

Detection and Prediction of Users Attitude Based on Real-time and Batch Sentiment Analysis of Facebook Comments

Hieu Tran and Maxim Shcherbakov

Volgograd State Technical Univeristy, Computer Aided Design Department,
Lenin avenue 28, 400005 Volgograd, Russia
trantrunghieu7492@gmail.com,maxim.shcherbakov@vstu.ru
<http://www.vstu.ru/en>

Abstract. The most of the people have their account on social networks (e.g. Facebook, Vkontakte) where they express their attitude to different situations and events. Facebook provides only the positive mark as a like button and share. However, it is important to know the position of a certain user on posts even though the opinion is negative. Positive, negative and neutral attitude can be extracted from the comments of users. Overall information about positive, negative and neutral opinion can bring understanding how people react in a position. Moreover, it is important to know how attitude is changing during the time period. The contribution of the paper is a new method based on sentiment text analysis for detection and prediction negative and positive patterns for Facebook comments which combines (i) real-time sentiment text analysis for pattern discovery and (ii) batch data processing for creating opinion forecasting algorithm. To perform forecast we propose two-steps algorithm where: (i) patterns are clustered using unsupervised clustering techniques and (ii) trend prediction is performed based on finding the nearest pattern from the certain cluster. Case studies show the efficiency and accuracy (Avg. MAE = 0.008) of the proposed method and its practical applicability. Also, we discovered three types of users attitude patterns and described them.

Keywords: Opinion Mining, Sentiment Analysis, Text Classification, Social networks, Facebook comments, Real-time and batch processing, Clustering analysis, Forecasting techniques

1 Introduction

Sentiment analysis of textual content is used for opinion mining of people who express their emotions and thoughts by text messages. New communication platforms such as social networks (e.g. Facebook or VKontakte) gives a new opportunity for better understanding information using natural language processing and sentiment analysis. According to the article published by zephoria.com in December 2015, nowadays Facebook has more than 1.55 billion monthly active

users [21]. These users write more 510 000 comments every minute and this is a source of large information on the Internet. Usually, these textual comments are the results of the reaction of people regarding recent news or happened events. Understanding of users attitude helps to know how a certain person or groups respond to the particular topic, and it serves to draw relevant conclusions or make efficient decisions based on feedback [8,20]. For example, in the political field. Assume there is news regarding a particular decision of the government in a certain country, which published in social networks by BBC or CNN. Based on examination of the textual comments, we can understand people positions and either a certain person supports this decision of the government or not.

From the business point of view, sentiment analysis helps companies to improve customer development process, enhance business intelligence systems, and change their marketing strategies to get more profit. Moreover, using this type of text analysis, the trend of peoples attitude to certain events or typical groups of events can be predicted. This foresight is valuable for proactive actions development for the future expected situation in every domain we refer to, such as politic, economic, business and so on. So, in fact, the questions is how to understand users behaviour and opinions according to processing textual comments in social network and how to predict either these opinions remain the same or will be changed in time? Opinions give the intuition about a person or customer preferences.

The main problem, which is considered in the current research, is how to understand positive or negative users opinion about published posts and news using sentimental text analysis. Are there any laws and consistent timewise patterns in users comments, and how to detect these patterns and predict them? The contribution of the paper is a new method based on sentiment text analysis for detection and prediction negative and positive patterns in Facebook comments.

The paper contains the following sections besides of introduction. The next section contains the literature review and analysis of the recent related works on sentiment text analysis. After it describes the main idea of the proposed method of Facebook comments sentiment analysis using a combination of the real-time and batch data processing. Results and discussion are covered in the last section.

2 Related works

Sentiment text analysis is a large but still growing research domain. An early, and still common, approach to sentiment analysis has been to use the called semantic orientation (SO) of terms as the basis for classification [11].

Turney showed that semantic orientation is useful for classifying more general product reviews [24]. The work suggests that product reviews may be classified easier than movie reviews, as product reviews tend to be shorter and be more explicitly evaluative. In [19] authors classified movie reviews used standard bag-of-words techniques with limited success. Twitter is a social network which represented as a sources of customer opinions to analyze. The early results of Twitter data sentiment analysis presented in work [2]. The authors of

the paper "Sentiment Analysis on Twitter" tentatively conclude that sentiment analysis for Twitter data is not that different from sentiment analysis for other genres [1]. In [14] authors used the similar way (Twitter API) to collect training data and then to perform a sentiment investigation. An analysis of Twitter followers reveals networks of users who are related by current news tops rather than by personal interactions. Furthermore, a database of sensor data from the reality mining corpus is used for dynamic social network analysis [6]. Besides social networks, common websites have a large data needed to investigated. As practical implementation, sentiment analysis was applied to the feedback data from Global Support Services survey [7]. It helps organizations determine the quality of services. Text information from social networks provides information about entities (people and organizations) and their corresponding relations and involvement in events [26]. Up to this time, there are not too many articles about real-time text data analysis. In the research [3] proposed the theory of big data stream analytics for language modeling for sentiment analysis but did not put into practice with a certain network. A group authors from the University of Southern California described a system for real-time analysis of public sentiment toward presidential candidates in the 2012 U.S. election as expressed on Twitter [25] but this system only works with the twitters, and could not reach to users replies.

Almost all of research is concerned about status, twitter, or post to analysis sentiment [12,15,13]. There is a small number of research papers mentioned the replies and comments analysis. However, replies to Twitter and Facebook comments have values of sentiment expressions from users. In [10], authors implemented analysis of Facebook posts [27] on the 2014 Thai general election. They did not use a large information from the comments of each post. Therefore, we focus on these comments which will be performed in this paper. Also in [4] presented prediction the vote percentage that individual candidate will receive on Singapore Presidential Election 2011 using Twitter data with census correction but they met issues with fake tweeter sentiment. And the source may be related to the scenario where the voters do not truly reflect their on-line sentiments from their choice of candidate. Several case studies have found that the online information has been quite successful acting as an indicator for electoral success [27]. There are some research on the use of twitter such as [23]. There is the Vietnamese research group do analysis based on Facebook comments, but they worked in a sphere of Vietnamese language [22].

Based on literature review we conclude, that sentiment analysis of text comments on the social network as a mechanism to understand the people's attitude is still an open question.

3 A method

3.1 General scheme

The main problem we would like to solve is the creation of the technique, which helps understand and predict user attitude expressed in Facebook comments

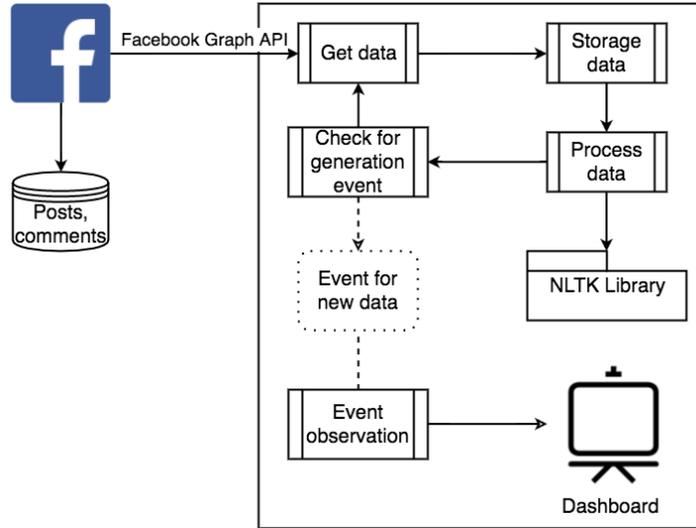


Fig. 1. A scheme for real-time stream processing

regarding published news or post. We propose the following method: (1) using collected from Facebook data we perform batch analysis to make the precise forecasting model to detect and predict negative and positive behaviour patterns; (2) using real-time text processing we detect a pattern and understand what is the current situation looks like; (3) understand how this situation will be developed based on pattern prediction; and (4) perform actions to change the pattern if it needed.

To be more concrete, we implement two main approaches for analysis sentiment in our method. The first approach is real-time comments analysis. The solution allows to analyze and update results right in time when data generated by users on Facebook. The second approach analyses stored data in batch mode. This is a cluster-based approach for the negative/positive attitude patterns detection and prediction. To evaluate the performance of forecasting, the median absolute deviation measure is used.

3.2 Real-time stream sentiment analysis and event generation

Real-time stream processing solution retrieves data from Facebook server continually and then, processes a data package in minor period of time almost real-time processing. The NLTK library is used for sentiment text analysis [16]. The results of data processing are checked by predefined users conditions. If it satisfies conditions, the solution creates an event to update dashboards status. Moreover, the real-time solution includes a procedure for listing to events. If a certain event occurred, the dashboard will be updated. Figure 1 presents a proposed scheme for real-time stream processing.

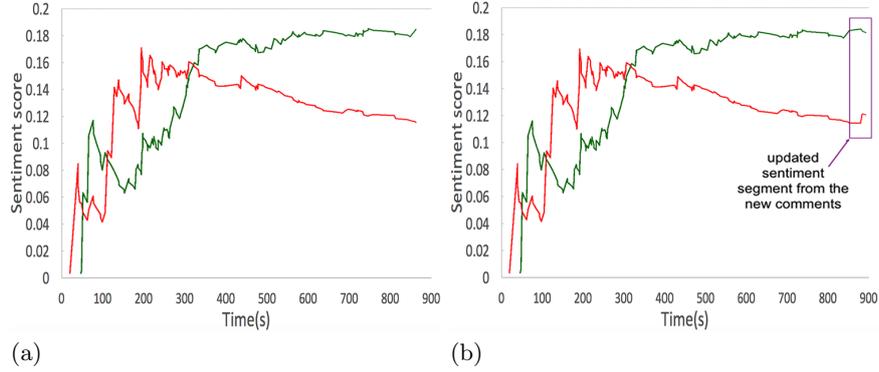


Fig. 2. Dashboard represents (a) negative/positive sentiment of the first 750 seconds of post’s life; (b) updated results after synthetic negative comment

Data retrieving. To keep a set of comments updated, a loop statement is used to receive comment data from Facebook. Time (T_u) between two loops statements is defined in advance. We set $T_u = 0.1$ seconds as it is nearly real-time and allows to store data continually updated. However, this parameter might be changed if needed. Graph API Facebook is used in every step of the loop to obtain all comments of the selected post and to transfer data for storing and further analysis. This approach provides assurance, that collected data is fresh.

Event generation. In the inner loop, when data is transferred for further processing, the comparison of a new segment of data, which has been obtained from Facebook recently with cached data, is made. It is crucial as it helps us find out the changes in data. If changes in data are discovered, the event needs to be generated. Observers will receive this event further.

Event Observation. To evaluate the performance of the method, the program implementation includes the procedure which allows observing any generated event. In occasion of the event reflecting changing data, the solution updates system automates re-analyze sentiment data. This involves the change of dashboard status. Furthermore, cache data is updated as well to be sure we have newest data from the server. For instance, we analyse the sentiment of comments on a post about a political topic entitled ‘Obama bans solitary confinement for juveniles and low-level offenders’ on CNN news channel on Facebook¹.

Figure 2 shows the dashboard representing the negative and positive attitude in different timestamps. The red line represents values of negative sentiment.

¹ The title of the post: Obama bans solitary confinement for juveniles and low-level offenders, <https://www.facebook.com/bbcnews/posts/10153348871732217>

The green line is the positives. These values are fall into the interval $[0; 1]$. Every moment t , we have a sum of negative, positive and neutral equals to 1:

$$V_t^{(p)} + V_t^{(n)} + V_t^{(u)} = 1, \quad (1)$$

where, $V^{(p)}$ - denotes the positive score, $V^{(n)}$ - denotes the negative score, $V^{(u)}$ - denotes the neutral score. In this paper, we do not place the neutral scores in the graph. The scale 'time' represents the time of comment posting as the interval from creating a post in seconds. The 'sentiment analysis' scale is the polarity value of comments.

The first figure explains the behavior for every 780 seconds of post life and the second one, reaction on the posted synthetic negative comment. Real-time analysis allows to detect current patterns and to compare obtained pattern with expected or required. However, finding and adjusting those references patterns depends on expert (human) intervention and due to high velocity and the variety of data, this procedure is very costly.

3.3 Batch data processing

As method includes pattern detection and forecasting models, another component of our system implements pattern discovery in batch mode via processing of high volumes of certain topic data collected over a long period of time. Comments in the form of textual content regarding the topic are collected from several popular pages. The number of posts is selected arbitrary (100 posts, 1000 posts, or even more). Our solution uses NLTK and the results are used for sentiment pattern detection and prediction. Figure 3 shows the scheme for batch processing of sentiment analysis.

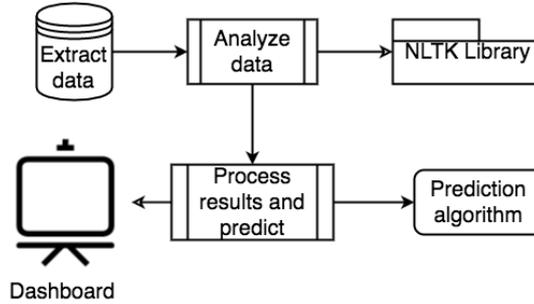


Fig. 3. Scheme for batch processing

Data collecting Implementation of batch data processing makes sense in the case of high volumes data. Firstly, we chose a topic, which is popular recently.

For each post, using Facebook Graph API, all comments have been collected during the first 30000 seconds. Data is stored in flat table format (e.g. CSV file) which is easy to save in distributed file system. The header of CSV file contains the following columns: [Datetime] [Topic] [Post] [Comment] [Positive] [Negative].

The first column contains a value in seconds when a comment was appeared, which is counted from initial posts appearance. Data type of the value is the integer. The second column contains the topics title. The third column is posts title. The fourth column is the content the text of a comment. The last two columns represent sentiment analysis scores (negative and positive) of comment on the post. The data type of scores is the float number.

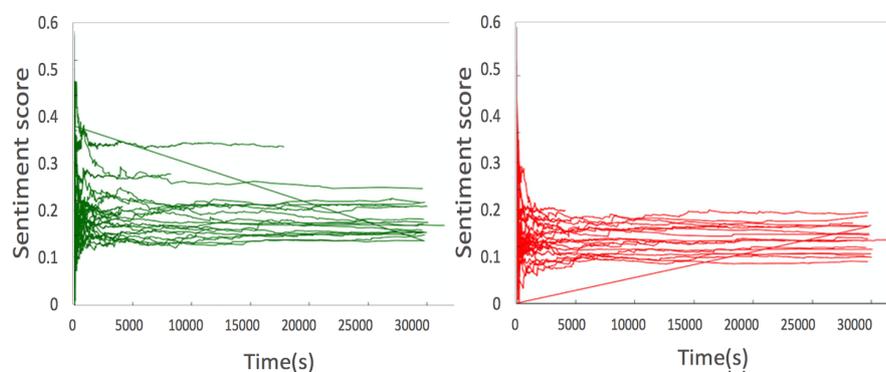


Fig. 4. Positive sentiment time series (left side) for a set of posts and negative sentiment time series (right side). Data about U.S. presidential election 2016 topics was obtained by CNN and BBC feeds

Analysis data. As it was mentioned above, the NLTK sentiment analysis was applied for each comment for the post. As the results, we obtain time series which describes peoples negative $V^{(n)}$ and positive $V^{(p)}$ scores during the time. Figure 4 expresses positive scores time series (left side) for a set of posts and the negative (right side).

Clustering data. The output of the previous step is a set of time series and the number of time series is equal to a number of examined posts. We use the fact, that some of the time series has a similar pattern and they could be arranged into 3 or 4 different groups. It allows defining a 'typical' behaviour of the people from a sentiment point of view. For instance, the certain post could have a high negative expression in the beginning and fade negative afterwards. The well-known technique for unsupervised grouping is clustering. In this paper, two machine learning approaches: k -means and MB-means were used to cluster sentiment of posts. Needless to say, k -means has been studied and applied in

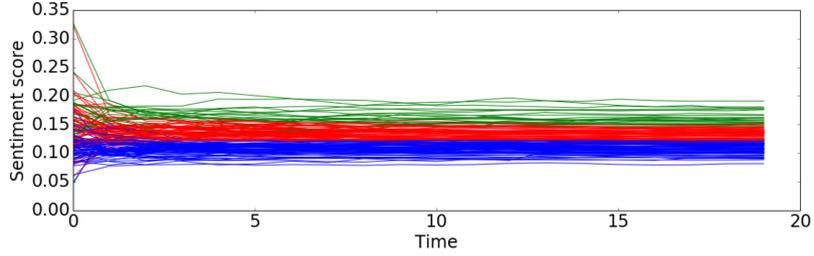


Fig. 5. The results of clustering using (a) k -means algorithm and (b) MB-means algorithm, where every plot has color according to the cluster

a wide range of domains, e.g. transportation network analysis [9], information security [28], pattern recognition [5], text classification [19] and many others domain.

We implement clustering, where initial time series are placed in three difference groups according to their characteristics. To perform clustering, each of posts expressed as a vector with 20 components (features): [id][Post][Period 1][Period 2]... [Period 20].

These features describe to our first 30,000 seconds of the posts. The values of 30,000 seconds have been chosen arbitrary and might be changed if needed. Every feature has the same period time, which is about 1,500 seconds. The value of each feature is the average value of all value in this period. All-time series (or vectors) have been divided into three clusters using k -means algorithm. The number of clusters picked up according to preliminary analysis. For comparison with k -means method, we also use the other algorithm to cluster which is named MB-means. Figure 5 shows the results of clustering, every time series has the color according to the cluster. However, in spite of difference clustering techniques, the results look quite similar.

We can describe three clusters in the following way.

- The first cluster which is red lines (based on k -means algorithm results) begins with high positive scores. It decreased quickly from the beginning to the third period, then it continuously went downward slowly. From the fifth period, it felt the lowest value. Then it remained stable in the next periods.
- With the second cluster which is green lines using k -means, we have a graph remained relatively stable from the beginning. Almost periods time, there were small increases or decreases.
- The last cluster is blue lines using k -means. It started with low positive scores. It grew up rapidly, and then it reached the peak of the score at the second period. Then there were slightly drop in the third period and leveled off. From the next period, it stayed constant and there are no more changes.

In conclusion, each cluster has characters itself obviously. It is the good signification to cluster testing data.

Prediction. Clustering allows defining typical patterns in people’s behaviour. The next task is prediction the trend development of peoples attitude on a post using data about the first reaction (or 5 values in vector representation). Note, we choice 5 values for the current research, but this parameter is subject to choice.

The prediction is performed in two steps. The first step including cluster detection procedure. Based on the first 5 features (or 5 values in vector representation) of given pattern which needs to be predicted, we select the appropriate cluster. The selected cluster is the cluster having the average profile for this 5 features with minimum deviation with our post (in terms of Euclidean distance). In the second step, we predict the trend development of peoples attitude on this post. The nearest (in terms of distance) time series from the certain cluster is found and use as the prediction for the rest time interval. The Mean Absolute Error (MAE) is used to evaluate the performance of the forecasting technique.

$$MAE = \frac{1}{n_{test}} \sum_{i=1}^{n_{test}} \left(\frac{\sum_{j=6}^{20} |h_j - h_j^*|}{15} \right)_i, \quad (2)$$

where, h_j – denotes the real value at the j -th timestamp, h_j^* – denotes the predicted value at the j -th timestamp, n_{test} – a number of time series included in test data set.

4 Results and discussion

To evaluate our approach we designed and implemented software solution using Python and Facebook API. For our experiments, we used the topic the United States presidential election 2016. All data was received from two famous new channels: BBC news and CNN on Facebook. We collected 200 posts about the mentioned above topic. For every post, using Facebook Graph API, all comments have been collected during the first 30,000 seconds. Approximately, the total number of comments is about 100,000 comments.

The results of real-time sentiment analysis indicate users emotions and thoughts in a real-time stream. Based on monitoring of dashboards, we can understand that the ratio of positive sentiment is higher that the negative sentiment. That means that the proportion of people who supports the news is more than the proportion who against. At the time when the intersection of the red line and the green line is observed, its the point that the trend of attitude is changed. It might the point that perhaps people are changing the attitude to the topic. Based on these points, we can create a line which is the general trend of peoples attitude. Also, it is possible to set up triggers indicates when an interested event occurs.

Using clustering techniques we are able to detect the most typical behaviour of the users and describes them. For instance, we observe that negative or positive estimations asymptotically approaching to a certain level and never exceed the threshold. Negative and positive attitude is fading during the time, and we

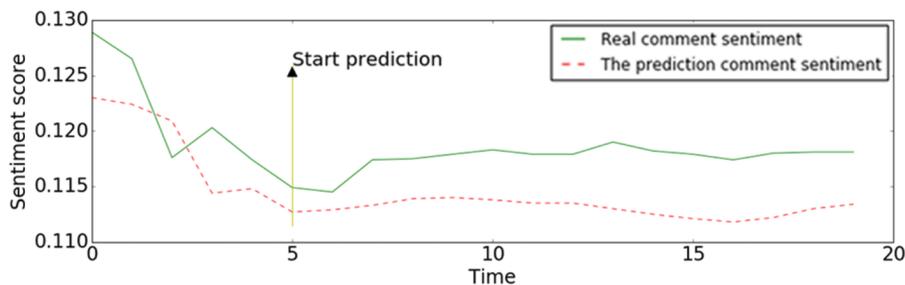


Fig. 6. Prediction trending development of positive sentiment

are able to estimate time of popularity of the post and advise actions to support popularity. Also, our technique allows predicting the trend development of peoples attitude. It could be a framework to detect the outliers in the comment of Facebooks community. To evaluate forecasting performance MAE has been applied as error measurement. To avoid a case where results obtained by chance, we developed cross validation (with folds = 10) and get average MAE = 0.008. Figure 6 gives the representation of results of trend forecasting. The green line is the real comment sentiment on this post. And the red line is the prediction line for the development of peoples attitude on the post.

5 Conclusion

In this study, we perform actions to understand users preferences and attitude based on sentiment analysis of Facebook comments and application of machine learning techniques. Detection of laws and consistent patterns in users comments published in time framework allows providers of services and sales to react in real time and be more proactive using trend prediction. We propose a new method based on sentiment text analysis for detection and prediction negative and positive patterns for Facebook comments which combines (i) real-time sentiment text analysis for pattern detection and (ii) batch data processing for creating forecasting algorithms. To perform forecast we propose two-steps algorithm where: (i) patterns are clustered using unsupervised clustering techniques and (ii) trend prediction is performed based on finding the nearest pattern from a certain cluster.

Based on the results, we found three types of user behavior in their opinion expression and find that our simple forecasting technique is very accurate. Proposed method can be readily used in practice by sales companies who can use the real-time approach for learning their customer attitude about products and making the assessment of a product. Some social and political organizations use to analyze community on a certain event such as The 2016 U.S. election. Our future work will be continued by focusing on improvement the model training,

improvement method of prediction using MAE. Besides, the next stage of our research will analyze a group of people such as Vietnamese and Russian [17,18].

Acknowledgments. The reported study was partially supported by RFBR, research project No. 16-37-60066 and research project MD-6964.2016.9.

References

1. Agarwal, A., Xie, B., Vovsha, I., Rambow, O., Passonneau, R.: Sentiment Analysis of Twitter Data. In: LSM '11 Proceedings of the Workshop on Languages in Social Media, pp. 30–38. Association for Computational Linguistics Stroudsburg, PA, USA (2011)
2. Birmingham, A., Smeaton, A.: Classifying sentiment in microblogs: is brevity an advantage is brevity an advantage? In: 19th ACM international conference on Information and knowledge management, pp. 1833–1836. ACM (2010)
3. Cheng, O.K.M., Lau, R.: Big Data Stream Analytics for Near Real-Time Sentiment Analysis. *Journal of Computer and Communications*. 3, 189–195 (2015)
4. Choy, M., Cheong, M.L.F., Laik, M.N., Shung, K.P.: A sentiment analysis of Singapore Presidential Election 2011 using Twitter data with census correction. arXiv:1108.5520 (2011)
5. Dash, M., Liu, H.: Feature selection for classification. *Intelligent Data Analysis*. 1(3), 131–156 (1997)
6. Eagle, N., Pentland, A.: Reality Mining: Sensing Complex Social Systems. *Personal and Ubiquitous Computing*. 10(4), pp.255–268 (2006)
7. Gamon, M.: Sentiment classification on customer feedback data: noisy data, large feature vectors, and the role of linguistic analysis. In: 20th international conference on Computational Linguistics, pp. 841–es, Association for Computational Linguistics. Stroudsburg, PA, USA (2004).
8. Godbole, N., Srinivasiah, M., Skiena, S.: Large-Scale Sentiment Analysis for News and Blogs. *ICWSM*. 7(21), 219–222 (2007)
9. Golubev, A., Chechetkin, I., Solnushkin, K.S., Sadovnikova, N., Parygin, D., Shcherbakov, M.V.: Strategway: web solutions for building public transportation routes using big geodata analysis. In: 17th International Conference on Information Integration and Web-based Applications & Services (iiWAS '15), pp. 91:1–91:4. ACM, New York, NY, USA (2015)
10. Grömping, M.: Echo Chambers: Partisan Facebook Groups during the 2014 Thai Election. *Asia Pacific Media Educator*. 24, 39–59 (2014)
11. Hatzivassiloglou, V., McKeown, K.R.: Predicting the semantic orientation of adjectives. In: 35th ACL and 8th EACL, pp. 174–181. ACL, Somerset, New Jersey (1997)
12. Honeycutt, C., Herring, S.C.: Beyond microblogging: Conversation and collaboration via Twitter. In: 42nd Hawaii International Conference on System Sciences, pp. 1-10 (2009)
13. Jansen, B.J., Zhang, M., Sobel, K., Chowdury, A.: Twitter power: Tweets as electronic word of mouth. *Journal of the American Society for Information Science and Technology*. 60, 1–20 (2009)
14. Kwak, H., Lee, C., Park, H., Moon, S.: What is Twitter, a Social Network or a News Media? In: 19th International Conference on World Wide Web, pp. 591-600. (2010)

15. Lau, R.Y.K., Xia, Y., Ye, Y. A Probabilistic Generative Model for Mining Cybercriminal Networks from Online Social Media. *IEEE Computational Intelligence Magazine*. 9, 31–43 (2014)
16. Natural Language Toolkit, <http://www.nltk.org/>
17. Nguyen, V.H., Nguyen, H.T., Snaes, V.: Normalization of Vietnamese tweets on Twitter. In: 2nd Euro-China Conference on Intelligent Data Analysis and Applications (ECC 2015), pp. 179–189. Springer International Publishing, Cham (2015).
18. Nguyen, H.T., Duong, P.H., Le, T.Q.: A Multifaceted Approach to Sentence Similarity. In: 4th International Symposium on Integrated Uncertainty in Knowledge Modelling and Decision Making (IUKM 2015), pp. 303–314. LNAI, Springer (2015)
19. Pang, B., Lee, L., Vaithyanathan, D.: Thumbs up? sentiment classification using machine learning techniques. In: ACL-02 conference on Empirical methods in natural language processing (EMNLP '02), Vol. 10, pp. 79–86. Association for Computational Linguistics, Stroudsburg, PA, USA (2002)
20. Pang, B., Lee, L.: Opinion mining and sentiment analysis. *Found. Trends Inf. Retr.* 2, 1–135 (2008).
21. The Top 20 Valuable Facebook Statistics Updated December 2015, <https://zephoria.com/top-15-valuable-facebook-statistics/>
22. Trinh, S., Nguyen, L., Vo, M., Do, P.: Lexicon-Based Sentiment Analysis of Facebook Comments in Vietnamese Language. *Recent Developments in Intelligent Information and Database Systems*. 642, 263–276 (2016)
23. Tumasjan, A., Sprenger, T.O., Sandner, P.G., Welpe, I.M.: Predicting Elections with Twitter: What 140 Characters Reveal about Political Sentiment. In: Fourth International AAAI Conference on Weblogs and Social Media, pp. 178–185. AAAI (2010)
24. Turney, P.D.: Thumbs up or thumbs down? semantic orientation applied to unsupervised classification of reviews. In: 40th Annual Meeting of the ACL (ACL02), pp. 417–424. ACL, Philadelphia, Pennsylvania (2002)
25. Wang, H., Can, D., Kazemzadeh, A., Bar, F., Narayanan, S.: A system for real-time Twitter sentiment analysis of 2012 U.S. presidential election cycle. In: the ACL 2012 System Demonstrations, pp. 115–120. Association for Computational Linguistics Stroudsburg, PA, USA (2012)
26. Weinstein, C., Campbell, W.M., Delaney, B.W., OLeary, G.: Modeling and Detection Techniques for Counter-Terror Social Network Analysis and Intent Recognition. In *IEEE Aerospace Conference*, pp. 1-16. IEEE Press (2009)
27. Williams, C., Gulati, G.: What is a Social Network Worth? Facebook and Vote Share in the 2008 Presidential Primaries. In: Annual Meeting of the American Political Science Association, pp. 1-17. Boston, MA (2008)
28. Zhang, Y., Dang, Y., Chen, H., Thurmond, M., Larson, C.: Automatic online news monitoring and classification for syndromic surveillance. *Decision Support Systems*. 47(4), 508–517 (2009)