sharif et al. proposed the Phrase Reinforcement algorithm to summarize tweet posts using a single tweet [15]. Later, Inouye and Kalita proposed a Hybrid TF-IDF algorithm and a Cluster-based algorithm to generate multiple post summaries [16]. In [17], Harabagiu and Hickl leveraged two relevance models for microblog summarization: an event structure model and a user behavior model. Takamura et al. [18] proposed a microblog summarization method based on the pmedian problem, which takes posted time of microblogs into consideration. Unfortunately, almost all existing document/microblog summarization methods mainly deal with small and static data sets, and rarely pay attention to efficiency and evolution

issues. There have also been studies on summarizing microblogs for some specific types of events, e.g., sports events. Shen et al. [12] proposed to identify the participants of events, and generate summaries based on sub-events detected from each participant. Chakrabarti and Punera [13] introduced a solution by learning the underlying hidden state representation of the event, which needs to learn from previous events (football games) with similar structure. In [14], Kubo et al. summarized events by exploiting “good reporters”, depending on event-specific keywords which need to be given in advance. In contrast, we aim to deal with general topic-relevant tweet streams without such prior knowledge. Moreover, their method stores all the tweets in each segment and

selects a single tweet as the summary, while our method maintains distilled information in TCVs to reduce storage/computation cost, and generates multiple tweet summaries in terms of content coverage and novelty. In addition to online summarization, our method also supports historical summarization by maintaining TCV snapshots.

4.3 Timeline Detection

The demand for analyzing massive contents in social medias fuels the developments in visualization techniques. Timeline is one of these techniques which can make analysis

tasks easier and faster. Diakopoulos and Shamma [7] made early efforts in this area, using timelines to explore the 2008 Presidential Debates by Twitter sentiment. Dork et al. [8] presented a timeline-based backchannel for conversations around events.

In [9], Yan et al. proposed the evolutionary timeline summarization (ETS) to compute evolution timelines similar ours, which consists of a series of time-stamped summaries. However, in [9], the dates of summaries are determined by a pre-defined timestamp set. In contrast, our method discovers the changing dates and generates timelines dynamically during the process of continuous summarization. Moreover, ETS does not focus on efficiency and scalability issues, which are very important in our streaming context.

Several systems detect important moments when rapid increases or “spikes” in status update volume happen. TwitInfo [10] developed an algorithm based on TCP congestion

detection, while Nichols et al. [11] employed a slope-based method to find spikes. After that, tweets from each moment are identified, and word clouds or summaries are selected.

Different from this two-step approach, our method detects topic evolution and produces summaries/timelines in an online fashion

5. ALGORITHM

A. Clustering algorithm:

Input: Cluster set

Output: Assigning cluster for new tweets

Step I- Collection of new tweet stream

Step II- Depending upon two attribute it create new cluster

1) maxsim (maximum similarity)

2) mbs (minimum boundary similarity)

Step III- If is less than then it creates new cluster.

Step IV- otherwise update new cluster

B. Summarization algorithm:

Input: Cluster Set

Output: Summarization according to rank

Step I- Building similarity graph for all tweet

Step II- Computing LexRank to know which tweets are top ranked

Step III- Adding tweets into summary according to equation

Step IV- Checking summary length till it reached to max size.

Selecting tweet globally

C. Topic Evolution Detection algorithm:

Input: A tweet stream binned by time units

Output: A timeline node set

Step I- Binning tweets by time.

Step II- Appending new timeline nodes whenever large variation detected. By using

```
While!stream.end()do
```

```
Bin ci=stream.next()
```

```
If hasLargeVariation()Then
```

```
TN.add(i);
```

D. Web Extraction Algorithm

Step I- Recognizing peer node

Step II- It align nodes in peer matrix to get a list of aligned nodes

childList=matrixAlignment(M)

Step III- Repetition of pattern detected starting with length 1

childList=repeatMining(childList,1)

Step IV- Optimal merging

6. CONCLUSION

We proposed a Sumblr which supports continuous tweet stream summarization. Sumblr uses a tweet stream clustering algorithm for compress tweets into TCV and manages them in an online way. Then, it uses a TCV (tweet cluster vector)-Rank summarization algorithm for generating online and historical summaries with random time durations. Also categorization will be done on summarized data. The topic evolution will be done automatically, permitting Sumblr to create dynamic timelines for tweet streams. We are also working same for other social site also.

REFERENCES

- [1] C. C. Aggarwal, J. Han, J. Wang, and P. S. Yu. "A framework for clustering evolving data streams." *ACMSIGMOD Conference* (2003): 81-92.
- [2] Radev, G. Erkan and D. R. "LexRank: Graph-based lexical centrality as salience in text summarization." *J. Artif. Int. Res.* 22 (2004): 457-479.
- [3] L. Gong, J. Zeng, and S. Zhang. "Text stream clustering algorithm based on adaptive feature selection." *Expert Syst. Appl.* 38 (2011): 1393-1399.
- [4] R. Yan, X. Wan, J. Otterbacher, L. Kong, X. Li, and Y. Zhang. "Evolutionary timeline summarization: A balanced optimization framework via iterative substitution." *34th Int. ACM SIGIR Conf. Res.* 2011, 745-754.
- [5] J. Nichols, J. Mahmud, and C. Drews. "Summarizing sporting events using twitter." *ACM Int. Conf. Intell* (2012): 189-198.
- [6] Zhenhua Wang, Lidian Shou, Ke Chen, Gang Chen and Sharad Mehrotra. "on summarization and timeline generation forevolutionary tweet stream." *IEEE* 27 (2015): 1301-1315.
- [7] N. A. Diakopoulos and D. A. Shamma, "Characterizing debate performance via aggregated twitter sentiment," in *Proc. SIGCHI Conf. Human Factors Comput. Syst.*, 2010, pp. 1195-1198.
- [8] M. Dork, D. Gruen, C. Williamson, and S. Carpendale, "A visual backchannel for large-scale events," *IEEE Trans. Vis. Comput. Graph.*, vol. 16, no. 6, pp. 1129-1138, Nov. 2010.
- [9] R. Yan, X. Wan, J. Otterbacher, L. Kong, X. Li, and Y. Zhang, "Evolutionary timeline summarization: A balanced optimization framework via iterative substitution," in *Proc. 34th Int. ACM SIGIR Conf. Res. Develop. Inf. Retrieval*, 2011, pp. 745-754.
- [10] A. Marcus, M. S. Bernstein, O. Badar, D. R. Karger, S. Madden, and R. C. Miller, "Twitinfo: Aggregating and visualizing microblogs for event exploration," in *Proc. SIGCHI Conf. Human Factors Comput. Syst.*, 2011, pp. 227-236.
- [11] J. Nichols, J. Mahmud, and C. Drews, "Summarizing sporting events using twitter," in *Proc. ACM Int. Conf. Intell. User Interfaces*, 2012, pp. 189-198.
- [12] 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.
- [13] D. Chakrabarti and K. Punera, "Event summarization using tweets," in *Proc. 5th Int. Conf. Weblogs Social Media*, 2011, pp. 66-73.
- [14] M. Kubo, R. Sasano, H. Takamura, and M. Okumura, "Generating live sports updates from twitter by finding good reporters," in *Proc. IEEE Int. Joint Conf. Web Intell. Agent Technol.*, 2013, pp. 527-534.