# Aggregating and Analyzing Articles and Comments on a News Website

David Ramamonjisoa[✉]

Faculty of Software and Information Science, Iwate Prefectural University,
IPU, Takizawa, Sugo 152-52, Japan
`david@iwate-pu.ac.jp`

**Abstract.** In the top news stories, the commenting activity is rising and falling until it stops. In some ongoing news stories such as disasters like the disappearance of flight MH370, global warming or climate change, political turmoil or economic crisis, this commenting activity cycle can repeat and last many years. To our knowledge, a study and analysis of those data does not exist up to now. There is a need to separate facts, opinions and junk within those comments data. In this paper, we present our framework for supporting readers in analyzing and visualizing facts, opinions and topics in the comments and its extension with comments aggregation and summarization for comments within several news articles for the same event. We added a time-series analysis and comments features such as surprising comments and a preferential threads attachment model.

**Keywords:** Articles and comments analysis · Comments visualization · Network analysis · Text mining

## 1 Introduction

User comments are a kind of user-generated content. Their purpose is to collect user feedback, but they have also been used to form a community discussing about any piece of information on the internet (news article, video, live talk show, music, picture, and so on). The commenting tool becomes social gathering software where commenters share their opinions, criticism, or extraneous information. Recent research consists of assisting the end-user to reading comments by providing succinct summary or useful comments according to some algorithms [1–3].

In the top news stories, the articles and commenting activity are rising and falling until they stop. In some unresolved events in the news such as disasters like the disappearance of flight MH370, global warming or climate change, political turmoil or economic crisis, this commenting activity cycle can repeat and last many years. To our knowledge, a study and analysis of those data has not been done. There is a need for an overview of the articles on the events and to separate facts, opinions and junk in the comments data. However, some studies have investigated the challenges in summarizing jointly the articles and comments covering the event, and generating socially-informed timelines [4]. The timelines are built by maximizing topic cohesion between the article and the comment summaries. The maximization uses an optimization algorithm that

allows the generation of a high quality article summary and comment summary via mutual reinforcement. A similar study on this method is described in the dual wing factor graph model between Web documents and their associated social contexts to generate summaries [5].

In this paper, we present our framework to support readers in analyzing and visualizing facts, topics, hidden structure in the comments and its extension with comment aggregation and summarization of comments within several news articles on the same event. A time-series analysis and comments specific features such as surprising comments and a preferential threads attachment model are also introduced. The results obtained from the experiments showed that the lift measure used for the surprising event feature combined with clustering and interesting features visualization can summarize and synthetize the event.

Firstly, we describe the comments modeling and system framework. Then, we present the features used for analysis such as clustering, classifying, summarizing and finding surprising comments. Thirdly, we detail the analysis models. Finally, we demonstrate the use of the models with a test data and discuss the results.

## 2 Aggregating News Articles and Comments

An event in the news (top story) can generate several articles and on the order of tens of thousands of comments on a single news website. The goal of this process is to gather all articles and comments from the event published in a single news website.

Let $E_k$ be some event, $A_k$ the associated articles and $T_k$ the topic category of the event. Then we have $E_k = \{T_k, A_k, C_k\}$.

$$T_k = \text{TopStories}|\text{World}|\text{Local}|\text{Entertainment}|\text{SciTech}$$
$$|\text{Business}|\text{Politics}|\text{Sports}|\text{Health}|\text{Products}$$

$A_k = \{a_{1k}, a_{2k}, \dots, a_{mk}\}$ where $a_{ik}$ is any article concerning the event

All comments are noted $C_k$ (as the comment corpus for the event). The model of the comment collection is described below:

$$C_k = \{c_i\}$$

where

$$c_i = (commentID_i, time_i, author_i, title_i, content_i, rating_i) \tag{1}$$

$$content_i = \begin{cases} \{s_{1i}, \dots, s_{ji}, \dots, s_{li}\}| \\ < reply-to > \\ < quotation > \end{cases} j \in [1..l] \tag{2}$$

$s_{ji}$ is the $j$-th sentence of the comment $c_i$ as a sequence of words $(w_1, .., w_k)$.

From the Eqs. 1 and 2, we can build a data matrix, a network graph or time series data from the comment corpus. We can also extend the equations to model comment properties such as the lift measure or topic model as described in this paper.

The aggregating module is part of the comments analysis framework depicted in Fig. 1. Comments are pre-processed by selecting and extracting relevant features. They can be summarized, compared or statistically examined. Pre-processed data are then used to discover some interesting patterns or regularities and then to visualize those findings for the end user.
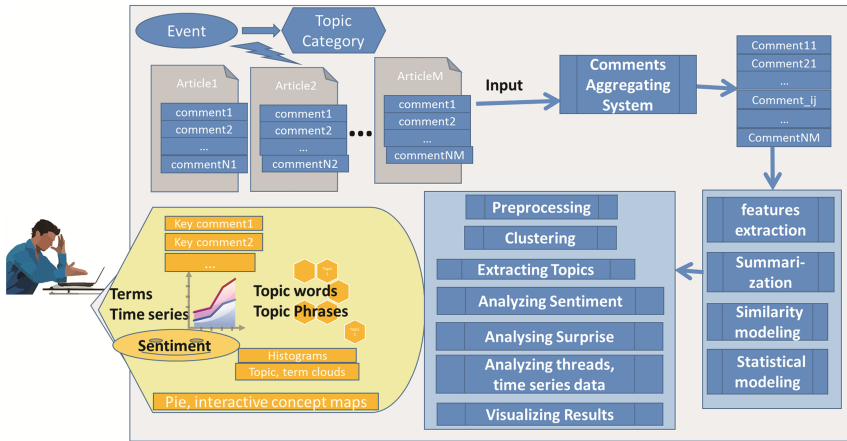


**Fig. 1.** Framework for the articles and comments aggregation and analysis.

## 3   Features on Articles and Comments

The feature selection and extraction module uses patterns of statistic relevance, sentence writing style relevance, natural language processing (NLP) techniques relevance and topic/user/sentiment modeling relevance. For example, in the work on social context summarization described in [4, 5] or on comments information retrieval in [2, 6], more than 10 features have been selected and applied. In general, those features are used to build a model for comments classification or clustering, then to summarize or visualize the performance of the model. Particularly, a system can use a filter to visualize the comments data by selecting the features in interest [7]. Those features in text data can be mixed with time series data to build a specific analyzer for financial data and news together [8].

In our framework, we are using the common features (statistics, NLP) and comments specific features. The common features are for example the number of words in each article or comment, the number of sentences, the readability, the similarity with articles, and the divergence with comment or article. In this paper, we focus on features that are used for the specific modeling described in Sect. 4 and other important ones for the framework such as preferences on specific threads and lift measure.

### 3.1    Vector Space Models

Each comment is transformed into a set of n-gram terms. The process requires the natural language processing on the original comment text. For each sentence in the comment, we tokenize it and remove unnecessary punctuation marks connected with the word. A heuristic process is used to remove redundant character to restore the original word. All upper-case letters are changed to lower-case letters. N-gram terms are extracted and selected to represent the comment. They are called the Bag of Words (BOW). Any 1-g, 2-g, or 3-g terms composed entirely with stopwords are deleted.

### 3.2    Lift to Measure Surprise Comments

The vector space model cannot provide information about the association between words. The Lift model is used on pairs of words and to compare the likelihood of the co-occurrence of words happening by chance to that observed in the actual corpus.

The best estimate of the probability of a word occurring in a corpus is given by the observed number of comments that contain that word divided by the total number of comments

$$P(w_i) = \frac{NumContaining(C, w_i)}{\#AllComments} \tag{3}$$

Lift can be generalized to all words $\{w_1, .., w_n\}$ in the comment.

$$Lift_1(\{w_1, \dots, w_n\}) = \frac{P(w_1 \cap w_2 \dots w_i \dots \cap w_n)}{\prod_{i=1}^{n} P(w_i)} \tag{4}$$

$Lift_2$ assumes that the co-occurrence of pairs of words is independent.

$$Lift_2(\{w_1, \dots, w_n\}) = \frac{\prod_{\forall i \neq j} P(w_i \cap w_j)}{\prod_{i=1}^{n} P(w_i)^{n-1}} \tag{5}$$

$Lift_3$ is the logarithm of the $Lift_2$ and useful when the number of comments is large.

$$Lift_3(\{w_1, .., w_n\}) = \left| \sum_{\forall i \neq j} log(P(w_i \cap w_j)) \right| - (n-1) \left| \sum_{\forall i} log(P(w_i)) \right| \tag{6}$$

$Lift_3$ calculates the surprise of a comment as the combination of the mutual information value of every pair of words in that comment using the corpus as a whole to calculate the relevant probabilities. This model is introduced and used for the twitter analysis in [9].

# 4   Analysis Models

The analysis of the event from articles and comment data is concerned with the detection of the burst period, clustering and synthetizing the topics in each burst to create a visual representation of the summarized data. In the following, we describe techniques used to realize those analyses.

## 4.1   Clustering Comments and Multi-dimensional Reduction Scaling

We used K-means clustering algorithm and the multi-dimensional reduction scaling (MDS) to implement this module. The similarity measure is the cosine similarity of the tf-idf vectors. The representation of those clusters in two-dimensional space is performed with the MDS.

## 4.2   Topic Model

Topics are the thematic summary of the comment collection. In other words, it answers the question what themes are those comments discussing.

The topic modeling is used to extract T topics out of the comments collection. That is, we have a set of comments $C_k = \{c_i\}i\varepsilon[1\ldots N]$ and a number of topics $T_k = \{t_i\}i\varepsilon[1\ldots m]$. A comment $c_i$ can be viewed by its topic distribution. For example, $\Pr(c_1 \in t_1) = 0.50$ and $\Pr(c_1 \in t_2) = 0.20$ and so on. The default topic modeling based on LDA is a soft clustering. It can be modified into hard clustering by considering each comment as belonging to a single topic (cluster) $t_r$,

$$r = argmax_r P(t_r|c) = argmax_r P(c|t_r)P(t_r) \tag{7}$$

where $r$ is the number of topics that has the maximum likelihood for each comment. Hence, the output of the LDA based topic clustering approach is an assignment from each comment to a cluster.

Key comments selection within clusters of comments is important for summarizing the contents.

For each topic obtained by the topic modeling, a set of comments are associated. We define the key comment as the top of the comments by ranking them within their clusters. The ranking method is realized by comparing each comment vector (a bag of words) to the list of words which form the topic vector. We use cosine distance for the comparison. The most similar to the topic is the key comment. The experiment on this analysis is first described in [10].

## 4.3   Preferential Attachment (PA) Model for Comments and Commenting an Article

A statistical analysis and comparison of the structure and evolution of the different discussion threads associated to the news websites and Wikipedia comments have shown

that the threads popularity network is following the preferential attachment model or the Yule process [11]. That is, a new comment is added to one of the existing comments with a probability proportional to their existing in-degree. The initial comment in the thread is most like or dislike one by the readers in the sense that the reaction and activity is high.

### 4.4  Visualization Techniques

We want to represent only the most important information in the results to the end-user. A quick overview for understanding the event timeline and story is necessary. Article titles contain the most concise information on the event, so we gather them and process the terms using cloud conversion. The structure of the conversion can also be visualized at glance with a specific network graph database tool.

For the topic model, we use the topic explorer tool [12].

## 5  Test Data and Experiments

### 5.1  Text Data Content

In this paper, we present an analysis of articles and comments on the unexplained disappearance of Malaysia Airlines Flight MH370, that occurred last March 2014, on the Guardian news website. This event was also studied in the work of Wang Lu et al. on timeline summary generation from different news websites such as CNN, NYT (New York Times) and BBC [4]. Their dataset included articles and comments published during the period from March 2014 until June 2014. Our analysis of the event started on March 2014 and ended at the end of September 2015. During that time, there are 368 articles, videos and audio clips published on the Guardian news website[1]. The total number of comments during the experiment period is 77900 and there are 5780 unique comments. The number of commenters is 1780. In Fig. 2, we can observe the comment ratings and user participation activity during the period. We only plotted unique comments by removing redundant comments cross-posted in several articles. Unique comments are selected on their first appearance chronologically in the comment corpus. Therefore, the interval between two time stamps in the horizontal axis, representing the comment time, is not linear. The time stamps in the horizontal axis indicate the vertical dotted line in the grid as the first comment is posted on '2014-3-10 18:41:00'. The peaks are during March and June 2014 and then decreasing in intensity by July 2015.

We can observe in Fig. 3 the repeating bursting and decaying of the article data during this event. The word cloud set to 100 terms on the top left in Fig. 3 represents the initial month of the event. There were 185 articles during the first month on the news website. Following the second Malaysian Airline Flight MH17 explosion, there were 6000 jobs cut at the company during August 2014 represented by the word cloud on the bottom left. The finding of the flaperon debris is shown by the word clouds on the right.

---

[1] http://www.theguardian.com/world/malaysia-airlines-flight-mh370.

Some experts on aviation, specifically the Boeing 777 s, commented in the early stages of the finding (July 2015) that there was no doubt that the debris was from the missing airplane. It took two months for the investigators to give the same conclusion.
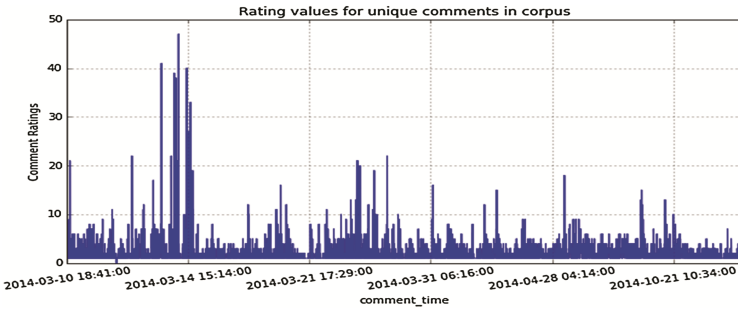


**Fig. 2.** Comment Ratings by users during the one and half year period

Figure 3 shows the number of articles each month and the generated keywords in important periods.
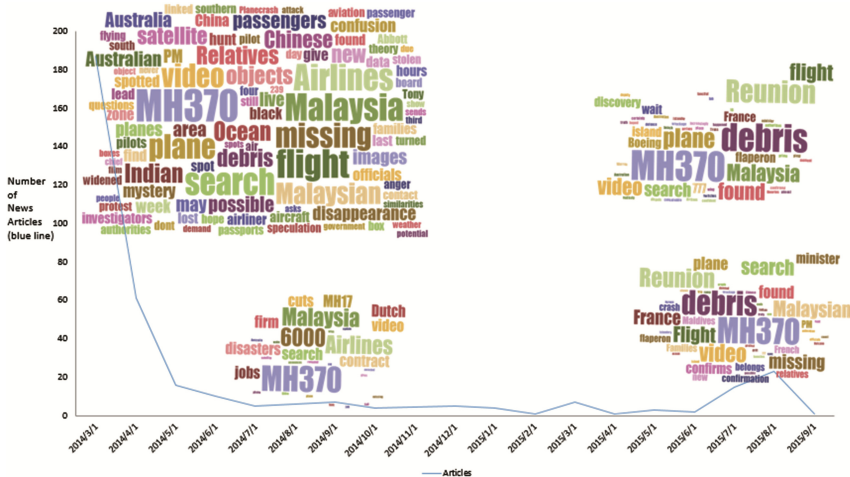


**Fig. 3.** Time series data representing the frequency of articles on the disappearing Malaysian Flight MH370 from March 2014 to September 2015, on the Guardian news website and the word clouds (top left: March 2014, bottom left: August 2014, top right: July 2014, bottom right: August 2014) are depicted.

Many unrelated comments were posted although the Guardian news website is moderated. The discussions on the comment board are analyzed with our tool from our previous research [10].

## 5.2   Clustering and Multidimensional Scaling (MDS) Results

In Figs. 4 and 5, the 368 articles in 10 clusters (represented by 10 colors) are plotted and respectively in five clusters during the 18 month period (18 points). The legend in each plot represents the main keywords as the top terms per cluster. The MDS results show the distance between clusters in a two dimensional space and clearly separate them. Clusters with the term 'Malaysia' are close to each other in the center and clusters with the terms 'search, plane, missing' are at the periphery in the 10 clusters plot. In the 5 clusters plot, the cluster with the term 'SES' (State Emergency Service) is singled out on the right, the cluster with the terms 'debris, Reunion' is at the bottom center and the cluster with terms 'MH17, jobs' is at the top center.
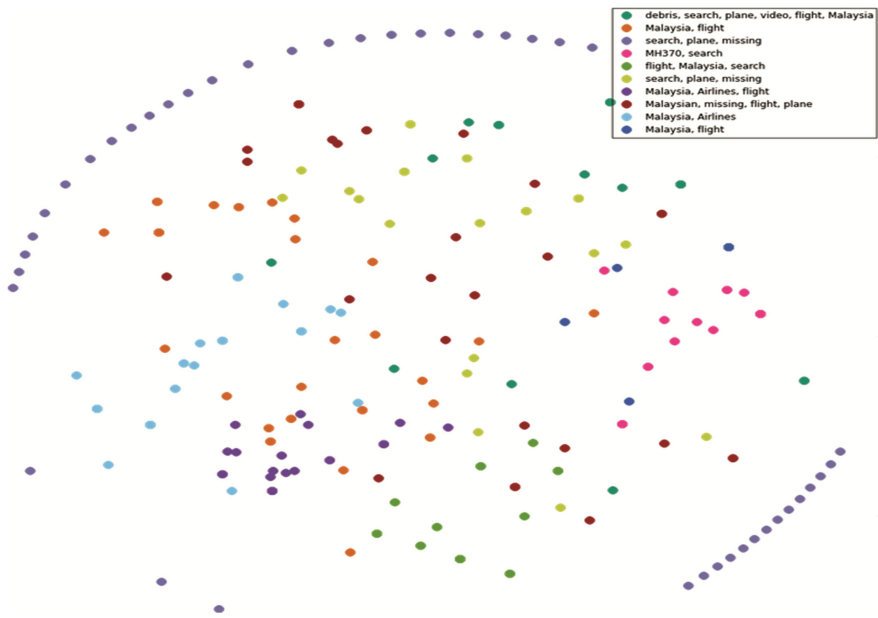


**Fig. 4.** 10 clusters in 2D plots obtained from the 368 articles on the disappearance of Air Malaysia Flight MH370 from March 2014 to September 2015. (Color figure online)

**Fig. 5.** 5 clusters in 2D plots obtained from the 18 months of data on the disappearance of Air Malaysia Flight MH370 from March 2014 to September 2015. (Color figure online)

### 5.3   Surprising Comments with the Lift Measure

We have implemented the surprise model with our own specifications and found some interesting patterns of comments with the model. For the simplicity of the computation, we assume that each sentence in the comment is composed only by the noun phrases (nps). Hence, we compute the occurrence of each pair of these to obtain the lift measure for each unique comment in the text comment data. The extraction of the noun-phrases is realized with the parsing and noun phrases extractor in NLTK [13].

Comments with great Lift3 value contain many sentences (num_sent) in general except the top ranked comment in Table 1 which is surprising. The top ranked comment has only six sentences but upon reading it, find that this comment has the most noun-phrases in it. The top ranked comment in terms of Lift measure is the itinerary of Flight MH370 according to Inmarsat [14].

Otherwise, the comments with $Lift_3$ equal to zero show one or two sentences and are either empty or one noun-phrases (see Table 2). Usually, comments with only one sentence are replies to some previous comments and their contents are brief. We can confirm that $Lift_3$ calculates the correlation of terms in comments by comparing the mutual information value of every pair of words in that comment against the corpus as a whole to calculate the relevant probabilities. The mutual information value of two random variables X and Y is zero if X is independent of Y. Figure 6 shows the values of the Lift3 measure during the period of our experiment.

**Table 1.** Top 5 of Lift3 computation results on the comments from "Malaysian Airline flight MH370" event and their related comment records

| | comment | lift3 | num_sent | nps |
|---|---|---|---|---|
| 4447 | \nTimeline of events, by the hour:\n16:30 UTC ... | 17614.684420 | 6 | [hour interval, acars, penang, strait, radar l... |
| 3895 | \nLets Get Real Here People... If They Are "Na... | 16821.676339 | 21 | [looks, ... .where, global scale your capabili... |
| 4699 | \nElectrical problems developed first from 17:... | 15018.411232 | 32 | [uniquely, fire self, radio compass, acars, vo... |
| 4488 | \nThe Cambodian Connection (cont. - Part 2)\nM... | 11902.477066 | 36 | [complete contradiction, aircraft 's whereabou... |
| 3272 | \nI don't believe it "diverted around Indonesi... | 11888.181240 | 33 | [2:03am local, low altitude, radio compass, us... |

| comment_time | comment_type | comment_id | comment_recommend | comment_content | comment_author | commented_article | lift3 |
|---|---|---|---|---|---|---|---|
| 2014-04-11 15:52:00 | {} | comment-34245297 | 2 | \nLets Get Real Here People... If They Are "Na... | hogan412 | 3z9mt | 16821.676339 |
| 2014-05-03 20:02:00 | {} | comment-35166215 | 7 | \nThe Cambodian Connection (cont. - Part 2)\nM... | CliveG7 | 3zpnt | 11902.477066 |
| 2014-05-04 13:10:00 | {} | comment-35191013 | 6 | \nTimeline of events, by the hour:\n16:30 UTC ... | EnoughIsEnoughUS | 3zpnt | 17614.684420 |
| 2014-06-28 18:05:00 | {u'reply_ref': u'#comment-37531499', u'reply_t... | comment-37540502 | 4 | \nI don't believe it "diverted around Indonesi... | Sy Gunson | 3qejk | 11888.181240 |
| 2014-10-05 08:58:00 | {u'reply_ref': u'#comment-41801533', u'reply_t... | comment-41804323 | 1 | \nElectrical problems developed first from 17:... | Sy Gunson | 426bm | 15018.411232 |

**Table 2.** Last 15 of Lift3 computation results on the comments

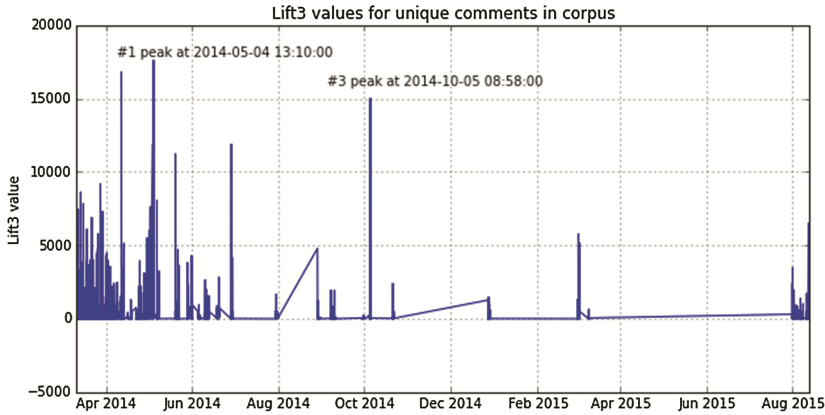| | comment | lift3 | num_sent | nps |
|---|---|---|---|---|
| 0 | \nThey did not have, nor do they need visas. T... | 4.401002 | 3 | [beijing, kind, chinese government] |
| 2 | \nTerrorists (and spies) use stolen passports.... | 7.875785 | 6 | [live, fact, terrorists] |
| 3 | \nYou can assume whatever you like - after all... | 0.000000 | 1 | [] |
| 4 | \nIt's a roaming thing, you call out from your... | 1.520780 | 2 | [network tries, ca n't] |
| 5 | \nBecause by understanding what has happened t... | 2.365878 | 3 | [aviation industry, good safety record] |
| 9 | \nThat is an excellent question.. There is no ... | 9.329782 | 2 | [excellent question.., 's database, interpol] |
| 11 | \nAnd refused to allow outside assistance unti... | 2.515640 | 2 | [google, expendable means] |
| 13 | \nI was surprised to hear that an aircraft can... | 2.842999 | 4 | [mid air, dangerous parts] |
| 14 | \nmust have been a group who went shopping and... | 0.000000 | 1 | [vip] |
| 15 | \nLooks like you might be right. So where does... | 0.000000 | 2 | [looks] |

**Fig. 6.** The variation of the Lift measure during the experiment period and the peak during the Inmarsat data publication on May 4th 2014.

### 5.4   Visualization of the Comments Network Graph

Figure 7 depicts the structure of conversations and the preferential attachments to some users' or comment posts. They are realized with Neo4j[2]. Neo4j is a tool which we can use to build a property graph from a database. The property graph contains connected entities (the *nodes*) which can hold any number of attributes (key-value pairs) and *relationships* which give meaning to the connection between two-node entities. A relationship always has a direction, a type, a start node and an end node. It can also have properties. Neo4j implements the property graph and optimizes the graph database to output query results faster independent of the size of the database. The comment corpus is transformed into a comment database where each comment is a record containing the properties described in Eq. (1). The modeling of this database in Neo4j consists of defining the nodes and the relationships. In our comments visualization, we want to observe users, comments and their relationships. Therefore, we define users as entities such as people with names and identifications. The Guardian news commenting facility does not allow anonymous people to post comments. Hence, each user must have a single user name and user identification. Comments are also entities in the graph with all properties they have in the record. We define the relationships that can be modeled as the interactions between users and comment entities. After the data modeling, we put the comment database in the Neo4j server and build the graph network describing the interactions of those entities.
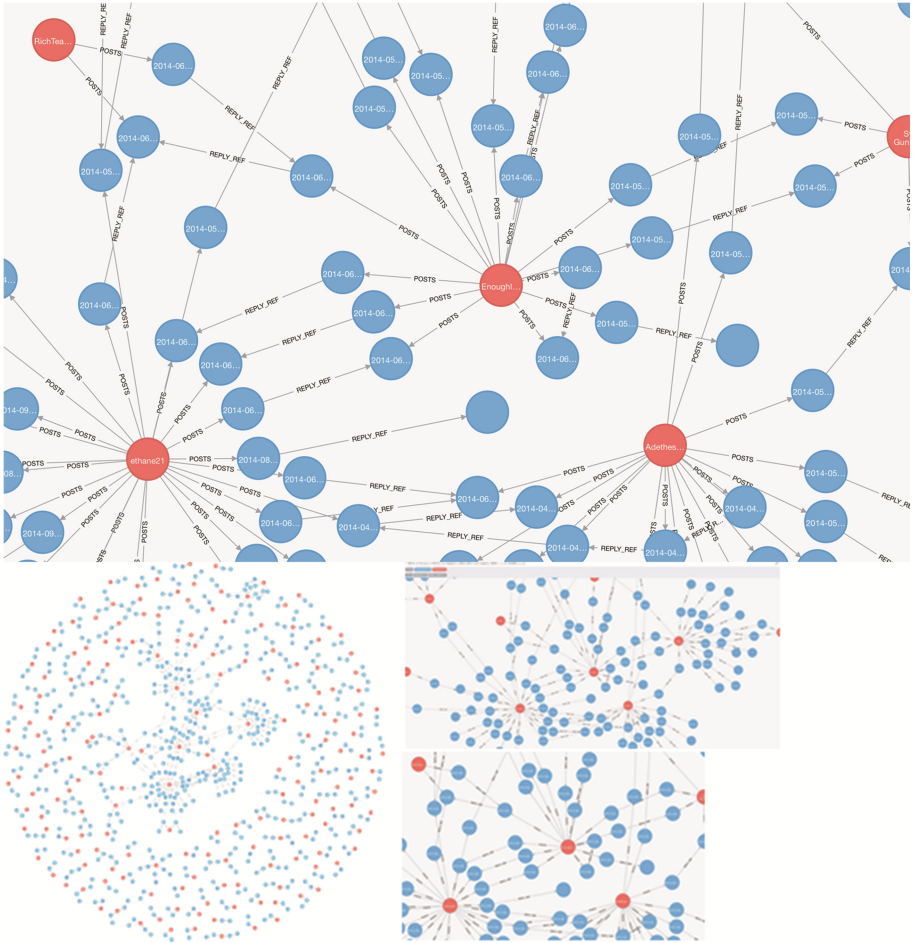
---

[2] http://neo4j.com.

**Fig. 7.** The comment graph network on the comments data with minimum rating equal to 3 (Color figure online)

In Fig. 7, users (red nodes) are labeled with their usernames. They can post several comments. The most active users are centered on the graph with the highest number of comments and replies. Comments are labeled with the time stamp (blue nodes). The three types of relationship are users posting comments, comments that are replies to other comments and comments quoted other users. By clicking a comment node, we can read its properties such as subject, content, and rating score. The graph can be zoomed in or out and queried in SQL-like format.

## 6   Discussions

In this paper, we introduced our framework for aggregating and analyzing articles and comments on news websites after events occur in the real world. We proposed several methods and models to perform the analysis.

We presented the case study on the Malaysia Airline Flight MH370 event. The summarization and clustering results are shown in Sect. 5. The truth regarding this event is not yet known. We did text mining and comment analysis in order to summarize and synthetize (clustering and graph network analysis) the vast amount of text data available on the event from the news website.

Surprisingly the Lift3 measure extracted the timeline of the event based on the top ranked comment in this experiment. By analyzing the features in this comment, we observed that it contains the highest number of name entities, a lower sentiment polarity and lower subjectivity tone.

The PA model is valid for the commenting behavior on articles discussing the event. Most of the articles do not have comments. Only the articles already commented on present active comments unless the moderator closed the commenting tool. Usually, the most popular articles are the ones most shared and commented on. This model will be used in the future to predict which articles to monitor using our comment analysis tool.

In the future, we will consider the optimization of the program to compute the $\text{Lift}_3$ values because computations of the joint probabilities of two terms are very expensive for comments with many sentences and corpuses with more than ten thousands comments for the duration of the event. The real-time summarization as in the IncreSTS [1] should also be taken into account. We will apply our framework to other domains such as finance data and text events.

## References

1. Liu, C.Y., Tseng, C.Y., Chen, M.S.: IncreSTS: towards real-time incremental short text summarization on comment streams from social network services. IEEE Trans. Knowl. Data Eng. **27**(11), 2986–3000 (2015). http://doi.ieee.org/10.1109/TKDE.2015.2405553. IEEE Press, New York
2. Momeni, E., Sageder, G.: An empirical analysis of characteristics of useful comments in social media. In: Proceedings of the Conference on Web Sciences, WebSci 2013, Paris, France, 02–04 May 2013, pp. 258–261. ACM, New York (2013). http://doi.acm.org/10.1145/2464464.2464490
3. Schinas, M., Papadopoulos, S., Kompatsiaris, Y., Mitkas, P.A.: Visual event summarization on social media using topic modelling and graph-based ranking algorithms. In: Proceedings of the 5th ACM International Conference on Multimedia Retrieval, ICMR 2015, Shanghai, China, 23–26 June 2015, pp. 203–210. ACM, New York (2015). http://doi.acm.org/10.1145/2671188.2749407
4. Wang, L., Cardie, C., Marchetti, G.: Socially-informed timeline generation for complex events. In: The Proceedings of the 2015 Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Denver, Colorado, USA, 31 May–5 June 2015, pp. 1055–1065. ACL Anthology publisher (2015). http://www.aclweb.org/anthology/N15-1112

5. Yang, Z., Cai, K., Tang, J., Zhang, L., Su, Z., Li, J.: Social context summarization. In: The Proceedings of the 34th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2011, Beijing, China, 24–28 July 2011, pp. 255–264. ACM, New York (2011)

6. Potthast, M., Stein, B., Loose, F., Becker, S.: Information retrieval in the commentsphere. ACM Trans. Intell. Syst. Technol. (TIST) **3**(4), 68 (2012). ACM, New York

7. Wanner, F., Ramm, T., Keim, D.A.: ForAVis – explorative user forum analysis. In: Proceedings of the International Conference on Web Intelligence, Mining and Semantics, WIMS 2011, Sogndal, Norway, 25–27 May 2011, pp. 1–10. ACM, New York (2011)

8. Wanner, F., Schreck, T., Jentner, W., Sharalieva, L., Keim, D.A.: Relating interesting quantitative time series patterns with text events and text features. In: Proceedings of the IS&T/SPIE Volume 9017, Electronic Imaging, International Society for Optics and Photonics, Visualization and Data Analysis 2014, San Francisco, USA, 2 February 2014, pp. 1–15. SPIE, USA (2014). doi:10.1117/12.2039639

9. Rooney, S.: Using lift as a practical measure of surprise in a document stream. In: The EBW 2015, Proceedings of the 3rd International Conference on E-Technologies and Business on the Web, Paris, France, 26–28 March 2015, pp. 7–12. SDIWC, Hong Kong (2015)

10. Ramamonjisoa, D., Murakami, R., Chakraborty, B.: A framework for comments analysis and visualization. Int. J. Digit. Inf. Wirel. Commun. (IJDIWC) **5**(3), 179–187 (2015). SDIWC, Hong Kong

11. Gomez, V., Kappen, H.J., Litvak, N., Kaltenbrunner, A.: A likelihood-based framework for the analysis of discussion threads. World Wide Web **16**((5-6)), 645–675 (2013). Social Networks and Social Web Mining

12. Murdock, J., Allen, C.: Visualization techniques for topic model checking. In: Proceedings of the 29th Conference on Artificial Intelligence, AAAI 2015, Austin Texas, USA, 25–30 January 2015, pp. 4284–4285. AAAI Press, USA (2015

13. Bird, S., Klein, E., Loper, E.: Natural Language Processing with Python, 1st edn. O'Reilly, Sebastopol (2009)

14. Ashton, C., Bruce, A.S., Colledge, G., Dickinson, M.: The search for MH370. J. Navig. **68**(1), 1–22 (2015). doi:10.1017/S037346331400068X. The Royal Institute of Navigation 2014