

Temporal TF-IDF: A High Performance Approach for Event Summarization in Twitter

Nasser Alsaedi, Pete Burnap and Omer Rana

Cardiff School of Computer Science & Informatics, Cardiff University

{AlsaediNM, P.Burnap, O.F.Rana}@cardiff.ac.uk

Abstract—In recent years, there has been increased interest in real-world event summarization using publicly accessible data made available through social networking services such as Twitter and Facebook. People use these outlets to communicate with others, express their opinion and commentate on a wide variety of real-world events. Due to the heterogeneity, the sheer volume of text and the fact that some messages are more informative than others, automatic summarization is a very challenging task. This paper presents three techniques for summarizing microblog documents by selecting the most representative posts for real-world events (clusters). In particular, we tackle the task of multilingual summarization in Twitter. We evaluate the generated summaries by comparing them to both human produced summaries and to the summarization results of similar leading summarization systems. Our results show that our proposed Temporal TF-IDF method outperforms all the other summarization systems for both the English and non-English corpora as they lead to informative summaries.

Keywords—microblogs summarization, Twitter, multilingual summarization, summary evaluation.

I. INTRODUCTION

Social networking services such as Twitter generate large volumes of content for many popular real-world events on a daily basis. Due to the number of posts published by hundreds of millions of users, digging through the noise and redundancy to extract and summarize the informative aspects of the content is a very challenging task. Moreover, the Twitter API only allows users to see the most recent posts on a topic in chronological order; it does not provide any relevancy based ordering of posts. This motivates the need for new automatic summarization systems that enable decision makers to be presented with informative summaries of user-generated content that support intelligence gathering and augment traditional situational awareness information sources. Such posts are likely to be multilingual in nature so the system also needs to be capable of handling this. Additionally these methods should be able to handle information flows in real-time - as events unfold.

Text summarization is the process of automatically creating a compressed version of a given text that provides useful information for the user [1]. Generally, there are two approaches in performing text summarization: extractive and abstractive. *Extractive* methods select a subset of words, phrases, or sentences from the original document to form a summary [2]. This contrasts with *abstractive* summarization, where the phrases in the text are rephrased may not appear in the input text [3].

Regarding summarizing Twitter streams, summarization (or tweet representation which is a form of extractive methods) can be viewed as a problem of automatically generating a condensed version of the most important content from one or more documents (tweets) [4]. Furthermore, the problem can also be defined as a ranking task where all tweets about a certain topic/event are ranked according to a weighting measure. Hence, the summarization is divided into two steps: (i) event detection and then (ii) tweet selection [5]–[7].

In this paper, we propose three techniques that focus on summarizing Twitter messages corresponding to events to improve event reasoning, visualization, and analytics. Our methods are language independent, satisfy real-time requirement and suitable with the huge quantity of the data. We use frequency-based method, voting approach and centrality-based technique to select messages that represent an event with high quality, strong relevance and are useful to people looking for information about an event. We evaluate our proposed techniques using a real-world dataset of Twitter messages according to ROUGE-1 [8] as well as metrics (quality, relevance and usefulness) [5]. Our methods are tested using English, Arabic and Japanese languages to test its applicability across multiple languages. We also compare their performance with number of recent and leading summarization systems including Becker et al. (centroid method) [5], Zubiaga et al. (sub-event detection then tweet selection method) [7], Xu et al. (graph-based approach) [9] and Hybrid TF-IDF (term frequency summarization approach) [4]. Finally we combined our three proposed summarization techniques in a visualization tool for event detection and decision making process.

II. RELATED WORK

There are numerous approaches for automatic summarization and detection of topics from social media. Many of them are inspired by the previous work on automatic text summarization. [10] have reviewed an extensive survey of text summarization techniques. In this paper we only focus on microblog event summarization systems.

The centroid-based method is one of the most popular extractive summarization methods. MEAD [11] is an implementation of the centroid-based method that scores sentences based on sentence-level and inter-sentence features, including cluster centroids, position, TF-IDF [Term Frequency - Inverse Document Frequency], etc. Moreover, MEAD is a

flexible platform for multi-document multi-lingual summarization which is publically available. Similarly, in [5] Becker et al. presented three centrality-based approaches (LexRank, Degree and Centroid) to select the high quality messages from clusters. Authors found that the Centroid approach, which computes the cosine similarity of the TF-IDF representation of each message to its associated event cluster centroid, outperformed other approaches based on three metrics: *quality*, *relevance*, and *usefulness*.

In feature-based approaches, a variety of statistical and linguistic features have been extensively investigated, for example, Sharifi et al. [12] proposed a phrase reinforcement (PR) algorithm to summarize the Twitter topic in one sentence. [13] extended this idea and generated journalistic summary for events in world cup games by employing phrase graph algorithm only on the longest sentence in each tweet. More fine-grained summarization was proposed by considering sub-events detection and combining the summaries extracted from each sub-topic (tweet selection, tweet ranking) [6], [7], [14]–[16]. For instance, Zubiaga et al. [7] explored the real-time summarization of scheduled events using a two-step system: (i) sub-event detection and (ii) tweet selection. The first step is based on peaks detection (reflected as peaks in the histogram of tweeting rates) with an enhancement of two ideas; the sudden increase in the tweeting rate and the outlier detection. The tweet selection step selects a representative tweet after ranking all tweets that were sent during the sub-event. They use the Kullback-Leibler divergence (KLD) weighting scheme for the tweet ranking.

Another approach is the graph-based LexRank which was introduced by [3]. The LexRank algorithm computes the relative importance of sentences in a document (or a set of documents). Then it creates an adjacency matrix among the textual units and finally computes the stationary distribution considering it to be a Markov chain. In their approach, they showed that the similarity graph of sentences provides a better view of important sentences compared to the centroid approach. The TextRank algorithm [17] is another graph-based approach that implements two unsupervised approaches for keyword and sentence extraction in order to find the most highly ranked sentences in a document using the PageRank algorithm [18]. Recently, Xu et al. in [9] extended the PageRank ranking algorithm and investigate a graph-based approach which leverages named entities, event phrases and their connections across tweets to create summaries of variable length for different topics. Moreover, [19] proposed a graph-based abstractive summarization scheme where bigrams extracted from the tweets are considered as the graph-nodes.

SumBasic [20] is a simple, yet high-performing summarization system based on term frequency. Authors empirically showed that words that occur more frequently across documents are more likely to appear in human generated multidocument summaries. Most recently, [4] developed a new method called “Hybrid TF-IDF”, which ranks tweet sentences using the TF-IDF scheme and produces better results than some of the leading traditional summarization approaches in

the context of microblogs summarization.

Several previous works have leveraged the importance of monitoring the event evolution. For example, [21] derived three features from timelines and used them in supervised learning to enhance multi-document summarization (MDS). [22] propose a language model with dynamic pseudo relevance feedback (PRF) to obtain relevant tweets, and then generate storylines via graph optimization. In [23], the authors proposed two topic models (Decay Topic Model (DTM) and Gaussian Decay Topic Model (GDTM)) that take advantage of temporal correlation in the data to extract relevant tweets for summarization. Other researchers have proposed models for the purpose of summarizing micro-blog events in Twitter including the use of Non-negative Matrix Factorization (NMF) [24], a structured retrieval approach [25], Structured Probabilistic Latent Semantic Analysis (PLSA) [26], and many more [23], [27], [28].

III. PROPOSED SUMMARIZATION APPROACHES

We propose three methods for summarizing a set of Twitter posts; Temporal TF-IDF, Retweet Voting Approach and Temporal Centroid Representation method. For all proposed methods, we use a one-hour time window based on the best temporal settings as described in [29]. The temporal TF-IDF is based on extracting the most highly weighted terms as determined by the TF-IDF weighting for two successive time frames. The voting method considers the highest number of retweets a post received in the time window as the criterion for finding the most representative post in a single time window. This method reflects users’ choices as they decide which message is the most ‘valuable’ by propagating it. The temporal centroid method selects posts that correspond to each cluster centroid as the summary of that cluster with respect to the time dimension. Next, we describe these methods and provide an analysis of the results.

A. Temporal TF-IDF

The algorithm is inspired by the fact that users tend to use similar words when describing a particular event as well as observations obtained from [30]:

- 1) High frequency words like stop-words occur in approximately the same percentage of documents no matter whether the document set is small or large and similarly, low frequency words like “murder” occur very rarely across small and large datasets.
- 2) The document frequency distribution of one corpus can be used to approximate another.

We propose a novel temporal Term Frequency - Inverse Document Frequency (TF-IDF) that generates a summary of top terms without the need of prior knowledge of the entire dataset, unlike the existing TF-IDF approach [31] and its variants. Temporal TF-IDF is based on the assumption that words which occur more frequently across documents over a particular interval (timeframe) have a higher probability of being selected for human created multi-document summaries than words that occur less frequently [20].

Typically, the TF-IDF approach requires knowing the frequency of a term in a document (TF) as well as the number of documents in which a term occurred at least once (DF). The need for *a priori* knowledge of the entire data set introduces significant challenge of using this approach where continuous data streams must be summarized in real-time as an event unfolds. In addition, the adopted scheme must be flexible to update frequently (every minute, 10 mins, hourly, 3 hours - depending on the time-frame size). Hence, the iterative calculation of term weights should be taken into consideration.

To overcome these limitations we introduce the temporal TF-IDF where we consider a set of posts in a cluster to be represented as a document. The total number of clusters equals the total number of documents which is a subset of the entire dataset or corpus. This reduces the overall computational complexity and overcomes the limitations of the TF-IDF based approaches in which the document set to be clustered must be known in advance. After the first cluster timeframe, we use clusters from the previous timeframe with the documents in the recent one to add more relevance and usefulness to our results such as emerging keyword. Consequently, we use the document frequency distribution of two timeframes instead of one, taking into account the changing event dynamic and narrative. We define the TF-IDF weighting scheme of a new document d for a collection C (from two clusters) as:

$$w_{ji} = \frac{1}{\text{norm}(d_i)} f_{ji} \times \log\left(1 + \frac{N}{N_j}\right)$$

where f_{ji} is the term frequency of word in document d_i and N_j is document frequency of word in a collection and N is the total number of documents in the collection. In order to avoid the bias caused by different document lengths, the length of each document vector is normalized so that it is of unit length $\text{norm}(d_i)$. This summarizer selects the most weighted post as summary as determined by the Temporal TF-IDF weighting.

B. Retweet Voting Approach

Many studies have illustrated the power of retweeting for many tasks such as predicting most influential users [32], identifying most knowledgeable posts [33], ranking and measuring information propagation [34] and analyzing network structure [34], [35]. Voting algorithms have been successfully implemented in many data mining applications [36]. Here we implement the highest number of retweets as a measure of representation task through voting algorithm. Voting algorithms have been used in many applications where in the context of social media may be considered taking into account the following features:

- The average length of a post.
- The total frequency of features in a post.
- Number of times of retweets, favorites and mentions.
- The inclusion of multimedia files such as photos, videos.

Using the retweet count (the number of times a tweet in a cluster has been retweeted) as ranking method in cluster has several benefits; first it represents the influence of a tweet beyond one-to-one interaction [32]. Second, retweeting serves as

a powerful tool to reinforce a message when not only one but a group of users repeat the same message [33]. Third, number of retweets is an indication of popularity [32] so in a way we are summarizing the cluster by the highest degree of agreement from users themselves. In addition we can generalize this method to be applied to other social networking sites such as Facebook (number of Share), Instagram (number of likes), Pinterest (number of Repins), etc. in one time-frame. We can also extend this approach to rank events/clusters by calculating the total number of retweets per cluster. However, using this method of representation suffers from many drawbacks:

- 1) The content of tweet is not always taken into consideration as many users retweet without even reading. For instance, most celebrities have a high number of retweets based on their popularity.
- 2) A tweet with high number of retweets might repeat over time as it receives the highest attention and Retweet Count generally increases with time. Thus, Retweet Score is not a comprehensive measure.

Many techniques including classification and clustering approaches have successfully been proposed and shown to distinguish between messages about real-world events and non-event messages [36], hence most of the celebrities' messages are removed unless they are related. To overcome the second drawback, we introduce a normalization factor where we calculate the Change of Retweet Score with time instead of the Retweet Score.

Retweet Score (rt) is defined as the ratio of the number of "retweets that a tweet gets (u_i) to the total number of retweets (u_{all}) of all posts in the target cluster. It is defined as,

$$rt = \frac{|\text{retweet}(u_i)|}{|\text{retweet}(u_{all})|}$$

Retweet Score Change is defined as the number of times a tweet has been retweeted in current time-frame (rt_{cur}) and is calculated by subtracting number of retweets count from previous time-frame (rt_{pr}) of that post.

$$rt \text{ change} = rt_{cur} - rt_{pr}$$

C. Temporal Centroid Method

The centroid approach takes into consideration a centrality measure of a tweet with respect to the overall topic of the cluster [2], [5]. It computes the cosine similarity of the TF-IDF representation of each message to its associated event cluster centroid, where each cluster term is associated with its average weight across all cluster messages. Then it selects the messages with the highest similarity value because they represent the average weight of all terms in clusters. The main idea behind this method (as it is based on frequency across all messages) is to identify posts that have high quality and most relevant to an entire cluster. The difference between our proposed centroid method and other centroid methods is that we include the time dimension. We select a post which has been a centroid for the longest time on average over a time-window rather than just taking the final centroid at

the end of that time-window. We believe that studying the temporal aspects of posts reveal additional information about their quality, relevance, and usefulness.

IV. EMPIRICAL EVALUATION

A. Datasets and Setup

Datasets: We used the Twitter Streaming API to collect around 2.7 Million tweets (2677937) posted from 26 November 2014 to 8 December 2014. This dataset was collected as part of our previous work on Twitter event identification [29].

Annotations: We selected top 10 event clusters per day, with an average of 320 posts per cluster. For each event cluster we selected the top-5 posts according to our proposed approaches (Temporal TF-IDF, Retweet voting, and Temporal centroid methods). We used three human annotators to label each post according to three desired goals as reported by [5]:

- 1) **Quality:** refers to the textual quality of the messages, which reflects how well they can be understood by a human. High-quality messages contain crisp, clear, and effective text that is easy to understand.
- 2) **Relevance:** how well a Twitter message reflects information related to its associated event. Highly relevant messages clearly describe their associated event.
- 3) **Usefulness:** the potential value of a post for someone who is interested in learning details about an event. Useful messages should provide some insight about the event, beyond simply stating that the event has occurred.

The annotators labeled each message on a scale of 1-4 for each attribute, where a score of 4 signifies high quality, strong relevance, and clear usefulness, and a score of 1 signifies low quality, no relevance, and no usefulness. A set of instructions and examples were given to annotators in order to perform the task as well as the assessments were done without reference to any model summaries. We have used crowdsourcing to annotate posts but this time we used the CrowdFlower system (<http://www.crowdflower.com>). Agreement between annotators was substantial to high, with kappa coefficient values = 0.92; 0.89; 0.61 for quality, relevance, and usefulness, respectively. After the annotators got familiar with the topics and the summarization task, each annotator was asked to summarize each cluster in order to generate the **gold standard** summaries.

Evaluation Methods: The similarity metric we use for evaluation and comparison between system summaries is the ROUGE metric proposed by [8]. The ROUGE metric counts the total number of matching n-grams (excluding stop-words) between the true summary and the summary generated from model. It has been shown that ROUGE scores correlate well with human judgments [8]. In this work, we use **ROUGE-1** scores as fitness function for measuring summarization quality because it showed the largest variation across other methods ROUGE-2, 3, 4, *L*, *SU* and *W* metrics [14], [15]. We have also evaluated the summarization techniques using metrics; quality, relevance and usefulness [5], [9].

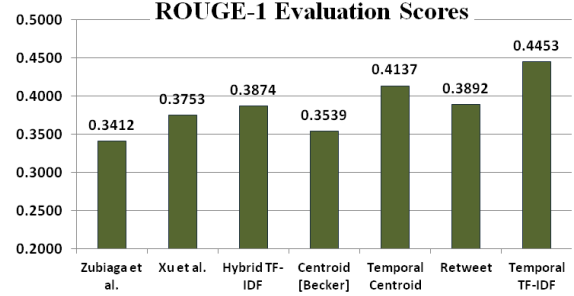


Fig. 1. Results of our proposed approaches against other summarization techniques.

B. Experimental Results

We conduct several experiments to evaluate different aspects of our summarization techniques. In the first experiment, we compare our proposed approaches to other recent and leading summarizers, including Becker et al. [5], Zubiaga et al. [7], Xu et al. [9] and Hybrid TF-IDF Summarizer [4]. We selected these baselines because they have been shown successful and effective for summarizing tweets as well as they represent different methods as described in section 2. We evaluate the different summarizers using the automatic ROUGE-1 evaluation. The values of the ROUGE-1 scores are presented in Figure 1.

Our approaches achieve good performance compared to other summarization methods. The Temporal TF-IDF adds more knowledge when determining both components (TF) and (IDF) for two timeframes. Our Centroid algorithm has achieved superior performance to the other approaches due to its inherent assumption that each cluster revolves around one central topic. In addition, the Retweet approach produces more satisfactory results than the baseline approaches summaries. Note that the ROUGE scores are based solely on the n-gram overlap between the system and reference summaries, which may not be the most appropriate measure for evaluating the event summaries. Hence further experiments are needed to investigate the proposed methods using more sophisticated evaluation measures.

The second experiment compares between the competing approaches according to user-perceived quality, relevance, and usefulness. Figure 2 summarizes the average performance of these approaches across all 50 test events.

All three of our proposed approaches receive high scores for quality where the Temporal TF-IDF produces the highest score. In other words, our approaches are able to select clear, informative summary according to human judgements. The Temporal TF-IDF technique also receives a high score for usefulness, indicating that its selected messages are useful with respect to the associated events. The Temporal TF-IDF takes in consideration two timeframes hence more details about an event are provided compared to other methods. The Temporal centroid and the Temporal TF-IDF, on average, select messages that are either somewhat relevant or highly

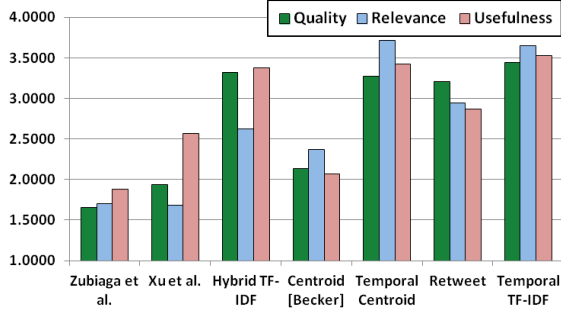


Fig. 2. Comparison of content selection techniques.

relevant which indicate that the Retweet voting approach is affected negatively toward the most influential users.

One noted problem with the Xu et al. method is that this graph-based summarization is sensitive to document length because its similarity estimation mainly depends on commonly used words, not words with high IDF scores, especially in short tweets, which increases the chances of selecting pointless information [14]. We also find that graph models tend to assign high salience scores to long tweets containing several #hashtags. Some of them are pointless tweets (low quality) and irrelevant (minimum relevance) to the topic.

Overall, it seems from the first two experiments that the simple frequency-based summarizers, namely Temporal TF-IDF and Hybrid TF-IDF, perform better than summarizers that incorporated more information or more complexity such as graph-based methods or centroid-based approaches. This possibly has much to do with the special nature of Twitter documents in which documents often have very little structure and have so few words that forming relationships between pairs of documents is not particularly helpful. Therefore, more complex relational models will probably not capture more topical information than frequency models and the added complexity of interrelationships did not help in summarizing Twitter posts. Moreover, Temporal TF-IDF outperforms Hybrid TF-IDF because our Temporal TF-IDF is more sensitive to time which can be clearly seen in both measures; quality and relevance.

In comparison to other methods, Zubiaga et al. and similar systems like [6], [14] are limited to scheduled events such as soccer games. They require the starting time in order for the system to start looking for new sub-events. The sub-event detection step (based on peaks detection) fails to detect important events/topics which reduces the chances of selecting valuable information (high quality, more relevant and very useful). This may explain why the Zubiaga et al. approach performance is lower than the results reported in this paper. Similarly, the centroid-based summarizers such as Becker et al. attempt to reduce redundancy but did so by clustering the documents first and then summarizing based on these clusters. However, their clustering approach did not seem to increase performance particularly in small-sized clusters, which have very little number of tweets per cluster. This can be

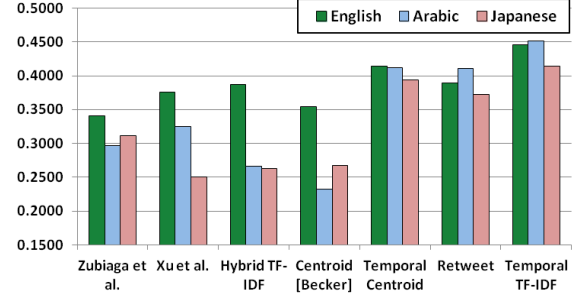


Fig. 3. ROUGE-1 results of various summarization techniques for different languages

clearly seen by the low quality and the insignificant usefulness measures in Figure 2.

For the third experiment, we generate subsets of our dataset to evaluate and compare the performance of the summarization systems using different languages. We randomly created three smaller subsets of English, Arabic and Japanese posts (number of posts are: 200, 500 and 40, respectively). We intentionally chose English, Arabic and Japanese because they belong to distinct language families (Indo-European, Semitic and Altaic languages, respectively). The results of the average ROUGE-1 values obtained for English, Arabic and Japanese corpora are shown in Figure 3.

The results in Figure 3 confirm findings from the first two experiments and in fact are consistent across all languages considered. The results in Figure 3 show that our proposed approaches outperform other summarizers for morphologically-rich languages such as Arabic and Japanese. For Japanese and Arabic, the performance of our Temporal TF-IDF and Temporal centroid supports that these methods handles well the variety of morphological phenomena present in these languages. The simple Retweet voting method achieves good results across all languages which support our claim that users' choices as they decide which message is the most representative of an event (cluster) and with no additional knowledge of that language can be used for the micro-blog summarization task.

C. Visualization Tool (Case study)

Visualizing complex relationships around real-time events in social streams is important to gain insight and use this for decision making. Hence, we combine the three proposed methods and implement them in a visualization tool for event detection and summarization. Our goal with the visualization tool is to facilitate discovery and increase interpretability of Twitter summaries for decision makers. We aim to support the exploratory and the event identification in a particular location while also giving the administrators easier access to search and the ability to efficiently explore online communities. We implement this interface in R tool (<http://www.r-project.org>) which is a free software environment for statistical computing and graphics. In particular, we use the package Gephi [37],

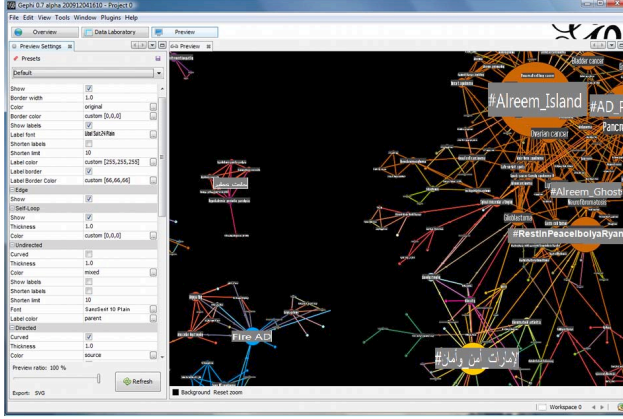


Fig. 4. Results of combining summarization approaches after classification-clustering framework

an open source graph visualization manipulation software (available at <http://gephi.org/>).

We use tweets for one day (3/12/2014) from our dataset to visualize and identify main events for that day as shown in Figure 4. Figure 4 features a number of nodes, each representing a particular tweet. The colour of each node represents the cluster that a tweet belongs to (Number of colours = Number of clusters (events)). The node size is dependent on two attributes: the TF-IDF value and the Change of Retweet Score the larger the node, the higher the retweet count for that post. Lines between nodes specify communication and relation between exchanged messages while also determine the centrality measures between messages and centroids.

The visualization tool is also able to visualize event-related updates, giving a comprehensive view of events throughout a predefined time period and the interpretation of these events. It supports the term-centric, temporal analytics of event-related information in Twitter. For creating the timeline we used the Annotated Time Line tool available in the Google Chart Tools (<https://developers.google.com/chart/>) as presented in figure 5.

In Figure 5, we present an example of the timeline of the number of tweets per hour (from the second dataset). Each of the peaks might be a candidate for an event, but because we employ a classification-clustering framework, only event-related updates are detected; some of the peaks are actually related to events in general and others indicate *disruptive events* events that threaten social safety and security, or could cause disruption to social order. Identifying disruptive events from social media stream is a useful source of information for improving situational awareness and decision support. The detected events and disruptive events are marked on the timeline, and are accompanied by a tag cloud description from different summarization systems.

V. CONCLUSION

The rate of information growth due to the social media content and the real-time requirement of many tools have

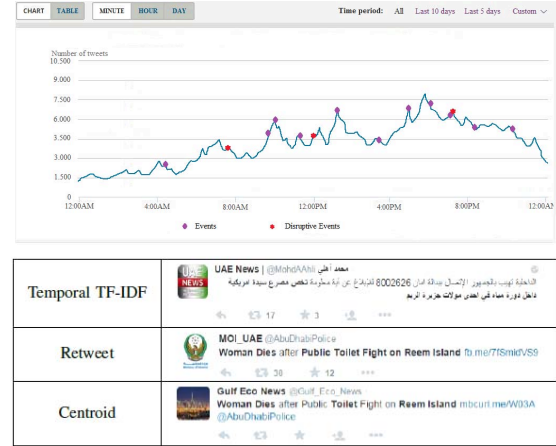


Fig. 5. The timeline of identified events and disruptive events with examples of summaries using different summarization techniques

called for a need to develop efficient summarization techniques. Here we implemented three summarization techniques; the Temporal TF-IDF, the Retweet voting approach and the Temporal centroid method. Based on results reported in this paper, the temporal frequency based method achieved the best results both in ROUGE scores and in human evaluation scores. The centroid representation also reflects the topic/event; hence the centroid representation performed well. Not far from them, user's choice (the retweet voting algorithm) achieved good results too which makes it among the best techniques for summarizing Twitter topics. Our evaluation also shows that our proposed methods perform well across a variety of language families, and we present here results that improve on current state-of-the-art for several noisy real-world datasets including a multilingual corpus. Finally we combined the three summarization techniques in a visualization to generate a meaningful real-time updates in order to facilitate the understanding and the exploration for the end-users and decision makers of ongoing events.

There are many interesting directions for future work. One of the main directions is to extend our investigation of the multilingual summarization and conduct more experiments on larger data sets as well as to add a range of other languages such as German, French, and Russian. We will study the effect of these different languages in greater detail in the future. Another direction is to produce multi-sentence or multi-post summaries or even to go further and form a coherent multi-sentence summary.

REFERENCES

- [1] D. Das and A. F. Martins, "A survey on automatic text summarization single-document summarization." *Literature Survey for the Language and Statistics II course at CMU, Pittsburg (2007)*, pp. 1–31, 2007.
- [2] N. Alsaedi, P. Burnap, and O. Rana, "Automatic summarization of real world events using twitter," *Proceedings of the 10th International AAAI Conference on Weblogs and Social Media (ICWSM'16)*, 2016.
- [3] G. Erkan and D. R. Radev, "Lexrank: Graph-based lexical centrality as salience in text summarization," *Journal of Artificial Intelligence Research*, vol. 22, no. 1, pp. 457–479, 2004.

- [4] D. Inouye and J. Kalita, "Comparing twitter summarization algorithms for multiple post summaries," *Proceedings of the 3rd IEEE International Conference on Social Computing (Socialcom'11)*, 2011.
- [5] H. Becker, M. Naaman, and L. Gravano, "Selecting quality twitter content for events," *Proceedings of the 5th International AAAI Conference on Weblogs and Social Media (ICWSM'11)*, 2011.
- [6] C. Shen, F. Liu, F. Weng, and T. Li, "A participant-based approach for event summarization using twitter streams," in *Human Language Technologies: Conference of the North American Chapter of the Association of Computational Linguistics, ACL'13*, 2013, pp. 1152–1162.
- [7] A. Zubiaga, D. Spina, E. Amigó, and J. Gonzalo, "Towards real-time summarization of scheduled events from twitter streams," in *Proceedings of the 23rd ACM Conference on Hypertext and Social Media*, ser. HT '12, 2012, pp. 319–320.
- [8] C.-Y. Lin and E. Hovy, "Automatic evaluation of summaries using n-gram co-occurrence statistics," in *Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology*, 2003, pp. 71–78.
- [9] W. Xu, R. Grishman, A. Meyers, and A. Ritter, "A preliminary study of tweet summarization using information extraction," in *Workshop on Language in Social Media (LASM 2013)*, Conference of the Association of Computational Linguistics, ACL'13, 2013, pp. 20–29.
- [10] A. Nenkova and K. McKeown, "A survey of text summarization techniques," in *Mining Text Data*. Springer, 2012.
- [11] D. R. Radev, S. Blair-Goldensohn, and Z. Zhang, "Experiments in single and multidocument summarization using mead," *First Document Understanding Conference*, 2001.
- [12] B. Sharifi, M.-A. Hutton, and J. Kalita, "Summarizing microblogs automatically," in *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, ser. HLT '10, 2010, pp. 685–688.
- [13] J. Nichols, J. Mahmud, and C. Drews, "Summarizing sporting events using twitter," in *Proceedings of the 2012 ACM International Conference on Intelligent User Interfaces*, ser. IUI '12, 2012, pp. 189–198.
- [14] D. Yajuan, C. Zhumin, W. Furu, Z. Ming, and H. Y. Shum, "Twitter topic summarization by ranking tweets using social influence and content quality," in *Proceedings of the 24th International Conference on Computational Linguistics (COLING' 12)*, 2012, pp. 763–780.
- [15] D. Chakrabarti and K. Punera, "Event summarization using tweets," in *Proceedings of the 5th International AAAI Conference on Weblogs and Social Media (ICWSM'11)*, 2011.
- [16] Z. Wang, L. Shou, K. Chen, G. Chen, and S. Mehrotra, "On summarization and timeline generation for evolutionary tweet streams," *IEEE Trans. Knowl. Data Eng.*, vol. 27, no. 5, pp. 1301–1315, 2015.
- [17] R. Mihalcea and P. Tarau, "TextRank: Bringing order into texts," *Proceedings of the 2004 conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 404–411, 2004.
- [18] S. Brin and L. Page, "The anatomy of a large-scale hypertextual web search engine," *Computer networks and ISDN systems*, vol. 30, no. 1-7, pp. 107–117, 1998.
- [19] A. Olariu, "Efficient online summarization of microblogging streams," in *Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics, EACL'14*, 2014, pp. 236–240.
- [20] L. Vanderwende, H. Suzukia, C. Brocketta, and A. Nenkova, "Beyond subbasic: Task-focused summarization with sentence simplification and lexical expansion," *Information Processing and Management*, vol. 43, no. 6, pp. 1606–1618, 2007.
- [21] J. Ng, Y. Chen, M. Kan, and Z. Li, "Exploiting timelines to enhance multi-document summarization," in *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics, ACL'14*, 2014, pp. 923–933.
- [22] C. Lin, C. Lin, J. Li, D. Wang, Y. Chen, and T. Li, "Generating event storylines from microblogs," in *Proceedings of the 21st ACM International Conference on Information and Knowledge Management*, ser. CIKM '12, 2012, pp. 175–184.
- [23] F. C. T. Chua and S. Asur, "Automatic summarization of events from social media," in *Proceedings of the 7th International AAAI Conference on Weblogs and Social Media (ICWSM'13)*, 2013.
- [24] X. Yang, A. Ghoting, Y. Ruan, and S. Parthasarathy, "A framework for summarizing and analyzing twitter feeds," in *Proceedings of the 18th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, ser. KDD '12, 2012.
- [25] D. Metzler, C. Cai, and E. Hovy, "Structured event retrieval over microblog archives," in *Proceedings of the 2012 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, ser. NAACL HLT '12, 2012.
- [26] Y. Lu, C. Zhai, and N. Sundaresan, "Rated aspect summarization of short comments," in *Proceedings of the 18th International Conference on World Wide Web*, ser. WWW '09, 2009, pp. 131–140.
- [27] S. Harabagiu and A. Hickl, "Relevance modeling for microblog summarization," *Proceedings of the 5th International AAAI Conference on Weblogs and Social Media (ICWSM'11)*, 2011.
- [28] B. O'Connor, M. Krieger, and D. Ahn, "Tweetmotif: Exploratory search and topic summarization for twitter," *Proceedings of the 4th International AAAI Conference on Weblogs and Social Media (ICWSM'10)*, 2010.
- [29] N. Alsaedi, P. Burnap, and O. Rana, "Identifying disruptive events from social media to enhance situational awareness," in *Proceedings of the 2015 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM 2015)*, 2015.
- [30] J. W. Reed, Y. Jiao, T. E. Potok, B. A. Klump, M. T. Elmore, and A. R. Hurson, "Tf-icf: A new term weighting scheme for clustering dynamic data streams," *Proceedings of the 5th International Conference on Machine Learning and Applications (ICMLA06)*, pp. 258–263, 2006.
- [31] G. Salton and C. Buckley, "Term-weighting approaches in automatic text retrieval," *Information Processing and Management*, vol. 24, no. 5, pp. 513–523, 1988.
- [32] M. Cha, H. Haddadi, F. Benevenuto, and K. P. Gummadi, "Measuring user influence in twitter: The million follower fallacy," *Proceedings of the 4th International AAAI Conference on Weblogs and Social Media (ICWSM'10)*, 2010.
- [33] S. Petrović, M. Osborne, and V. Lavrenko, "Rt to win! predicting message propagation in twitter," *Proceedings of the 5th International AAAI Conference on Weblogs and Social Media (ICWSM'11)*, 2011.
- [34] J. Cheng, L. A. Adamic, P. A. Dow, J. Kleinberg, and J. Leskovec, "Can cascades be predicted?" *Proceedings of the 23th International Conference on World Wide Web Conference Committee (IW3C2)*, pp. 925–936, 2014.
- [35] H. Kwak, C. Lee, H. Park, and S. Moon, "What is twitter , a social network or a news media?" *Proceedings of the 19th International Conference on World Wide Web Conference Committee (IW3C2)*, pp. 591–600, 2010.
- [36] N. Alsaedi, P. Burnap, and O. Rana, "A combined classification-clustering framework for identifying disruptive events," in *Proceedings of 7th ASE International Conference on Social Computing (Social-Com'14)*, 2014.
- [37] G. G. Vega, orge Fabrega, and J. Kunst, "rgexf: An r package to build gexf graph files," in *The Comprehensive R Archive Network*, 2012.