

# Analysis of Unstructured Streaming Twitter Data

Dr. Sandeep M. Chaware<sup>1</sup>, Miss. Avanti Pandit<sup>2</sup>, Miss. Ashwini Pawar<sup>3</sup>, Miss Pratiksha Tambe<sup>4</sup>, Miss Nirmal Thakare<sup>5</sup>.

<sup>1</sup>Head of the Dept. <sup>2,3,4,5</sup>B. E Students,

<sup>1,2,3,4,5</sup>Department of Computer Science Engineering, Marathwada Mitra Mandal's College of Engineering, Karvenagar, Pune, India

**Abstract-** Microblogging offerings have revolutionized the way human beings trade facts. Confronted with the ever-growing numbers of tweets with multimedia contents and trending topics, it's far proper to offer visualized summarization to assist users to quickly hold close the essence of topics. While existing works normally attention on text-based strategies best, summarization of a couple of unstructured data (e.g., text and image) are scarcely explored. In this paper, proposes a tweet summarization framework to automatically generate visualized summaries for unstructured trending topics. Specifically, a model, termed multimodal-LDA (MMLDA), is proposed to find subtopics from microblogs (which are unstructured which may contain text, images, audios and videos) by means of exploring the correlations amongst different unstructured media kinds. Based on the records accomplished from MMLDA, a multimedia summarizer is designed to one by one pick out representative textual and visual samples and then form a complete visualized summary. The contribution work is, to extract microblogs which shows unstructured data as images, text, video posts. We evaluate the widespread experiments on an actual-international Twitter microblog dataset to illustrate the prevalence of our proposed technique towards the modern-day processes.

**Keywords-** Microblog, Unstructured, Summarization, Trending Topic, Twitter, Social Media, OCR, MMLDA

## I. INTRODUCTION

Users are allowed to share multimedia content on platforms, such as news, images and video links. With the wide availability of information sources, rapid information propagation and ease of use, microblogging has quickly become one of the most important media for sharing, distributing and consuming interesting contents, such as the trending topics. Currently, some microblogging platforms, such as Twitter, offer users the list of (manually created) hot trending topics, together with a set of related microblogs in each trend. Such service offers a potentially useful way to help users to conveniently gain a quick and concise impression of the current hot topics. In addition, users may obtain further understanding of the topics by browsing the related microblogs. However, due to the tremendous volume of tweets and the lack of effective summarization mechanism in

existing trending topic services, users are often confronted with incomplete, irrelevant and duplicate information, which makes it difficult for users to capture the essence of a topic. Therefore, it would be of great benefit if an effective mechanism can be provided to automatically mine and summarize subtopics (i.e., divisions of a main topic) from tweets related to a given topic.

## A. Motivation

1. The analysis method to automatically generate visualized summaries for trending topics and analyzing the summary.
2. Images can supplement the textual content with additional information, especially in the circumstance of tweets or posts, where the text lacks sufficient expressive power as aforementioned.
3. Multimedia contents can facilitate subtopic discovery.
4. Incorporating concrete multimedia exemplars into summarization can assist users to gain a more visualized understanding of interesting topics and/or subtopics.

## II. RELATED WORK

The author proposes a novel matrix factorization approach for extractive summarization, leveraging the success of collaborative filtering. First to consider representation learning of a joint embedding for text and images in timeline summarization [1]. A multimedia microblog summarization framework to automatically generate visualized summaries for trending topics. Specifically, a novel generative probabilistic model, termed multimodal-LDA (MMLDA) is proposed to discover subtopics from microblogs by exploring the correlations among different media types [2]. In this model they were organizing the messy microblogs into structured subtopics and also generating the high quality textual summary [2].

Novel idea of using the context sensitive document indexing is proposed to improve the sentence extraction-based document summarization task. In this paper, the team proposes a context sensitive document indexing model based on the Bernoulli model of randomness [3]. The team argue that for some highly structured and recurring events, such as sports, it is better to use more sophisticated techniques to summarize the relevant tweets. The problem of summarizing event-tweets and give a solution based on learning the

underlying hidden state representation of the event via Hidden Markov Models [4].

The interface for browsing Twitter streams name as “Eddi” which clusters tweets by topics trending within the user’s own feed. An algorithm for topic detection and a topic-oriented user interface for social information streams such as Twitter feed. BenchmarkTweet Topic against other topic detection approaches, and compares Eddi to a typical chronological interface for consuming Twitter feeds [5]. The author proposed methods to compute quality, diversity and coverage properties using multidimensional content and context data. The proposed metrics which will evaluate the photo summaries based on their representation of the larger corpus and the ability to satisfy user’s information needs [6].

There are some limitations of worked carried out earlier, which are:

- Only focus on summarizing synchronous multi-modal content.
- Computation is expensive because of algorithm time complexity.
- Real time summary generation is not done.
- When concept is small the previous strategies are not applicable.

#### A. OPEN ISSUES

Text summarization is performed for the purposes of saving users time by reducing the amount of content to read. However, text summarization has also been performed for purposes such as reducing the number of features required for classifying or clustering documents. Some microblogging platforms offer users the list of (manually created) hot trending topics, together with a set of related microblogs in each trend. Such service offers a potentially useful way to help users to conveniently gain a quick and concise impression of the current hot topics. In addition, users may obtain further understanding of the topics by browsing the related tweets. However, due to the tremendous volume of microblogs and the lack of effective summarization mechanism in existing trending topic services, users are often confronted with incomplete, irrelevant and duplicate information, which makes it difficult for users to capture the essence of a topic. Some of the issues are given below:

- The lack of effective summarization mechanism.
- Users are often faces with incomplete, irrelevant and duplicate information due to existing trending topics services.
- It makes difficult for users to capture the essence of a topic.

### III. PROPOSED SYSTEM OVERVIEW

Traditional documents that contain only textual objects, microblogs constitute of multiple media types, such as image

and text. In this paper proposes a novel framework to summarize and analyze multimedia microblogs for trending topics. Specifically, first proposes a novel generative probabilistic model, called multimodal-LDA (MMLDA), to partition the microblogs relevant to the same topic into different subtopics. MMLDA model is capable of not merely capturing the intrinsic correlation between visual and textual information of microblogs, but also estimating the general distribution as well as subtopic-specific distribution under a trending topic. For text summarization, specifies three criteria, namely coverage, significance and diversity to measure the summarization quality. For visual summarization, a two-step process is devised to automatically select the most representative images: 1) images within a subtopic are grouped by spectral clustering; and 2) images in each group are ranked by a manifold ranking algorithm and the top-ranked image is selected as representative. The Fig. 1 shows the architecture of unstructured twitter data summarization system. The processes of generating textual and visual summaries for each subtopic, by utilizing the reinforced textual/visual distributional information. Then, the textual and visual summaries are aggregated to form a comprehensive multimedia summary.

The contribution work is, to extract microblogs / tweets which show the video posts. These videos are separated in two parts: one is number of images and second is speech transcriptions. We evaluate the widespread experiments on an actual-international Twitter microblog dataset to illustrate the prevalence of our proposed technique towards the modern-day processes.

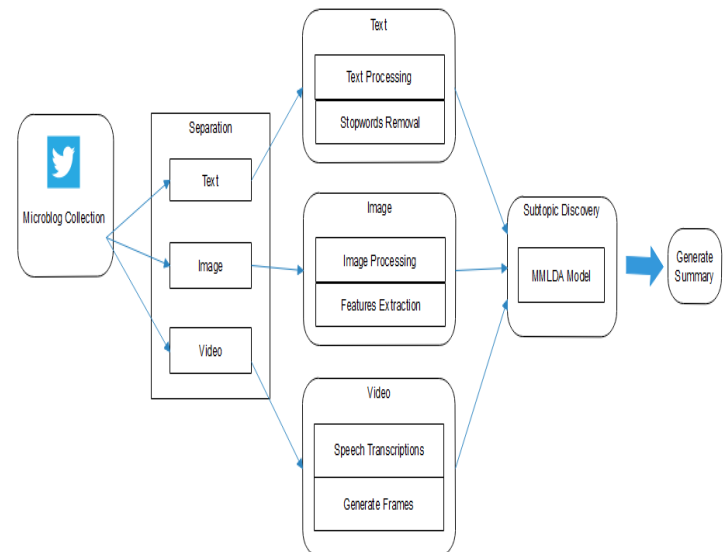


Fig.1: System Architecture

The fig.1 shows the components of system architecture. The components in the above fig are described below:

#### 1. Microblog Collection:

This module will extract the tweets from Twitter API for further processing. These tweets are unstructured, it may contain text, audio, video. We extract up to 1500 tweets for processing.

#### 2. Separation:

The data or tweets that are fetched from the Twitter API are separated at this stage. The separation is categorized as images, text and video, as per category processing module is apply on that data.

#### 3. Processing:

- a. Text Processing: Text is processed stopwords are removed in this, misspelled words are corrected.
- b. Image Processing: Image is processed for extracting the text from the images, feature extraction technique is used for extracting the text from the images.
- c. Video Processing: The videos are separated in two parts: one is number of images and second is speech transcriptions. Separated images again provide to image processing module to analyze text from image. Audio is provided to speech transcription module to generate text from it.

#### 4. Subtopic Discovery:

This phase uses the multimodal-LDA algorithm for summarization of the output of processed data which comes from the text processing, image processing and video processing.

#### 5. Summary:

Summary is generated and displayed.

#### Advantages

- It provides to automatically mine and summarize subtopics (i.e., divisions of a main topic) from microblogs related to a given topic.
- Microblogs comprise of multiple media types, such as image and text and video.
- Multimedia contents can facilitate subtopic discovery.
- Well organizing the unstructured microblogs into structured subtopics.
- Generating high quality textual summary at subtopic level.
- Selecting images relevant to subtopic that can best represent the textual contents.

#### IV. CONCLUSIONS & FUTURE WORK

In this paper, we proposed a multimedia microblog summarization method to automatically generate visualized summaries for trending topics. This system concludes that unstructured convert into structured data and stored. On the basis of that data results the summary will be generated and analysis is done.

Microblogs comprise of multiple media types, such as image and text and video. Specifically, a novel multimodal-LDA

(MMLDA) model was proposed to discover various subtopics as well as the subtopic content distribution from microblogs, which explores the correlation among different media types. Based on MMLDA, a summarizer is elaborated to generate both textual and visual summaries. The model will be able to organize the unstructured microblogs into structured subtopics. Generating high quality textual summary at subtopic level.

#### V. REFERENCES

- [1]. W. Y. Wang, Y. Mehdad, D. R. Radev, and A. Stent, "A low-rank approximation approach to learning joint embeddings of news stories and images for timeline summarization," in NAACL-HLT, 2016, pp. 58–68.
- [2]. Y. Yang, and T.-S. Chua, "Multimedia summarization for trending topics in microblogs," in CIKM. ACM, 2013, pp. 1807–1812.
- [3]. P. Goyal, L. Behera, and T. M. McGinnity, "A context-based word indexing model for document summarization," IEEE Transactions on Knowledge & Data Engineering, vol. 25, no. 8, pp. 1693–1705, 2013.
- [4]. D. Chakrabarti and K. Punera, "Event summarization using tweets," in Proc. 5th Int. AAAI Conf. Weblogs Social Media, 2011, pp. 66–73.
- [5]. M. S. Bernstein, B. Suh, L. Hong, J. Chen, S. Kairam, and E. H. Chi, "Eddi: Interactive topic-based browsing of social status streams," in Proc. 23rd Annu. ACM Symp. User Interface Softw. Technol., 2010, pp. 303–312.
- [6]. P. Sinha, S. Mehrotra, and R. Jain, "Summarization of personal photologs using multidimensional content and context," in Proc. 1st ACM Int. Conf. Multimedia Retrieval, 2011, p. 4.
- [7]. JingwenBian, YangYang, Hanwang Zhang, Tat-Seng Chua, "Multimedia Summarization for Social Events in Microblog Streams", IEEE transactions on multimedia, Vol. 17, No. 2, February 2015
- [8]. JingwenBian, YangYang, Tat-Seng Chua, "Multimedia Summarization for Trending topics in Microblogs", in CIKM '13 Proceedings of the 22nd ACM international conference on Information & Knowledge Management, Pages 1807-1812 San Francisco, California, USA — October 27 - November 01, 2013.
- [9]. H. Lin and J. Bilmes, "Multi-document summarization via budgeted maximization of submodular functions," in Proc. Human Lang. Technol.: Annu. Conf. North Amer. Chapter Assoc. Comput. Linguistics, 2010, pp. 912–920.