# SOCIAL DATA VISUALIZATION SYSTEM FOR UNDERSTANDING DIFFUSION PATTERNS ON TWITTER: A CASE STUDY ON KOREAN ENTERPRISES

Dosam Hwang

*Department of Computer Engineering, Yeungnam University*
*Gyeongsan, Korea 712-749*
*e-mail:* `dshwang@yu.ac.kr`


Jai E. Jung*

*Department of Computer Engineering, Chung-Ang University*
*Seoul, Korea*
*&*
*Department of Information System, Universiti Malaya*
*Kuala Lumpur, Malaysia*
*e-mail:* `j2jung@gmail.com`


Seungbo Park

*Institute of Media Content, Dankook University*
*Yongin, Korea, 448-701*
*e-mail:* `molaal@naver.com`


Hien T. Nguyen

*Faculty of Information Technology, Ton Duc Thang University*
*Ho Chi Minh City, Vietnam*
*e-mail:* `hien@tdt.edu.vn`

**Abstract.** Online social media have been playing an important role of creating and diffusing information to many users. It means the users can get cognitive influence to the other users. Thus, it is important to understand how the information can be diffused by interactions among users through online social media. In this paper, we design a social media monitoring system (called "TweetPulse") which can analyze and show meaningful diffusion patterns (DP) among the users. Particularly, Tweet-Pulse focuses on visualizing information diffusion in Twitter, given a certain time duration. Also, this work has investigated the relationships 1) between DP and event detecting, 2) between DP and emotional words, and 3) between DP and the number of followers of the users. Thereby, to understand the continuous patterns of the information diffusion, we propose two different types of analytic methods, which are 1) macroscopic approach and 2) microscopic approach. For evaluating the proposed method, we have collected and preprocessed the dataset during about 4 months (14 March 2012 to 12 July 2012). As a conclusion, TweetPulse has helped users to easily understand DP from a large scale dataset streaming through Twitter.

**Keywords:** Information visualization, diffusion patterns, marketing strategies, Twitter

## 1 INTRODUCTION

Online social media (e.g., MySpace, Twitter, and Facebook) have been regarded as an important channel where information can be efficiently diffused to other users. This information diffusion can be accelerated (in opposite, decelerated), according to a number of situations. For example, the information about tsunami and earthquakes has been significantly distributed to the people in the certain areas. Thus, information diffusion is an important social phenomenon that we have to recognize and understand.

These social networking systems have supported collaborative tagging service and RSS-based technologies for users, and also they have become the factors which give important influences on diffusing the information through online social system [1, 2, 20]. For example, in Twitter, retweet (in short, RT) function has the effect of another information diffusion which is related to own opinion; also, hash tagging can attract more attention for creating and pulling attention from others [3, 4].

In fact, we have already received a lot of information from traditional media, e.g., television or radio. However, with the advent of new communication technology and internet development, users are able to collect new information more easily. Unlike the existing one-sided mass media, the current social media can support users to generate information and share with each other more widely [5, 21, 22].

---

* Corresponding author

Users can communicate with friends via online social network services like Twitter and Facebook. Moreover, they can exchange various opinions with strangers, and receive news from celebrity [6]. It is no longer a surprising social process that social network services bring news quicker than the traditional media. Actually, the subway in Seoul had been stopped on the rush time because of accident. The people who were in the subway uploaded the news on SNS and these tweets were RT to other users; this prevented rush time delays. There are ever more people who are contacting on SNS with the commuting time.

## 1.1 Backgrounds on Online Social Networking Services

Social networking service (SNS) is an efficient tool for sharing tremendous amount of information and disseminating information quickly [7]. In particular, Twitter can only represent one-side social relationships (called following and follower), while other online social networking services can represent two-side social relationships. In order to establish online relationships, it does not require the consent of the opponents. It means that social relationships in Twitter can be formed more easily than in other SNS. Next, Twitter supplies Twitter API to access a variety of information, and developers can easily implement third party applications by collecting from Twitter. Also, Twitter can supply re-transmission feature of the tweets, called retweet (RT). If a user creates a tweet, it will be delivered to the followers of the user, and the followers who received tweet can also retweet the same tweets to their followers for various purposes. In this case, the tweet can deliver the information like word-of-mouth process (from one user to his/her followers, and again, from his/her followers to their followers). Because there is no limit of the number of retweeting, more and more people can get the tweet. In terms of usability, it is convenient to retweet to their followers by clicking a button one time for rapid information diffusion on Twitter [8, 9, 10]. Thus, in this study, we study online information diffusion on Twitter by finding the pattern of the spread of information from Twitter and also the causes of this pattern. For this purpose, we have developed a software system (called TweetPulse) to show the relationships of diffusion of information with emotional language on Twitter.

## 1.2 Outline of This Paper

The outline of this paper is as follows. In Section 2, we introduce the related work of SNS diffusion and trend, and describe the information diffusion on the SNS. Section 3 discusses TweetPulse and its diffusion speed and diffusion convergence rate. Section 4 addresses the detecting events, and how long is the diffusion continued. In Section 5, we present the simulation to show the impact of events in the information diffusion and what information is propagated speedily. Finally, the conclusion of this work is given in Section 6.

## 2 RELATED WORK

### 2.1 Information Diffusion and Trend Analysis with SNS

SNS is a service to support users to build a human network on the web. As society advances, it is urged to express personal information and feelings grows stronger. Thus, the way of building social relationships changed to make friends on the Internet and websites like SNS is able to develop. Now more and more users use SNS to support new information. Because the SNS's information speed is very high, concise, and accurate, it is very advantageous for people in modern society [11, 12]. Unlike the traditional media's one-sided nature the SNS is based on a variety of digital media and thus it enables the contents, consumption and production by the user; thus, it opens a new chapter of society and culture. SNS can provide the opportunity for the formation of new connections through a variety of ways to communicate, various discussions on social issues on SNS, and leads the citizens to a diversity of social participation. SNS has characteristics such as speed, personality, information openness and easiness and can be classified into eight types depending on the capability of providing, such as profile-based, business-based, blog-based, vertical, collaborative, communication-centric, topic based, and micro blogging [13, 18, 23].

### 2.2 Information Diffusion on the SNS

Just a few years ago, we find the information necessary for academic or business on portal site. However, we are able to quickly and easily share various kinds of information created in real time through SNS with several members who form communities [16, 14]. With SNS explosive diffusion, information sharing activity among SNS users is increasing [15]. With the development of the SNS, the website shows changes of information diffusion on the SNS. Especially it has become popular, as shown in Figure 1 which shows politician interest index, the preference, keyword changes like TrueStory.

Also, the Pulse-K is a social media monitoring and analysis service developed by Conan technology. As shown in Figure 2, Pulse-K[1] provides media and emotional on the strip real-time which keywords to monitor, and users are able to see how to change the reaction of the people on those keywords.

Figure 3 shows a social search tool which can provide original search results from Twitter and Blog; it also can check the maximum diffusion message. Social metrics insight[2] provides crisis management and issues detected by services through sensitivity word analysis to help companies. SNS such as Twitter or Blog help also companies to communicate with consumers to promote new products.

---

[1] `http://www.pulsek.com/`
[2] `http://insight.some.co.kr/`

Figure 1. Snapshot of diffusion of politicians related keywords in True Story

## 3 MODELING SOCIAL PULSE

The information can be collected from Twitter diffusion. The aggregated information can be visualized on the chart so that users can easily understand the diffusion patterns (increasing or decreasing) on TweetPulse. Once a set of information of each user is projected from Twitter, it can be easily interpreted as a discrete signal which is a sequence of events over time. More importantly, we can set a time



Figure 2. Snapshot of Pulse-K services

Figure 3. Snapshot of social metrics insight

window (of size $\delta$), since we need to understand how the information is diffused in Twitter [17, 19].

Assume that $t \in T$, a TweetPulse is composed of two parts; a time window $w_{\tau_t} = [\tau_t, \tau_t + \delta]$, and a pitch (a height of the pulse). Especially, the pitch can be explained as either the number of users who have applied the same information or the number of information. Thus, it can be formalized by

$$P_t = \{\langle w_{\tau_t}, \pi_{\tau_t}\rangle | \tau_t \in \tau\} \tag{1}$$

where $\pi_{\tau_t}$ is the pitch of a time window and can be computed by

$$\pi_{\tau_t} = \left| \{r_k | u_i \times t_j \times r_k \times \tau_{t_j} \in F_\tau, t_j = t, \tau_{t_j} \in [\tau_t, \tau_t + \delta]\} \right| \tag{2}$$

where $u_i \in U$ and $r_k \in R$. More interesting work is to find relationships between given particular information. For example, if two tags are supported by more resources at the same time, we determine two sets of information have greater correlation than others.

### 3.1 Diffusion Speed

We can comprehend the speed of information diffusion, i.e., how quickly a tweet is propagated to other users in certain time duration. Since a tweet $t$ has been firstly

used at $\tau_t^0$, we can measure the pitch of the TweetPulse $\pi_{\tau_t}$. Given a tweet $t$ and its TweetPulse $P_t$, the speed of this diffusion can be measured by

$$S_t = \max_{P_t} \frac{\pi_{\tau_t}}{\delta} \tag{3}$$

where $\delta$ indicates a size of the time window.

The speed simply indicates the maximum diffusion power of the corresponding tweet. We can measure the speed, when the TweetPulse of the tag shows the highest pitch. In other words, if we can construct a cumulative curve of the TweetPulse, we can easily find the speed at the steepest slope.

## 3.2 Convergence Rate of Diffusion

We can also measure the convergence rate of the propagated tweet, i.e., how quickly the tweet is spread to most of users. Given a tag $t$ and its TweetPulse $P_t$, the convergence rate of its propagation can be measured by a temporal duration

$$C_t = |\tau_t^\Omega - \tau_t^0| \tag{4}$$

where $\tau_t^\Omega$ is the time ending the TweetPulse (i.e., $\pi_{\tau_t^\Omega} = 0$). It means that after $\tau_t^\Omega$, there is no meaningful TweetPulse. Similar to the diffusion speed, convergence rate is also an important indicator for measuring diffusion power.

## 4 DETECTING AND ANALYZING EVENTS

### 4.1 Event Detection

If at any time a large rate of diffusion is shown on the chart, maybe this is due to any key causes (social issues or events), and we will call it events detecting. Figure 4 is a graph showing the comparison of TweetPulse for change of counts searched by two keywords (Politicians) from 12 April 2012 to 29 April 2012.

We can see in some time-zones that diffusion speed is faster than in other time-zones for each keyword, and we are able to surmise at the time-zones which showed the highest diffusion speed that the time-zones had some social issue or events. On 12 April 2012, relationship with a certain politician "$P_1$" has social issue with tweets, e.g., "$tweet_1$", and the tweet has influenced diffusion speed of the keyword of $P_1$. On 29 April 2012, the social event which "$tweet_2$" has influenced the diffusion speed of the keyword of another politician "$P_2$". As this views a big diffusion speed at the time-zones, it will probably have events detecting about the keyword.

Also, in Figure 5, there is a time-zone when the diffusion speed is maximum on 12 April 2012. It means at that time there were some events relationships with the keyword of 'movie1'.

Thus, Figure 6 shows the snapshot before the time-zone when the tweet was propagated at maximum diffusion speed. An actress appears in movie "$M_1$" and
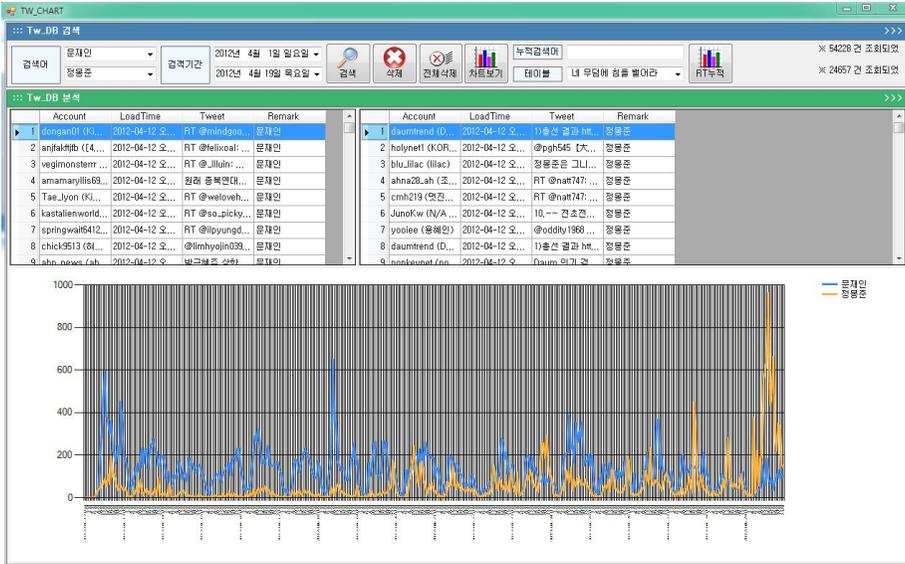
Figure 4. Detecting events on TweetPulse by politicians

writes "$tweet_3$" on the Twitter. This tweet includes the keyword of "$M_1$", and the users who follow the actress RT this message, so on 12 April 2012 the keyword of "$M_1$" has propagated fast.

## 4.2 Continuous Time of the Information Diffusion

Figure 7 shows highest diffusion speed in some time-zones. At 9 AM on 16 April, there is a high diffusion speed compared with other time-zones of keyword "$P_2$" related to a certain politician. Also, at 11 AM on 21 April, a keyword of another politician "$P_3$" has a higher diffusion speed than other time-zones.

In this paper, we have also studied diffusion continuance time when an event has happened to one keyword. According to how long a tweet has been retweeted, we can check it using two types of analytics:

1. macroscopic and
2. microscopic analytic methods.

### 4.2.1 Macroscopic Analytic Method

First we can check it on the monitor of events detecting on TweetPulse. We can call it *macroscopic* analytic method. For example, in Figure 7, we can see that "Ahn" keyword at the time-zone of 16 April 2012 (9 AM) has an event detected, also "Kim" keywords have event detected at 21 April 2012 (1 AM) and the diffusion
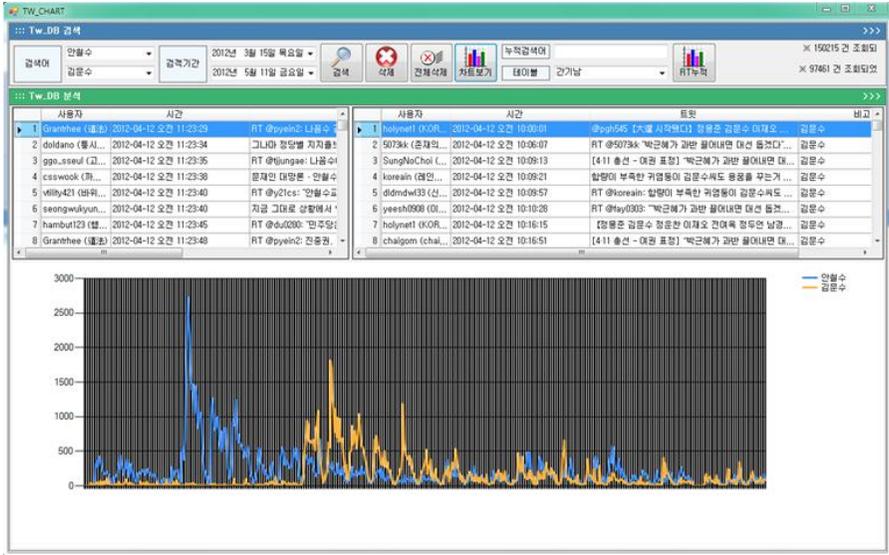
Figure 5. Detecting events on TweetPulse by movie title



Figure 6. DB data of "$M_1$"

phenomenon lasts for 3 to 5 days; but maybe there are other events at the same time. For example, one event happened on 21 April 2012, and it just continued for a day, and another event happened on the next day.

### 4.2.2 Microscopic Analytic Method

To discover how long the diffusion phenomenon holds details, we have to use the microscopic analytic method. By help of this method, we will check the continuance time of the information diffusion. Thus, we can see it on Tweet cumulative analysis program. There will be many tweets related with a keyword to be uploaded every day. We study the emergence of some important events, so we just select the ones with over 80 tweets propagated.

As shown in Figure 8, on 16 April, the important event which has caused significant diffusion is "$tweet_4$". Compared in Figure 8, we cannot see the "$tweet_4$" for

Figure 7. Detecting events on TweetPulse by politicians

another two days except the $16^{th}$. However, we can see the tweet whose "$tweet_4$" lasted for 3 days.

Next we want to compare the tweets of "Kim". As shown in Figure 9, on 21 April, the important event which caused tweet diffusion is "$tweet_5$". Figure 9 shows that "$tweet_5$" has lasted for 2 days, and on 23 April, the tweet "comment" has been diffused more quickly than other tweets. It means not only one event happened, but also several events happened together. Also, most of the events have been just diffused during 1 to 2 days, while only a few events have lasted for 3 to 5 days.

## 5 EXPERIMENTATION

### 5.1 System Architecture

To understand the TweetPulse, the simulation system is made up of two software modules:

1. data collection function and
2. analysis tool for understanding diffusion phenomenon.

The analysis tool for understanding diffusion phenomenon is also composed of two parts; the first one is Tweet counter analysis program and the second is Tweet cumulative analysis program.

As shown in Figure 10, the data collecting program divides the table for each keyword to collect data in the database by Tweet API. Tweet counter analysis

Figure 8. Tweet cumulative analysis of "Ahn" on a) 16 April, b) 17 April, and c) 18 April

Figure 9. Tweet cumulative analysis of "Kim" on a) 21 April, b) 22 April, and c) 23 April
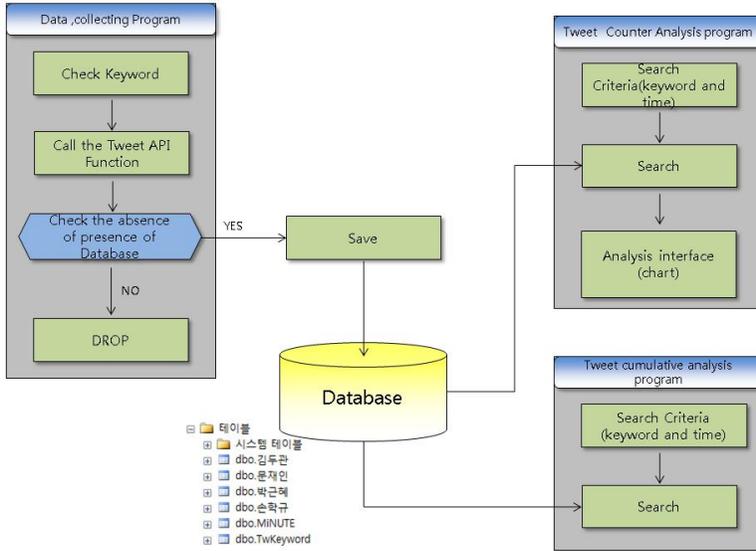
Figure 10. System architecture

program shows the tweet counter on the time during which the tweet information was collected in the database. Tweet cumulative analysis program shows the tweet information which was collected in the database by accumulating during one day.

## 5.2 Data Collection

For experimentation, we have targeted 8 politicians and 16 movie titles for collecting datasets and the keywords. In this study, we have collected data during 4 months (14 March 2012 to 12 July 2012). Thus, we have collected about 4 000 000 tweets. With this, we have constructed a vast database, and were able to search for 4 months data to cover the flaw which search tweet just for 7 days in Twitter.

## 5.3 Experimental Environment

In this paper, we have collected data for 4 months (14 March 2012 to 12 July 2012). We have developed a monitoring program, as shown in Figure 11.

This system is developed by using C# and MS-SQL. It can save the tweets which were searched using 8 keywords at special interval simultaneously. The most superior point in this system is that it can be run and collect data in PC without using the web page. The advantages of this program are in the ability to speculate what event has happened at the time-zones when a large diffusion speed is shown. Also, we can guess at the time-zone when there is a large diffusion speed that some events happened. More importantly, tweet cumulative analysis program is a monitor
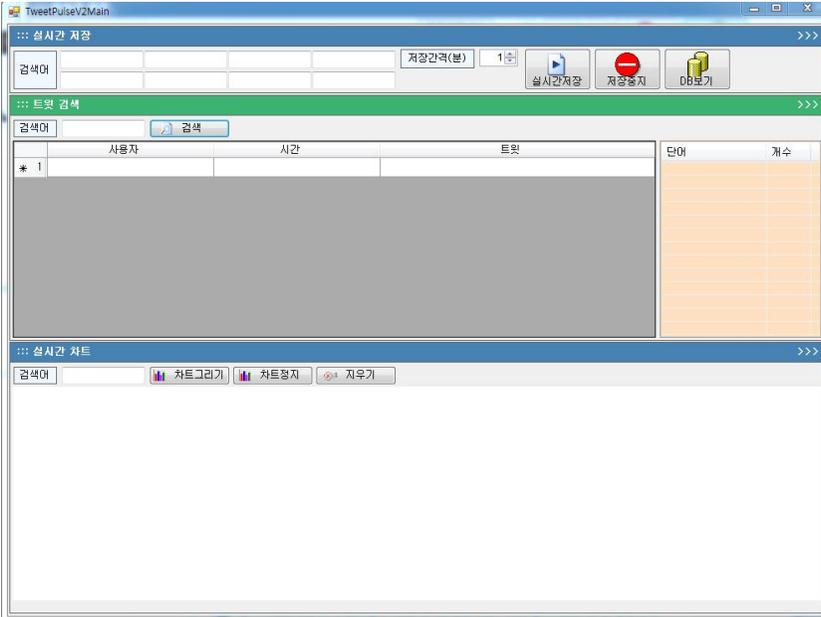
Figure 11. Snapshot on GUI of Tweet collection software

program made to study the relationship of diffusion phenomena and perceptual language of tweets. The program is also used for microscopic analysis continuance time of the information diffusion.

Data cleansing such as de-duplication is an important issue for understanding information diffusion patterns. But, in these systems, the meaning of the diffusion of information is not accurate. Tweet search on 100 tweets per second is used in the data collecting program. Every 1 second returns 100 tweets, thus there will be collect duplication. Due to this problem, the process of duplication removal is necessary. Given a data which want to save $t$, the data can be formalized by

$$t = \langle RT, U, T, S, sURL \rangle. \tag{5}$$

As shown in the Figure 12, before saving the data in DB which collect 100 data every 1 second to check user, upload-time, and tweet in DB, when $(U_1 = U_2) \wedge (S_1 = S_2) \wedge (T_1 = T_2)$, in other words $t_1 = t_2$, it means the data is already saved in database, so it can be dropped. In opposite, if $t_1 \neq t_2$, the data will be saved in DB.

Figure 13 is a chart which cumulates tweets for more than 30 days. The tweets which view significant diffusion speed include emotional words. The tweets which include emotional words diffused faster and more abundantly.
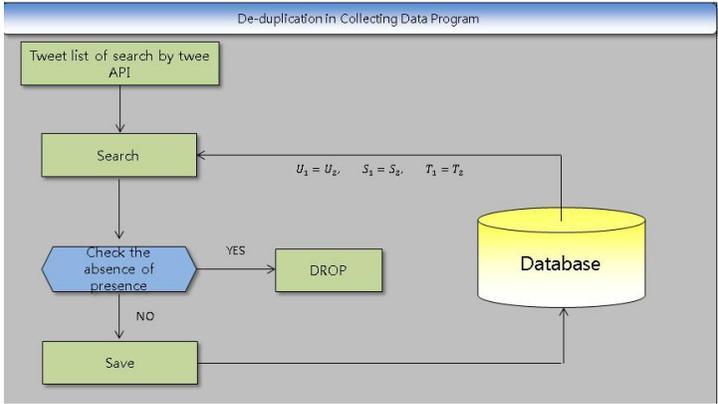
Figure 12. System architecture of the duplication removal process
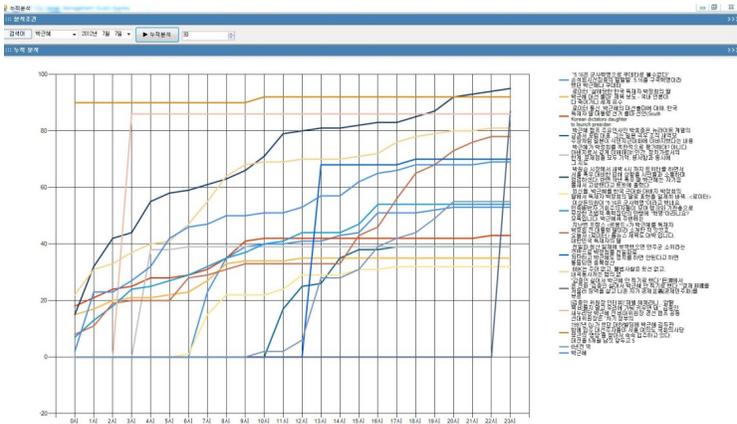


Figure 13. Tweet cumulative analysis program

## 6 CONCLUDING REMARKS AND FUTURE WORK

In this paper, we designed and implemented a model which can analyze the diffusion of information efficiently by finding the TweetPulse, and showed the rate of diffusion using visual interface. First of all, the relationship between information diffusion and some events keyword were analyzed, and then the number of followers of the user who write a tweet and its effect on information diffusion. In our next researches, we shall try to find what feelings are added when a tweet is diffused.

**Acknowledgement**

**REFERENCES**

[1] YANG, J.—COUNTS, S.: Comparing Information Diffusion Structure in Weblogs and Microblogs. In: W. W. Cohen, S. Gosling (Eds.): Proceedings of the 4th International Conference on Weblogs and Social Media (ICWSM 2010), May 23–26, AAAI Press, Washington, DC, USA 2010.

[2] PHAM, X. H.—JUNG, J. J.—HWANG, D.: Beating Social Pulse: Understanding Information Propagation via Online Social Tagging Systems. Journal of Universal Computer Science, Vol. 18, 2012, No. 8, pp. 1022–1031.

[3] CHA, M.—MISLOVE, A.—GUMMADI, K. P.: A Measurement-Driven Analysis of Information Propagation in the Flickr Social Network. In: J. Quemada, G. Léon, Y. S. Maarek, W. Nejdl (Eds.): Proceedings of the 18th International Conference on World Wide Web (WWW 2009), Madrid, Spain, April 20–24, ACM 2009, pp. 721–730.

[4] JUNG, J. J.: CONTEXTGRID: A Contextual Mashup-Based Collaborative Browsing System. Information Systems Frontiers, Vol. 14, 2012, No. 4, pp. 953–961.

[5] LENHART, A.: Adults and Social Network Websites.

[6] HU, M.—LIU, B.: Mining and Summarizing Customer Reviews. In: W. Kim, R. Kohavi, J. Gehrke, W. DuMouchel (Eds.): Proceedings of the 10th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, August 22–25, ACM, Seattle, Washington, USA 2004, pp. 168–177.

[7] JUNG, J. J.: Attribute Selection-Based Recommendation Framework for Short-Head User Group: An Empirical Study by Movielens and IMDb. Expert Systems with Applications, Vol. 39, 2012, No. 4, pp. 4049–4054.

[8] SIGANOS, G.—FALOUTSOS, M.—FALOUTSOS, P.—FALOUTSOS, C.: Power Laws and the As-Level Internet Topology. IEEE/ACM Transactions on Networking, Vol. 11, 2003, No. 4, pp. 514–524.

[9] SCHNEIDER, A.—JACKSON, R.—BAUM, N.: Social Media Networking: Facebook and Twitter. The Journal of Medical Practice Management, Vol. 26, 2010, No. 3, pp. 156–157.

[10] JUNG, J. J.: Evolutionary Approach for Semantic-Based Query Sampling in Large-Scale Information Sources. Information Sciences, Vol. 182, 2012, pp. 30–39.

[11] BOYD, D.—GOLDER, S.—LOTAN, G.: Tweet, Tweet, Retweet: Conversational Aspects of Retweeting on Twitter. In: Proceedings of the 43rd Hawaii International Conference on Systems Science (HICSS-43), IEEE Computer Society, Koloa, Kauai, HI, USA 2010, pp. 1–10.

[12] KRISHNAMURTHY, B.—WILLINGER, W.—GILL, P.—ARLITT, M. F.: A Socratic Method for Validation of Measurement-Based Networking Research. Computer Communications, Vol. 34, 2011, No. 1, pp. 43–53.

[13] VAITHYANATHAN, S.: The Value of Social Media Data in Enterprise Applications. In: K. S. Candan, Y. Chen, R. T. Snodgrass, L. Gravano, A. Fuxman (Eds.): Proceedings of the ACM SIGMOD International Conference on Management of Data (SIGMOD 2012), May 20–24, ACM, Scottsdale, AZ, USA 2012, pp. 755–756.

[14] JUNG, K.—HEO, W.—CHEN, W.: IRIE: Scalable and Robust Influence Maximization in Social Networks. In: Proceedings of the 12th IEEE International Conference on Data Mining (ICDM 2012), Brussels, Belgium 2012.

[15] JUNG, J. J.: Semantic Optimization of Query Transformation in a Large-Scale Peer-to-Peer Network. Neurocomputing, Vol. 88, 2012, pp. 36–41.

[16] VALLER, N.—PRAKASH, B. A.—TONG, H.—FALOUTSOS, M.—FALOUTSOS, C.: Epidemic Spread in Mobile Ad Hoc Networks: Determining the Tipping Point. In: J. Domingo-Pascual, P. Manzoni, S. Palazzo, A. Pont, C. M. Scoglio (Eds.): Proceedings of the 10th International IFIP TC 6 Networking Conference (NETWORKING 2011), May 9–13, Vol. 6640 of Lecture Notes in Computer Science, Springer 2011, pp. 266–280.

[17] JUNG, J. J.: Boosting Social Collaborations Based on Contextual Synchronization: An Empirical Study. Expert Systems with Applications, Vol. 38, 2011, No. 5, pp. 4809–4815.

[18] ELLISON, N. B.—STEINFIELD, C.—LAMPE, C.: The Benefits of Facebook "Friends": Social Capital and College Students' Use of Online Social Network Sites. Journal of Computer-Mediated Communication, Vol. 12, 2007, No. 4, pp. 1143–1168.

[19] JUNG, J. J.: Discovering Community of Lingual Practice for Matching Multilingual Tags from Folksonomies. Computer Journal, Vol. 55, 2012, No. 3, pp. 337–346.

[20] VOSSEN, G.: Big Data as the New Enabler in Business and Other Intelligence. Vietnam Journal of Computer Science, Vol. 1, 2014, No. 1, pp. 3–14.

[21] PHAM, X. H.—NGUYEN, T. T.—JUNG, J. J.—NGUYEN, N. T.: ⟨A, V⟩-SPEAR: A New Method for Expert Based Recommendation Systems. Cybernetics and Systems, Vol. 45, 2014, No. 2, pp. 165–179.

[22] NGUYEN, D. T.—JUNG, J. J.: Privacy-Preserving Discovery of Topic-Based Events from Social Sensor Signals: An Experimental Study on Twitter. Scientific World Journal, Vol. 2014, 2014, Article ID 204785.

[23] JUNG, J. J.: Semantic Wiki-Based Knowledge Management System by Interleaving Ontology Mapping Tool. International Journal on Software Engineering and Knowledge Engineering, Vol. 23, 2013, No. 1, pp. 51–63.
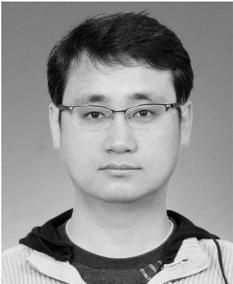
**Dosam HWANG** received his Ph. D. degree in natural language processing at Kyoto University. In 1987, he was chosen the best researcher at Korea Institute of Science and Technology. Since 2005 he has been the Professor at Yeungnam University in Korea. Currently he is working on natural language processing, machine translation, ontology, semantic web, and information retrieval. He has served for a number of international conferences and a technical committee for ISO.

**Jai E. Jung** is a corresponding author of this paper. He is an Associate Professor in Chung-Ang University, Korea, since September 2014. Before joining Chung-Ang University, he was a faculty member in Yeungnam University from 2007. He was a postdoctoral researcher in INRIA Rhône-Alpes, France in 2006, and a visiting scientist in Fraunhofer Institut (FIRST) in Berlin, Germany in 2004. He received the B. Eng. in computer science and mechanical engineering from Inha University in 1999. He received M. Sc. and Ph. D. degrees in computer and information engineering from Inha University in 2002 and 2005, respectively. His research topics include knowledge engineering on social networks by using machine learning, semantic Web mining, and ambient intelligence.

**Seungbo Park** is a research professor in Kyung Hee University, Korea. He received the B. Sc. and the M. Sc. degrees in Electrical Engineering from Inha University, Korea, in 1995 and 1997. He received the Ph. D. degree in Computer Information Engineering from Inha University, Korea. His research interests include video story analyzing, semantic contents, video knowledge representation, social network analysis, and AI. He worked at Daewoo Electronics as an engineering researcher.

**Hien T. Nguyen** has served as the Dean, Faculty of Information Technology, Ton Duc Thang University, Vietnam, since 2012. He received Ph. D. in computer science, M. Sc. in computer science, and B. Eng. in computer engineering from Ho Chi Minh City University of Technology, Vietnam, in 2011, 2005, and 2002, respectively. His research interests include job shop scheduling, information extraction, natural language processing, web/text mining, semantic web, and social network analysis and mining.