

1 **SOCIAL MEDIA DATA ANALYSIS FOR TRAFFIC INCIDENT DETECTION AND**  
2 **MANAGEMENT**

3  
4  
5 **Kaiqun Fu**

6 Department of Computer Science  
7 Virginia Polytechnic Institute and State University  
8 7054 Haycock Rd, Falls Church, VA 22043  
9 Tel: 703-538-8310; Email: [fukaiqun@vt.edu](mailto:fukaiqun@vt.edu)

10  
11 **Rakesh Nune, Corresponding Author**

12 District of Columbia Department of Transportation  
13 55 M Street, SE, Washington DC20003  
14 Tel: 571-265-4047; Email: [Rakesh.Nune@dc.gov](mailto:Rakesh.Nune@dc.gov)

15  
16 **Jason X. Tao**

17 District of Columbia Department of Transportation  
18 55 M Street, SE, Washington DC20003  
19 Tel: 202-671-1489; Email: [jason.tao@dc.gov](mailto:jason.tao@dc.gov)

20  
21  
22 Word count: 2,527 word text + 6 tables/figures x 250 words (each) = 4,027  
23  
24  
25  
26  
27  
28

29 July 31, 2014

1

2 **ABSTRACT**

3 Social media and online services with user-posted content have generated a staggering amount of  
4 information that has potential applications in various areas such as sentiment analysis and  
5 emergency management. In this research, we explored the feasibility of using social media data for  
6 detecting traffic incidents and collecting supplemental incident information. A comprehensive  
7 approach has been developed to extract and analyze real-time traffic related twitter data for  
8 incident management purpose. The developed approach consists of three steps: (1) Development  
9 of traffic incident related key words and their association rules; (2) Extraction of real-time tweets  
10 with influential word sets; and (3) ranking and selection of traffic related tweets. The developed  
11 approach was implemented at District of Columbia Department of Transportation for incident  
12 management. Data validation has been conducted against the real-world incident database. The  
13 preliminary results of analysis have shown that social media data is promising for early incident  
14 detection and can be used as supplemental source for incident data collection.

15

16

17

18 *Keywords:* Social media analysis, Incident Detection, Traffic Incident Management, Tweet  
19 extraction, Term frequency

20

## 1 INTRODUCTION

2 Early detection of traffic incidents is critical to reduce the impact of incidents on the traffic  
3 conditions. There are multiple incident detection sources which include reports of roadway  
4 operation patrollers, citizen calls, CCTV monitoring and alerts from vehicle detection stations. In  
5 the District of Columbia, the CCTV system is the most important means to detect and verify traffic  
6 incidents. Since the current CCTV system does not provide 100% coverage in the District, delay in  
7 incident detection may take place especially in those areas not covered by the CCTV system.  
8 Recently, we conducted a pilot project to explore the application of social media data analysis in  
9 incident detection and traffic management.

10

11 In the past decade, as social media (e.g., Twitter, Instagram and Facebook) became more and more  
12 popular, social media data has been collected and used in various applications. For instance,  
13 Bollen *et al.* (1) performed sentiment analysis on large scale Twitter data to predict the daily  
14 directions of stock market. Schmidt (2) explored the social media data to predict and track disease  
15 outbreaks. Sabra (3) explored the benefits of integrating social media tools into emergency  
16 communications and monitoring the social media contents during an emergency or disaster. Schulz  
17 *et al.* (4) developed an approach for efficient processing and storing of social media data for  
18 emergency management purpose. Applying social media analysis in traffic management is a  
19 relatively new concept. Today, many government organizations and news companies are using  
20 social media such as Twitter to communicate with the public and commuters on traffic incidents.  
21 On the other hand, many commuters report incidents, special events or congestion levels that on  
22 route through the social media. Therefore, the contents of popular social media contain abundant  
23 information on traffic incidents and roadway congestion. It is feasible to extract and analyze the  
24 social media data for detecting traffic incidents and obtaining supplemental incident information.

25

26 In this pilot project, we developed a comprehensive approach to extract and analyze real-time  
27 traffic related twitter data for incident management purpose. The developed approach consists of  
28 three steps: First, a set of keywords were extracted from historical traffic related tweets generated  
29 by influential users, then a comprehensive algorithm is applied to establish the association rules of  
30 keywords; Secondly, real-time tweets are extracted through Twitter API with influential word sets;  
31 At the last step, ranking and selection of the extracted tweets are conducted to capture the most  
32 relevant incident tweets. A detailed methodology is presented in the next section, followed by a  
33 description of implementation of the developed approach and data validation. Conclusions are  
34 drawn in the last section.

35

## 36 METHODOLOGY

37 In order to successfully extract traffic incident related tweets, it is necessary to establish a set of  
38 keywords that will be sent over twitter to query data through the twitter timeline API [5]. Through  
39 preliminary analysis, we have identified four influential Twitter users who actively post traffic  
40 information. These four account names are “WTOPtraffic”, “VaDOT”, “drgridlock” and  
41 “DCPoliceDep”, which are held by WTOP, Virginia DOT, The Washington Post and District of  
42 Columbia Police Department, respectively. The tweets from these four users are used to develop

1 the set of keywords. In this research, a data collector is built to collect 3200 recent tweets posted by  
2 four influential users that compose a document set, denoted as  $T_1$ .

3  
4 A statistic method called “term frequency–inverse document frequency” (or simply as *tf-idf*) (6) is  
5 applied to document set  $T_1$  to determine the importance of each word in the set. *tf-idf* is the product  
6 of two statistics, term frequency and inverse document frequency. For word  $t$  and document  $d$ , the  
7 term frequency  $tf(t,d)$  is defined as:

$$8 \quad tf(t, d) = \frac{f(t,d)}{\max\{f(w,d): w \in d\}} \quad (1)$$

10 where  $f(t, d)$  is the number of times that term  $t$  occurs in document  $d$ .

11  
12 The inverse document frequency, or *idf*, is a measure of how much information the word provides,  
13 that is, whether the term is common or rare across all documents. *idf* is defined as

$$14 \quad idf(t, D) = \log\left(\frac{N}{|\{d \in D : t \in d\}|}\right) \quad (2)$$

15  
16 where  $D$  is a corpus of all document,  $N$  is the number of documents in the corpus, and  
17  $|\{d \in D : t \in d\}|$  is the number of documents where word  $t$  appears. In our case, each tweet is a  
18 document, and all tweets together compose a corpus. Since all words to be analyzed are from the  
19 extracted tweets,  $|\{d \in D : t \in d\}|$  will never be zero. We choose a total of 50 words with highest *tf-idf*  
20 weights as the traffic incident related key words. These key words are shown in Table 1.  
21  
22

23  
24 **TABLE 1 List of Traffic Incident Related Keywords**

26 #dctrffic	#mdtraffic	#vattraffic	crash	due
delays	st	traffic	right	ave
27 lane	lanes	rd	nw	accident
buses	left	md	nb	bridge
28 sb	earlier	street	near	blocked
loop	congestion	expect	va	dc
29 update	road	work	following	closed
open	vehicle	inner	car	killed
30 new	get	minute	directions	close
31 schedule	police	beltway	operating	us

32  
33 Use  $Q$  to denote the set of all above keywords. The keywords list  $Q$  can be used as input to the  
34 Twitter Search API to crawl all the tweets that match the queries.

35

1 Through the preliminary analysis, it is found that if a single keyword such as “crash” or “lane” is  
2 used in a query, there are enormous noises in extracted tweets. In addition, the number of tweets  
3 from the data crawler will easily reach the assigned rate limit. In order to solve these issues, we  
4 applied the Apriori algorithm (7) to keyword list Q to develop the association rules of keywords  
5 and then construct Twitter queries using a combination of a few keywords (or called “wordset”)  
6 rather than single keyword.

7  
8 The Apriori Algorithm is usually applied in the field of transaction mining to establish frequent  
9 itemsets for boolean association rules. The algorithm proceeds by identifying the frequent  
10 individual items in the database and extending them to larger and larger item sets as long as those  
11 item sets appear sufficiently often in the database. In general, the concept of the Apriori Algorithm  
12 can be expressed in the following steps:

13 Step 1: Identity frequent itemsets that are the sets of item with minimum support (denoted  
14 by  $L_k$  for  $k$ th-Itemset);

15 Step 2: Apply the Apriori property: Any subset of frequent itemset must also be frequent.  
16 That means all subset of frequent itemset must have minimum support;

17 Step 3: Extend the length of itemset by the join operation: To find  $L_{k+1}$ , a set of candidate  
18  $(k+1)$ -itemsets is generated by joining  $L_k$  with itself.

19 In this case, each keyword is viewed as an item and thus a wordset is viewed as an itemset in the  
20 field of transaction mining. The minimum support level is set to 0.02 in this implementation. It is  
21 also found that when wordsets of 3 keywords are used, the algorithm is able to filter most noise in  
22 tweets and the results are optimized.

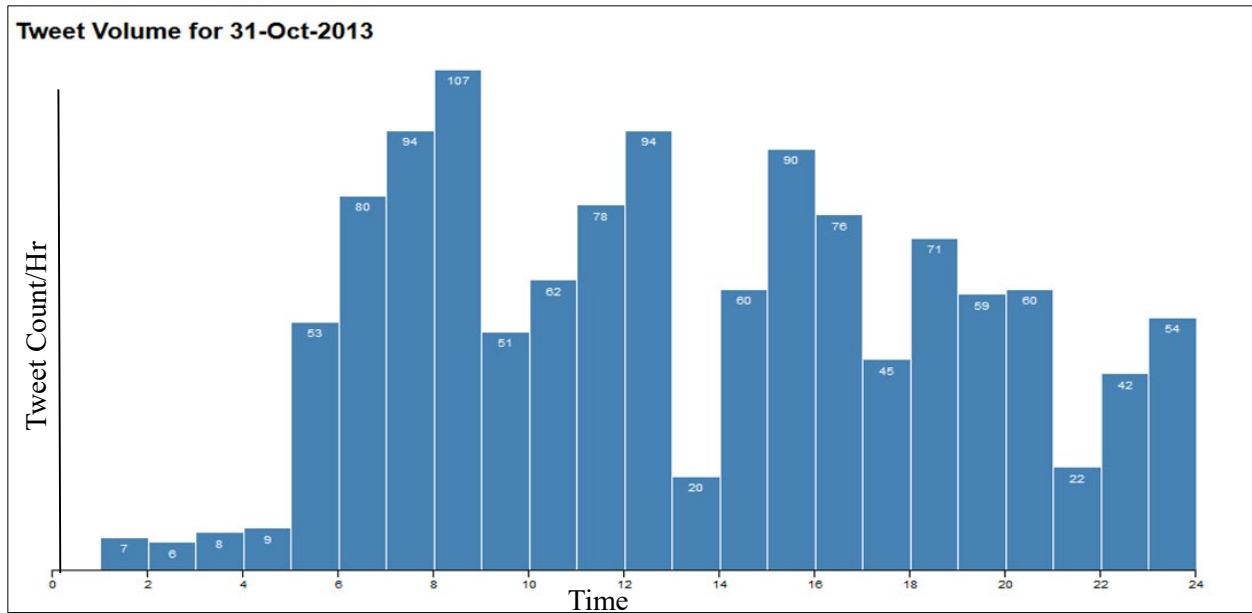
23  
24 In order to highlight the most relevant tweets from the extracted dataset, all tweets are ranked  
25 based on the tf-idf weights of the contained keywords. For each tweet, a score is calculated by  
26 adding the weights of all keywords appearing in the tweet. All extracted tweets are ranked by their  
27 scores and stored into the database.

## 28 29 **IMPLEMENTATION AND ANALYSIS**

30  
31 The above described approach for Twitter data extraction and analysis was implemented at District  
32 of Columbia Department of Transportation in late 2013. An intuitive user interface is built to  
33 display the latest traffic tweets along with visualizations tools like histogram. The main front end  
34 user interface was implemented using bootstrap (8) - a HTML5 and CSS3 based web applications  
35 framework and backend with Django (9). The implemented user interface allows the users to  
36 search for the incident or congestion related tweets on particular dates and times. This feature  
37 helps the operators at traffic management center to obtain supplemental incident information  
38 which may be missing from the existing incident profiles. Tweets volume histograms for a regular  
39 work day and a weekend day are displayed in Figures 1 and 2. The figures clearly shows that in a  
40 regular work day, more traffic related tweets are posted during morning and afternoon rush hours.  
41 In a weekend day, most traffic related tweets are posted in the evening time from 5:00PM to

1 8:00PM. This is consistent with the travel rates and incident rates over daily time.

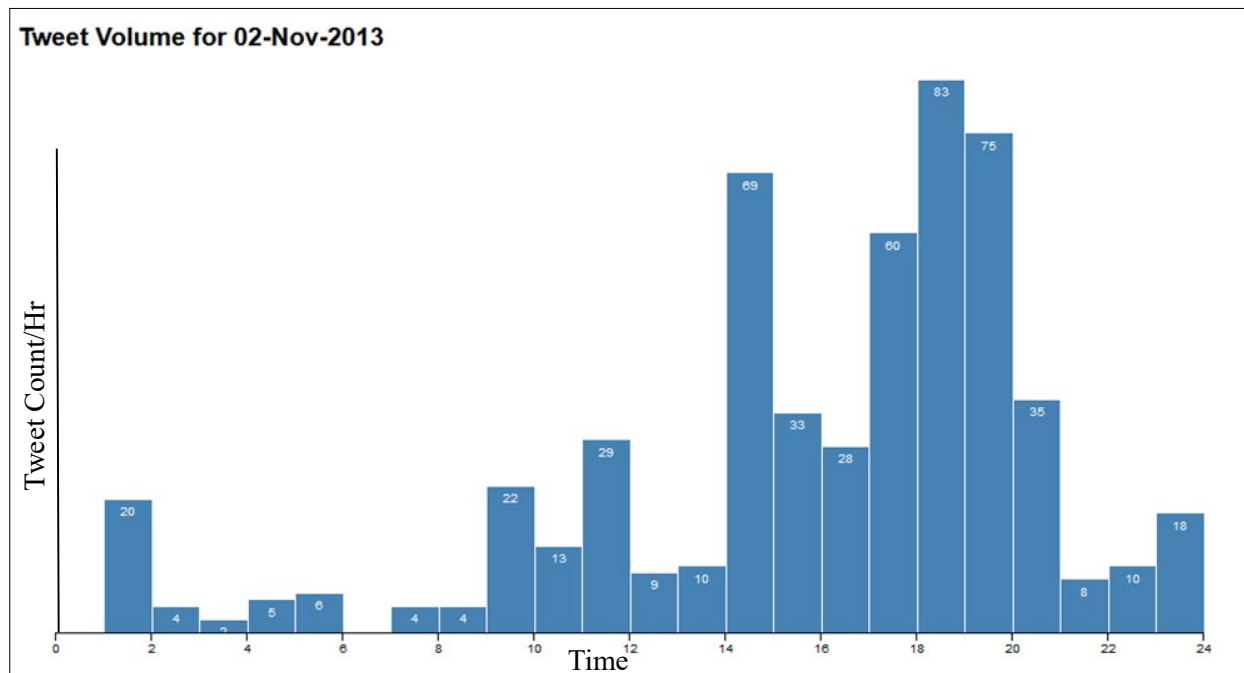
2



3

4 **FIGURE 1 Tweets Volume Histogram on a Regular Work Day**

5



6

7 **FIGURE 2 Tweets Volume Histogram on a Weekend Day**

8

9 The twitter data is unstructured and number of fields varies depending on the tweets unlike a  
10 predefined model such as relational databases. The number of variables may changes depending

1 on users' profiles and input. To deal with the scalability issue and storage requirements, NoSQL  
 2 database MongoDB [10] is used to store and process the extracted tweet. Figure 3 shows the  
 3 database schema used in NoSQL.

4

Key	Value	Type
▲ (4) 476575427474243584	{ 16 fields }	Object
▣ _id	476575427474243584	Int64
▣ lang	en	String
▣ favorited	false	Boolean
▣ pj_UTctime	1402459410	Int32
▷ (3) retweeted_status	{ 11 fields }	Object
▣ truncated	false	Boolean
▣ text	West / in-bound lanes on NY Ave @ N-Capitol are close...	String
▣ created_at	Wed Jun 11 04:03:30 +0000 2014	String
▷ (1) pj_scoreObj	Array [1]	Array
▣ retweeted	false	Boolean
▷ (1) user_mentions	Array [1]	Array
▣ source	<a href="https://about.twitter.com/products/tweetdeck"...	String
▷ (1) user	{ 22 fields }	Object
▷ (6) pj_kwlist	Array [6]	Array
▣ retweet_count	1	Int32
▣ id	476575427474243584	Int64

5

6 **FIGURE 3 The Database Schema in NoSQL**

7  
 8 Validation has been conducted by comparing the real-time tweet contents extracted by the  
 9 developed algorithm with the real-time incident data collected through other sources (e.g., CCTV  
 10 monitoring, roadway patroller reports and citizen's call). The results of analysis have shown that  
 11 the algorithm is able to capture almost all vehicle crashes or disabled vehicles on the roadways that  
 12 cause congestion. One of the important findings is that the Twitter data contains much more  
 13 numbers of incidents than that collected through the CapTOP system – a traffic management  
 14 platform used by the District Department of Transportation.

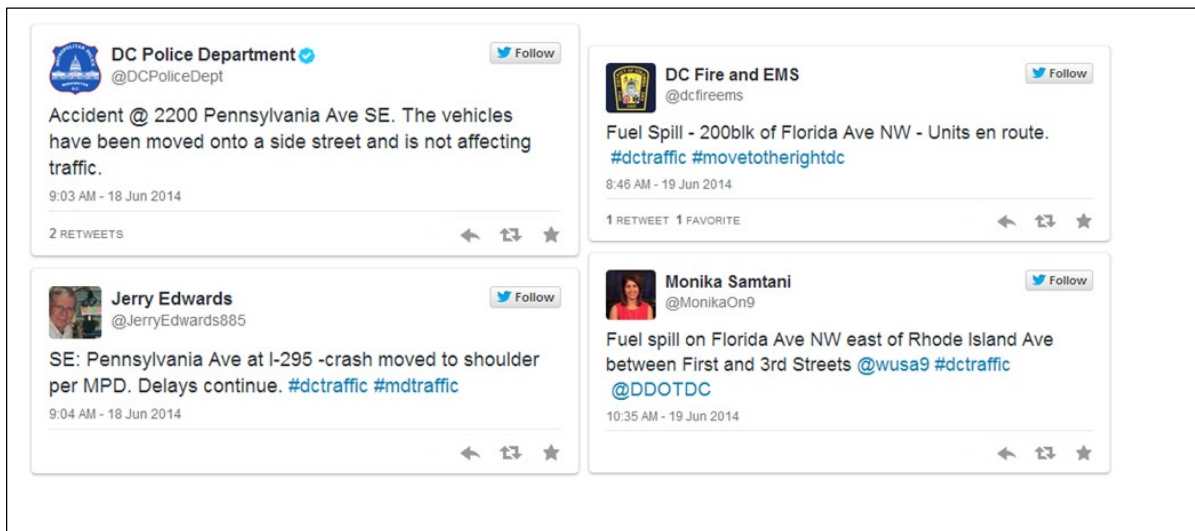
15  
 16 In order to illustrate the efficiency of the developed approach, we compared some of the detected  
 17 incidents that took place on June 8, 2014 and June 9, 2014 between the CapTOP database and the  
 18 extracted tweets. Table 2 lists four incident records in the CapTOP database that were reported by  
 19 roadway operation patrollers. These records include incident ID, detection date and time, incident  
 20 location information (Street names) and incident descriptions and operators comments for these  
 21 incidents. It can be seen that for most of the above records, the incident description and operator  
 22 comments do not provide sufficient information for these incidents.

23  
 24  
 25  
 26  
 27

1 **TABLE 2 Incident Records in the DDOT Incident Database**

Incident Id	Date	Time	Street Id (A)	Street Id (B)	Description	Comment
69673	6/18/2014	09:08	Branch Ave	Pennsylvania Ave	DC eg 3526	NULL
69724	6/19/2014	10:50	First ST	Florida Ave	NULL	Nothing found ROP 4 responded to incident.
69734	6/19/2014	16:59	Park Street and W Street	Branch Ave	NULL	NULL
69735	6/19/2014	17:10	Eastern Ave	I-295	NULL	Slow traffic condition nb 295 prior to crash in Maryland side. Camera out of order

2  
3 Figures 4(a) and 4(b) show the corresponding tweets for each incident listed in Table 2. From the  
4 Figure 4(a) and 4(b) one can see that the posted time for each tweet is significantly earlier than that  
5 recorded in the CapTOP database. More importantly, from these tweets users can capture much  
6 more detailed information about the incidents. Such supplemented information is not only  
7 traditionally textual, but also pictorial. For example, one of the tweet results for incident number  
8 69734, posted by DC Police Department, contains a detailed screen shot of Google Map with the  
9 incident location highlighted. Such visual and user friendly information presentation can enhance  
10 the DDOT incident database with the concept of *Big Data*.

27 **FIGURE 1a Traffic Incident Related Tweets Extracted from the Algorithm**

28

29

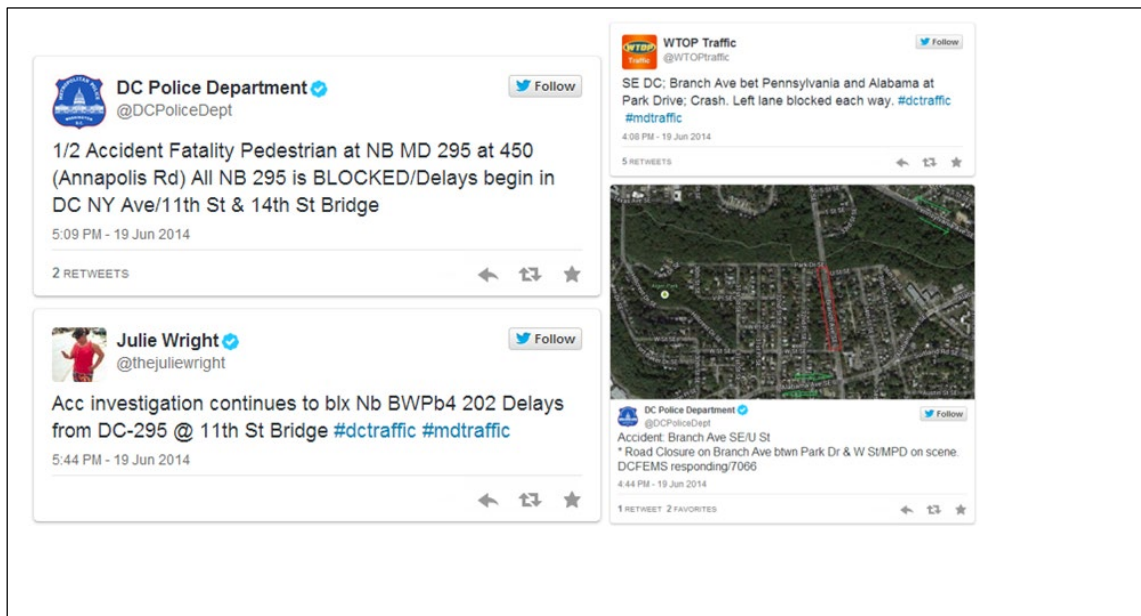
30

31

32

33





**FIGURE 2b Traffic Incident Related Tweets Extracted from the Algorithm**

## CONCLUSIONS AND FUTURE WORK

In this research, a comprehensive approach has been developed to extract and analyze the traffic incident related tweets for detecting traffic incidents and collecting supplemental incident information. The validation and analysis has showed the promising application of the approach in incident management.

The twitter data crawling and scoring process proposed in our algorithm utilizes *tf-idf* calculation to determine the weights for each keyword. The algorithm used currently is static, which means once the weight for each keyword is calculated; the weight will be assigned to the keyword for all following tweet score calculations. Although this method can provide better algorithm performance, it cannot model the actual transportation events. Since the appearance of transportation events are random, and more importantly independent to each other. In other words, knowledge from the previous events gives no inference to the next. Due to this characteristic of the transportation events, it is optimal to update the keywords weights dynamically. This scheme is planned to be done in our future works, a relatively short time section will be defined and the weights for keywords will be calculated and updated in this time section. Then the updated weights will be used in the tweet score calculation for the next time section.

All the code developed for above analysis is available for download from DDOT-TOA Github page (<http://ddot-toa.github.io/>).

1 **REFERENCES**

- 2 1. Bollen, J., H. Mao, and X. Zeng. Twitter mood predicts the stock market. In *Journal of*  
3 *Computational Science*, Vol. 2(1), 2011, pp. 1-8.
- 4 2. Schmidt, C. W., Trending Now: Using Social Media to Predict and Track Disease Outbreaks. In  
5 *Environmental Health Perspectives*; 2012, Vol. 120 Issue 1, pp30-33.
- 6 3. Sabra, J. Emergency Management 2.0: Integrating social media in emergency  
7 communications. In *Journal of Emergency Management*, Vol. 9(4), July/August, 2011, pp.  
8 15- 18.
- 9 4. Schulz, A., J. Ortmann, F. Probst, Getting user-generated content structured: Overcoming  
10 information overload in emergency management. In *Proceedings of 2012 IEEE Global*  
11 *Humanitarian Technology Conference (GHTC 2012)*, 2012, pp.1-10.
- 12 5. Makice K., *Twitter API: Up and Running: Learn How to Build Applications with the*  
13 *Twitter API*. O'Reilly Media, 2009, ISBN-10: 0596154615.
- 14 6. Rajaraman, A., and J. D. Ullman (2011). *Mining of Massive Datasets*. Cambridge  
15 University Press, 2011, pp. 1–17.
- 16 7. Agrawal, R. and R. Srikant, Fast Algorithms for Mining Association Rules in Large  
17 Databases, in *Proceedings of the 20th International Conference on Very Large Data Bases*.  
18 Morgan Kaufmann Publishers Inc., 1994, pp. 487-499.
- 19 8. Spurlock, J. *Bootstrap*. O'Reilly Media, 2013, ISBN-13: 978-1449343910
- 20 9. Alchin M. *Pro Django*, Apress, 2013, ISBN-13: 978-1430258094
- 21 10. Hoberman, S. *Data Modeling for MongoDB Paperback (First edition)*. Technics  
22 Publications, 2014, ISBN-13: 978-1935504702

23  
24