

1 **INFERRING ACTIVITIES FROM SOCIAL MEDIA DATA**

2  
3

4 **Emmanouil Chaniotakis**

5 Centre for Research and Technology Hellas – Hellenic Institute of Transport  
6 6<sup>th</sup> km Charilaou-Thermi Rd., 57001, Thermi, Thessaloniki, Greece  
7 Tel: +30 2310 498 495 Fax: +302310 498 269; Email: [chaniotakis@certh.gr](mailto:chaniotakis@certh.gr)

8

9 **Constantinos Antoniou**

10 Professor, Chair of Transportation Systems Engineering  
11 Technical University of Munich  
12 Arcisstraße 21, D-80333, Munich, Germany

13 e-mail: [c.antoniou@tum.de](mailto:c.antoniou@tum.de)

14

15 **Georgia Aifadopoulou**

16 Centre for Research and Technology Hellas – Hellenic Institute of Transport  
17 6<sup>th</sup> km Charilaou-Thermi Rd., 57001, Thermi, Thessaloniki, Greece  
18 Tel: +30 2310 498 451 Fax: +302310 498 269; Email: [gea@certh.gr](mailto:gea@certh.gr)

19

20 **Loukas Dimitriou**

21 Department of Civil and Environmental Engineering, University of Cyprus,  
22 75 Kallipoleos Str., P.O. Box 20537, 1678 Nicosia, Cyprus  
23 Tel: +35722892286; E-mail: [lucdimit@ucy.ac.cy](mailto:lucdimit@ucy.ac.cy)

24

25

26 Word count: 4,486 words text + 9 tables/figures x 250 words (each) = 6736 words

27

28

29

30

31

32

33 Submission Date: 01-08-2016

34

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21

## **ABSTRACT**

Social Media have been found to produce an unprecedented amount of information that can be extracted and used in transportation research with one of the most promising areas to be the inference of individuals' activities. While most studies in the literature focus on the direct use of Social Media data, this study presents an efficient framework for the inference of users' activities from Social Media data following a user-centric approach. The framework is applied to data from Twitter, combined with inferred data from Foursquare that contain information about the type of location visited. Then, the user data is classified using a density-based spatial classification algorithm that allows for the definition of commonly visited locations and the individual-based data is augmented with the known activity definition from Foursquare. Based on the known activities and the Twitter text, a set of classification algorithms is applied for the inference of activities of all tweets. The results are discussed upon the types of activities recognized and its classification performance. The classification results allow for a wide application of the framework in the exploration of the individuals' activity space.

*Keywords:* Social Media, Activity Space, Data Enrichment, Data Generalization, Classification, Data Collection

## 1 INTRODUCTION

2 From the rise of Social Media, they have received attention from the scientific community, forming  
3 a potentially new stream of research in many fields. The reasons behind this attention can be  
4 summarized upon the unprecedented amount of information that can be extracted and the  
5 opportunities that Social Media platform use can provide on direct communication with users. The  
6 statistics of Social Media use are astonishing: as of April 2016, on Twitter around 250,000 tweets  
7 are posted each minute; while almost 300,000 statuses are updated and about 136,000 photos are  
8 uploaded on Facebook. Studies indicate that around 76% of the online adults use Social Media,  
9 with 83% of them to be using their smart-phones in order to access and interact in Social Media  
10 platforms (pewinternet.com). Another interesting fact is that according to online ranking websites  
11 (alexa.com) Facebook and Twitter are within the 10 most visited websites (3<sup>rd</sup> and 10<sup>th</sup> position  
12 respectively in April 2016) while most of the rest most visited websites (4 out of 8) are search  
13 engines. Within the Social Media research context, several definitions have arose with the most  
14 inclusive to be found to be the one from Kaplan and Haenlein (1) who define Social Media as “a  
15 group of Internet-based applications that are built on the ideological and technological foundations  
16 of web 2.0, and allow the creation and exchange of user generated content”. Probably for the first  
17 time and with the use of Social Media data-sets and platforms, scientists are able to sense peoples’  
18 thoughts, feelings and actions and interact with them in platforms characterized by ease of use and  
19 capabilities of equal and public exchange of opinions.

20 Focusing on the use of data, the information posted by users in Social Media that can be  
21 exploited for research and application, is diverse (2) and depends on the Social Media platform  
22 functionalities and data availability. However, as studies in transportation illustrate (3, 4), the  
23 existence of spatial and temporal characteristics can by default provide information concerning the  
24 transportation system. Not limited to that, textual and user related information (available either  
25 from posts or from the user profile) enrich the information provided and increase the potential to  
26 be utilized towards newer solutions to transportation related problems (5, 6). Finally, posts in  
27 Social Media can include integrated links (URL) from other Social Media platforms such as  
28 Foursquare, Instagram and Facebook, further enriching the information to be acquired.

29 One of the major transportation modelling necessity for which the incorporation of Social  
30 Media is appealing is the development of activity-based models. Although at the frontier of  
31 transportation demand modelling, activity-based models present many challenges especially  
32 concerning the data required for representing individuals activities (7). In this direction data from  
33 Social Media can be utilized to provide the information on the activities users perform.

34 Pertinent literature suggests a wide spectrum of Social Media analysis in transportation.  
35 However, to the authors’ knowledge, only a few focus on users’ activities (8–10), with most of the  
36 studies to focus on the extraction of information following what we will refer to as a post-centric  
37 approach: In post-centric analysis, each Social Media post is viewed as an entity, ignoring in most  
38 cases the user perspective. This is believed to disorient from the necessity to merge the online and  
39 offline space that could potentially offer useful information towards a better representation of the  
40 transportation system and specifically individuals’ activities. On the other hand, user-centric  
41 approaches can provide valuable interpretations on elemental research questions, such as why a  
42 user posts on a Social Media platform, which activities are users performing and how online  
43 activity is linked with users’ offline life.

44 This study provides an efficient user-centric framework for the inference of users’  
45 activities from Social Media data. The framework is applied in data from Twitter, inferred with  
46 data from Foursquare that contains information about the type of location visited. Then, the users’  
47 data is classified using a density based spatial classification algorithm that allows for definition of

1 commonly visited location and the individual based data is inferred with the known activity  
2 definition from Foursquare. Based on the known activities and the Twitter text, a set of  
3 classification algorithms is applied for the exploitation of activities from Twitter text that discussed  
4 upon the types of activities recognized and its performance.

5 The remainder of the text is organized as follows: first the related work on Social Media  
6 for transportation is presented, focusing on the potential uses of Social Media in transportation  
7 research. Next, in the Methods section the data collection and data enrichment methodologies  
8 proposed are presented and discussed, while in the Analysis section the results of its application  
9 are presented. Finally, in the Conclusions section, a discussion takes place on the applicability of  
10 the proposed method including some insights on the applicable future work.

## 11 **RELATED WORK**

12 Information from Social Media for transportation is mainly examined on the spatial component  
13 (geotag) that posts sometimes include (11, 12) and in some cases on the ways that the Social Media  
14 platforms are used for direct communication with customers (13). The use of the Social Media  
15 data can be categorized in two categories: the use of the streaming data that Social Media  
16 commonly provide or the use of the historical Social Media data.

17 Streaming data in the majority of studies originates from Twitter and has been used for  
18 identification of events mainly focusing on social events or system disruptions (12, 14, 15). Data  
19 can also originate from other Social Media platforms oriented towards the provision of information  
20 concerning social events (16). In these works, the proposed methodologies use the textual and  
21 spatial information and in most cases, aggregation is performed on the spatial data and text  
22 identification or topic modelling is performed on the textual information, in order to augment  
23 incidents that can be detected.

24 Historical data from Social Media is mainly exploited upon the enhancement of transport  
25 models and the extraction of information concerning the transportation system. The spectrum of  
26 applications is wide. Some studies explore the spatial and temporal mobility patterns (17, 18). The  
27 scale of those analysis can vary from a global level (17) to the very local level of a specific location  
28 in a city (18). In most cases models are defined to explain the dynamics of online posts in regards  
29 to the offline locations. Additionally, the investigation of the applicability of data from Social  
30 Media for travel demand modelling is explored (4, 8). This research stream steers in some studies  
31 towards the definition of the urban settings and related characteristics, such as Points of Interest  
32 (POIs), Urban Boundaries, Land uses (19, 20), while some studies explore the potential of either  
33 directly deriving or inferring Origin Destination matrices from Social Media (21, 22) or the  
34 identification of user activities (8–10).

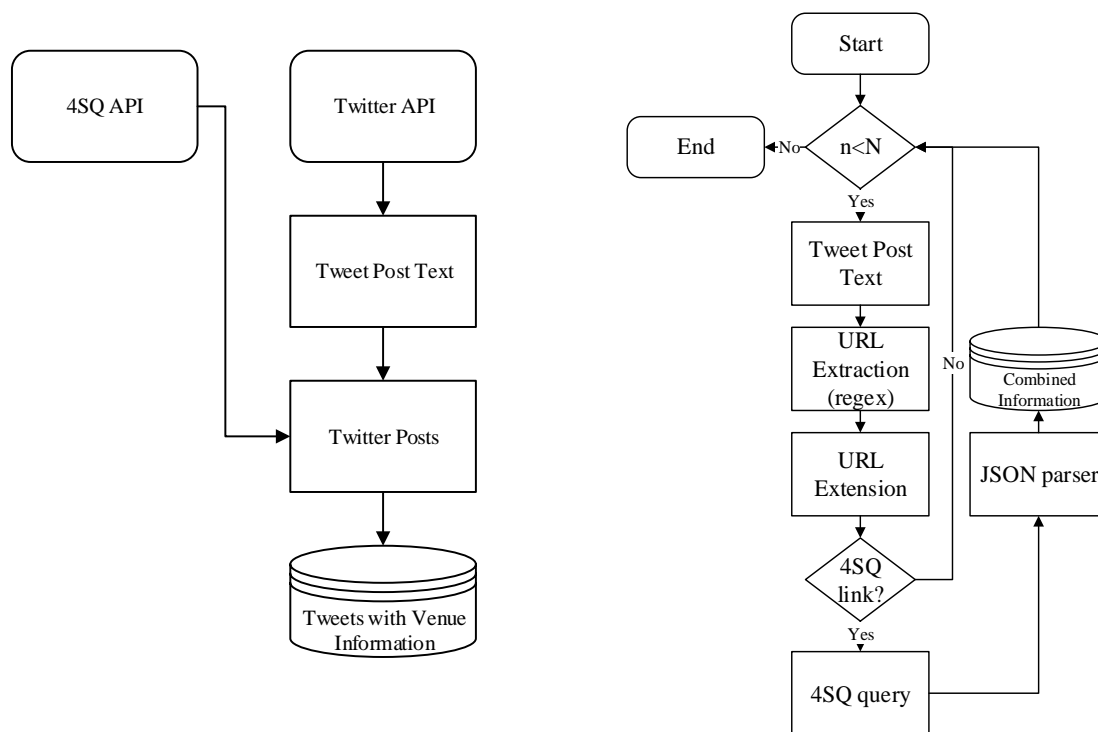
35 This work builds on the finding of Lee et al. (8) and Hasan and Ukkusuri (10) on activity  
36 investigation, from Social Media. The former investigated the activity space of Twitter users and  
37 implemented Latent Class Analysis for the definition of major activities location aiming at  
38 identifying among others the optimal data collection period. The latter proposed a topic-modelling  
39 framework for the definition of multi-day activity patterns of individuals to identify user-specific  
40 activity patterns and the users, who mostly contribute to the definition of such a pattern. Although  
41 both works focus on users' activities, they only use a very small fraction of information available  
42 in Social Media due to constrains towards either the use of only Twitter-originated data or the use  
43 of only Foursquare-originated data. Also, activity inference in the Lee et al. (8) does not focus on  
44 the extraction of the type of activity, while Hasan and Ukkusuri (10) follow an activity-centric  
45 approach which does not identify locations of high density, something that can be very useful in  
46 terms of inferring a complete per user activity space.  
47

## 1 METHODS

### 2 Data Collection and Data Enrichment

3 The Social Media data examined in this work originates from Twitter and was collected using the  
 4 Twitter Application Programming Interface (API) service, for the collection of geotagged tweets.  
 5 A two-phase data collection methodology is applied, where first data is conventionally collected  
 6 by performing Twitter API queries for the last 200 geotagged tweets. This first data collection  
 7 phase allows the creation of a user database; essentially forming a users' sample. Then for each  
 8 user,  $n$  number of tweets are collected by performing queries to acquire their timeline (11). It  
 9 should be noted that the application of the second phase of the data collection methodology from  
 10 Twitter is not bounded in a specific area as it collects all the tweets that individuals post.  
 11

12 Based on the Twitter dataset, extraction of information from Foursquare posts integrated  
 13 in Twitter is performed to achieve the enrichment of activities-related information. In order to  
 14 parse the information, the methodology presented in Figure 1 is applied. This includes the  
 15 extraction of URL patterns using regular expressions, the extension of shortened links to long links  
 16 and Foursquare queries. The queries' response is of JavaScript Object Notation (JSON) format and  
 17 it can be either a venue, or a check-in. In some cases it requires for the link to be (within the  
 18 Foursquare query) resolved. Finally, the response is parsed in tabular form and stored with the  
 19 initial tweet information. It should be noted, that the required queries are subject to the limitations  
 20 that the Twitter and Foursquare APIs impose. At the time this methodology was applied (June  
 21 2016), the maximum number of queries was 500 queries per hour. The application of the  
 22 methodology for enriching twitter data was performed in R (23) using related packages  
 23 (ThinkToStartR, stringr, rjson, RCurl, httr). The only exception was the URL extension component  
 24 that was applied in Java.  
 25



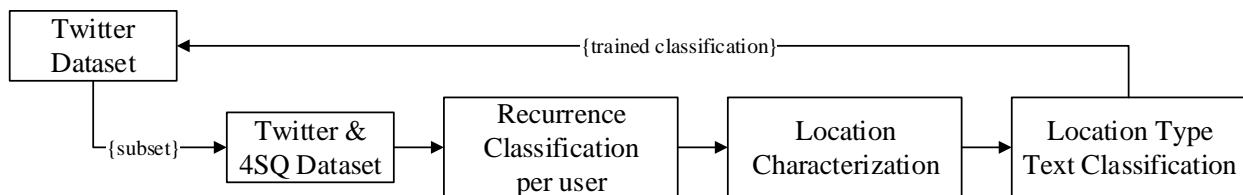
26  
 27 **FIGURE 1 Twitter data enrichment from Foursquare: high-level methodology (left);**  
 28 **detailed parsing methodology (right).**

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25

## User-Centric Activity Enrichment

The resulting dataset from the data collection and dataset definition methodology is used for the extraction of information on the nature of activities users perform and post about in Social Media. The methodology is based on the notion that there are places that users visit frequently to perform activities (e.g. home, work) namely recurring activities, and other places which are scarcely visited or only visited once for non-recurring activities (e.g. leisure). This distinction of places visited and –as a consequence– performed activities is considered of high importance, for the data generalizations that are required aiming at both data reduction – in terms of the identification of locations that are commonly visited – and enrichment – in terms of generalization of activities identified.

The high level representation of the suggested methodology is presented in Figure 2. First, locations are characterized based on their spatial distance and their density, in order to specify locations that are visited by each user more than a specified number of times. The classification algorithm used is the Density-based Spatial Clustering of Applications with Noise (DBSCAN) (24), which characterizes locations visited less than a specified number of times and under a specified density as noise. This is very convenient for the purposes of the study, as it allows for the distinction of locations that are frequently visited and can be all together characterized by only one Foursquare post, in a robust and fast way. The classified locations receive a characterization based on the venue categories they belong to. The characterized – based on recurrence – tweets combined with venue and temporal characteristic are then used for the definition of document-term matrices specified to each type of activities that can be performed in venues. This allows for direct classification of each tweet describing an activity, in an activity cluster using classification algorithms. It should be noted that not all tweets are used to describe activities.



26  
27  
28  
29

**FIGURE 2 User-Centric Activity Enrichment methodology (high level)**

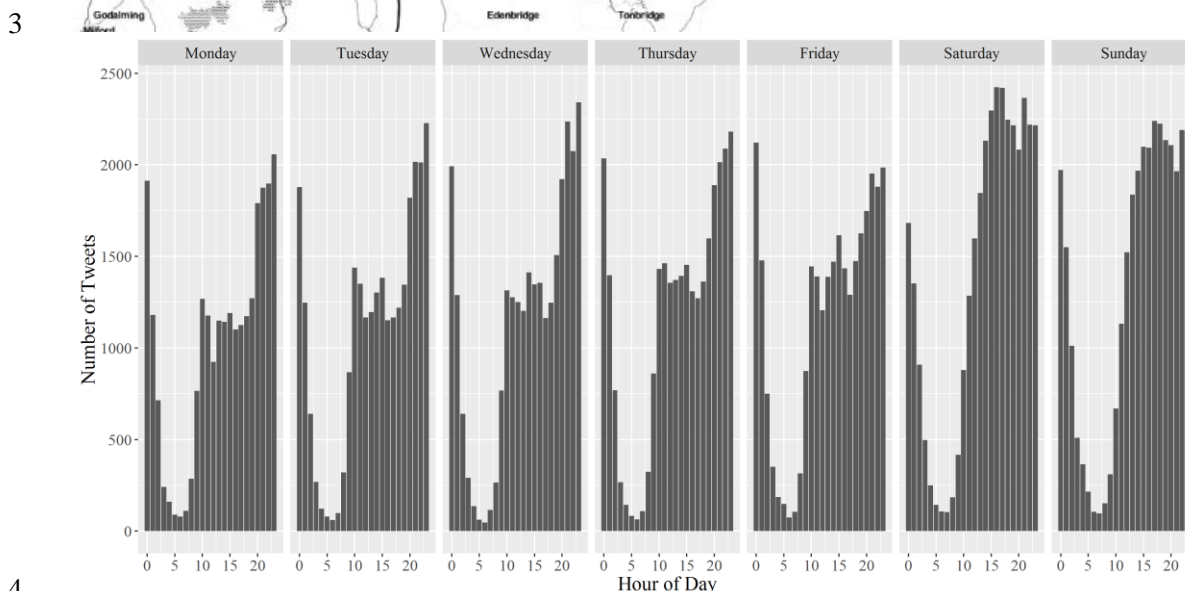
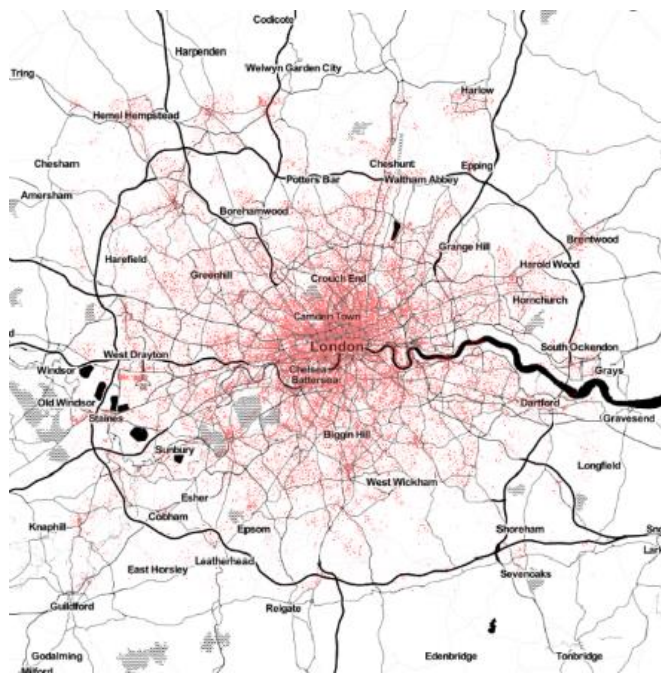
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43

## Dataset Definition

The first phase twitter data collection method was used for continues data collection in a period of one year for the greater London area. The resulted dataset consists of 482,883 unique users and approx. 4.5 million tweets. For a random sample of those users (90,000 users) the last 200 tweets were collected using the second phase of the Twitter data collection methodology. In total, the database included 11,060,814 tweets. From those, 8,141,996 were tweeted in the greater London area. From the geotagged tweets in the greater London area, 3,764,230 (46.2%) included a link that could be parsed with venue information, 220,118 of which originated from Foursquare (2.7%).

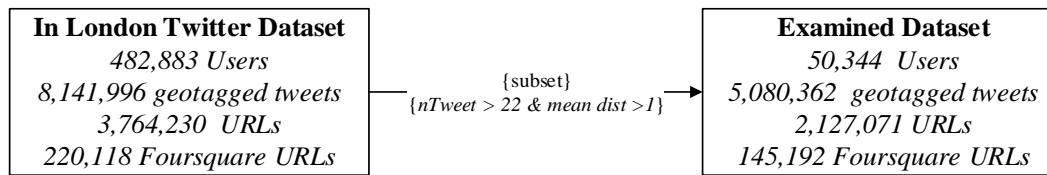
The spatial and temporal characteristics of the sampled tweets are presented in Figure 3. Concerning the spatial distribution, it is clearly evidenced that the urban structures can be identified, as tweets seem to also be concentrated in the urban areas around the city of London. On the temporal level, the distributions presented indicate a peak among evening hours, while the highest number of unique tweets takes place during the weekends and specifically Saturday afternoon. This trend is found to be similar in other cities examined in the literature (3, 11)

1 indicating a tendency towards posting during non-working hours.  
 2



4  
 5 **FIGURE 3. Spatial (top) and temporal (bottom) characteristics of collected Twitter data**  
 6 **for London**

7  
 8 As clarified in the literature, data collected from Twitter contains a large number of users  
 9 who have either only a few tweets or post automated messages from the same location (4). For that  
 10 reason, a subset of users characterized based on their above average posting activity and a threshold  
 11 of average twitter trip length (average direct distance between tweets of a user's determined as the  
 12 average of the distance matrix) was selected, for the application of the User-Centric Activity  
 13 Enrichment methodology. The subsetting data flow and the subset procedure is presented in the  
 14 Figure 4. The average number of tweets was found to be 22 and the average distance threshold  
 15 was defined to be 1 km. The subset included in total 50,344 users who have posted in total  
 16 5,080,362 geotagged tweets.



4 **FIGURE 4 Subset characteristics, in the greater London area and concerning users with**  
 5 **above average Social Media use**

6  
7 **ANALYSIS**

8  
9 **User Characteristics and Recurrence Classification**

10 The methodological framework presented above, was applied for the London city dataset. The  
 11 selected users with above average posting activity were examined upon their posting  
 12 characteristics and the clustering characteristics. First, on the posting characteristics the average  
 13 number of posts each user was posting was found to be 101.8 with a standard deviation of 311.4,  
 14 indicating that there is a number of users that are posting a very large number of tweets. The  
 15 standard deviation of the percentage of tweets posted per day was found to be 0.8% suggesting a  
 16 close to uniform distribution and indicating that on average the users selected are frequent Twitter  
 17 users. From the initially selected 50,344 a fraction of 4,185 users (8.3%) were found to have posted  
 18 at least one tweet that include a foursquare link. The average percentage of posts including  
 19 activities for these users was found to be 39.0%.

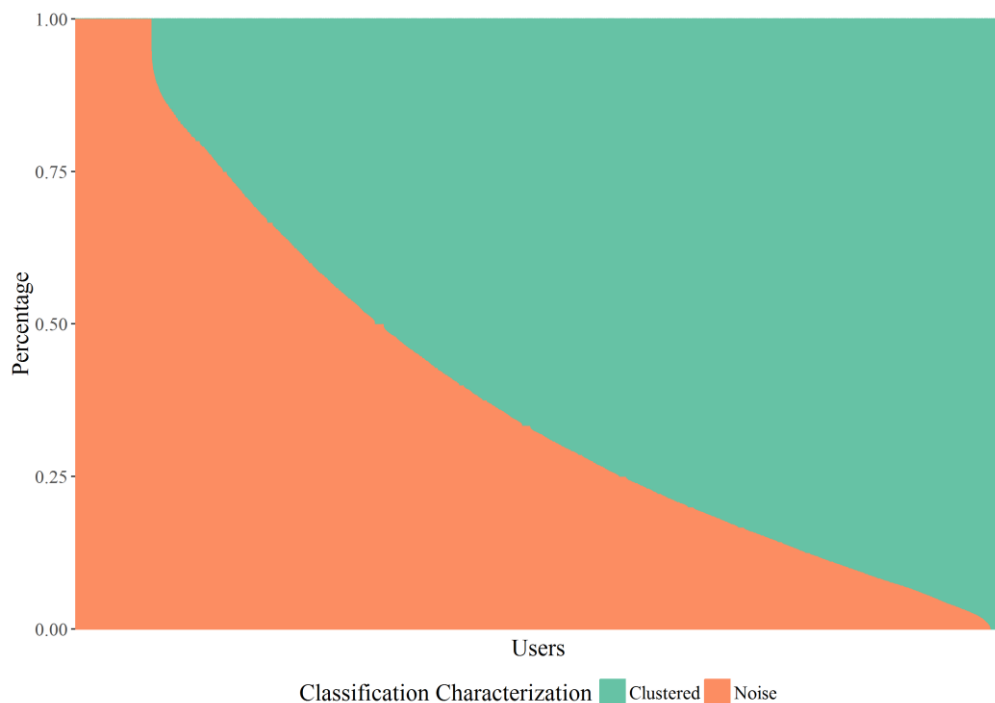
20 On the classification for location recurrence, the analysis was performed in R using the  
 21 *dbscan* library which implements the algorithm presented by Ester et al. (24) for the tweets of each  
 22 user. The parameters were selected after examination of various settings taking into account the  
 23 GPS accuracy (25) and the number of tweets that each individual posts. For this reason, the  
 24 neighborhood of a point parameter (*Eps*) was defined to be 0.002 and the minimum number of  
 25 points to be 5. On the following figure (Figure 5), an example of the classification results is  
 26 presented for two cases: one for which the locations visited cannot be clustered and a second one  
 27 for which the algorithm identifies two classes.



30 **FIGURE 5 Example of classification using Density Based Spatial Clustering of Applications**  
 31 **with Noise: no cluster identified (left) and two clusters identified (right)**



1 The average number of clusters was found to be 2.41 with a standard deviation of 3.18.  
 2 The percentile distributions of the locations that belong in a cluster and those characterized as  
 3 noise per user are presented in Figure 6. It is indicated that a significant amount of users tend to  
 4 only post in locations that were not be included in any cluster. This can be used in profiling of  
 5 individuals who use Social Media based on the type of things that they tend to post.  
 6



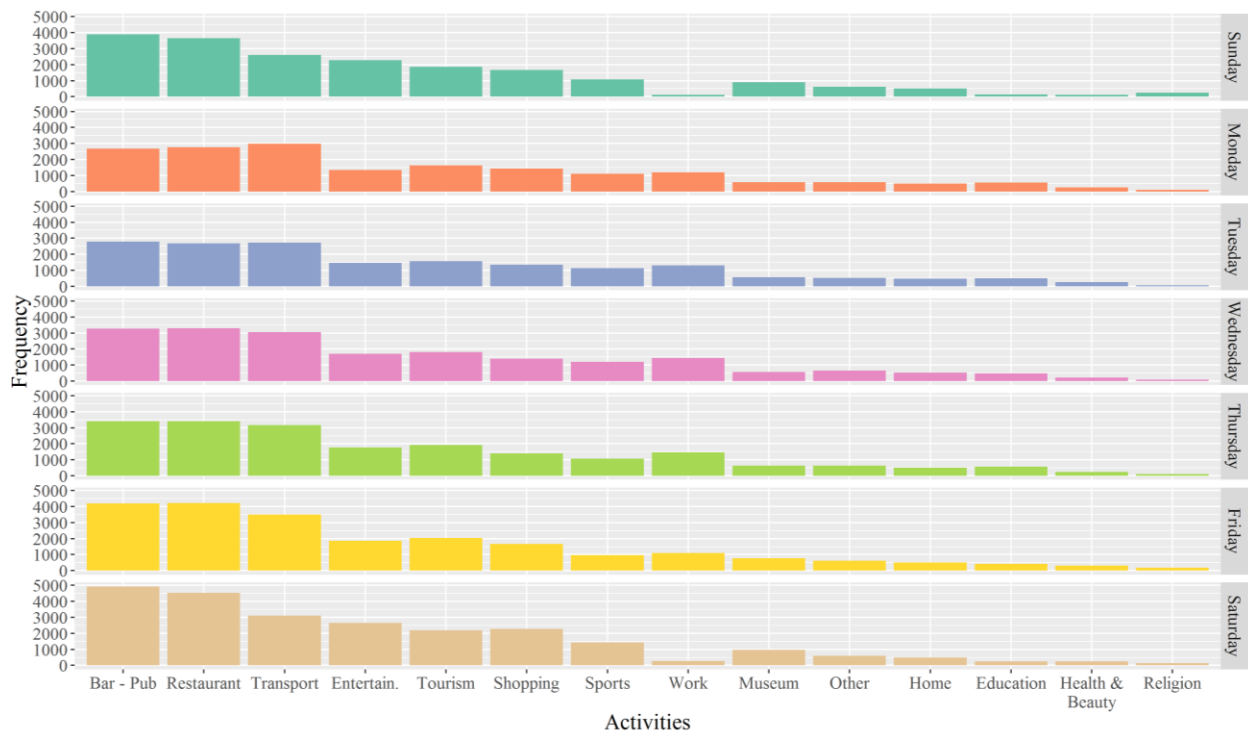
7  
 8  
 9 **FIGURE 6 Example of classification using Density Based Spatial Clustering of Applications**  
 10 **with Noise: no cluster identified (left) and two clusters identified (right)**  
 11

### 12 Activity Data Enrichment

13 The categories provided by Foursquare on the locations visited by users were manually aggregated  
 14 in 14 representative categories and examined upon their temporal characteristics for verification.  
 15 As it is illustrated in Figure 7 there was a tendency towards leisure activities that was present for  
 16 all days of a week. Activities such as education and work were clearly higher represented during  
 17 week days and less represented during weekend days. The activities with the highest representation  
 18 were found to be the “Bar – Pub” and “Restaurant” activities.

19 The application of the data enrichment methodology increases the information concerning  
 20 the activities performed by individuals. First the identified as clustered destinations enriched were  
 21 enriched based on location characterization (activity). More specifically, for the dataset of the  
 22 above average users defined; from the initial number of Foursquare posts the number of posts that  
 23 were assigned in recurrence clusters was found to be 172,675. This is based on the initial 65,806  
 24 tweets that included known location characterization to be associated with a known cluster. It  
 25 should be noted that this data enrichment step further enriched the classification process presented  
 26 below, as it allows for posts that were not meant to denote visitation at a specific location to be  
 27 associated with a cluster location characterization and as a consequence an activity.

28 The last step of the data enrichment methodology is the training of classification model  
 29 that can represent the activity space based on the Twitter text that individuals post. This approach



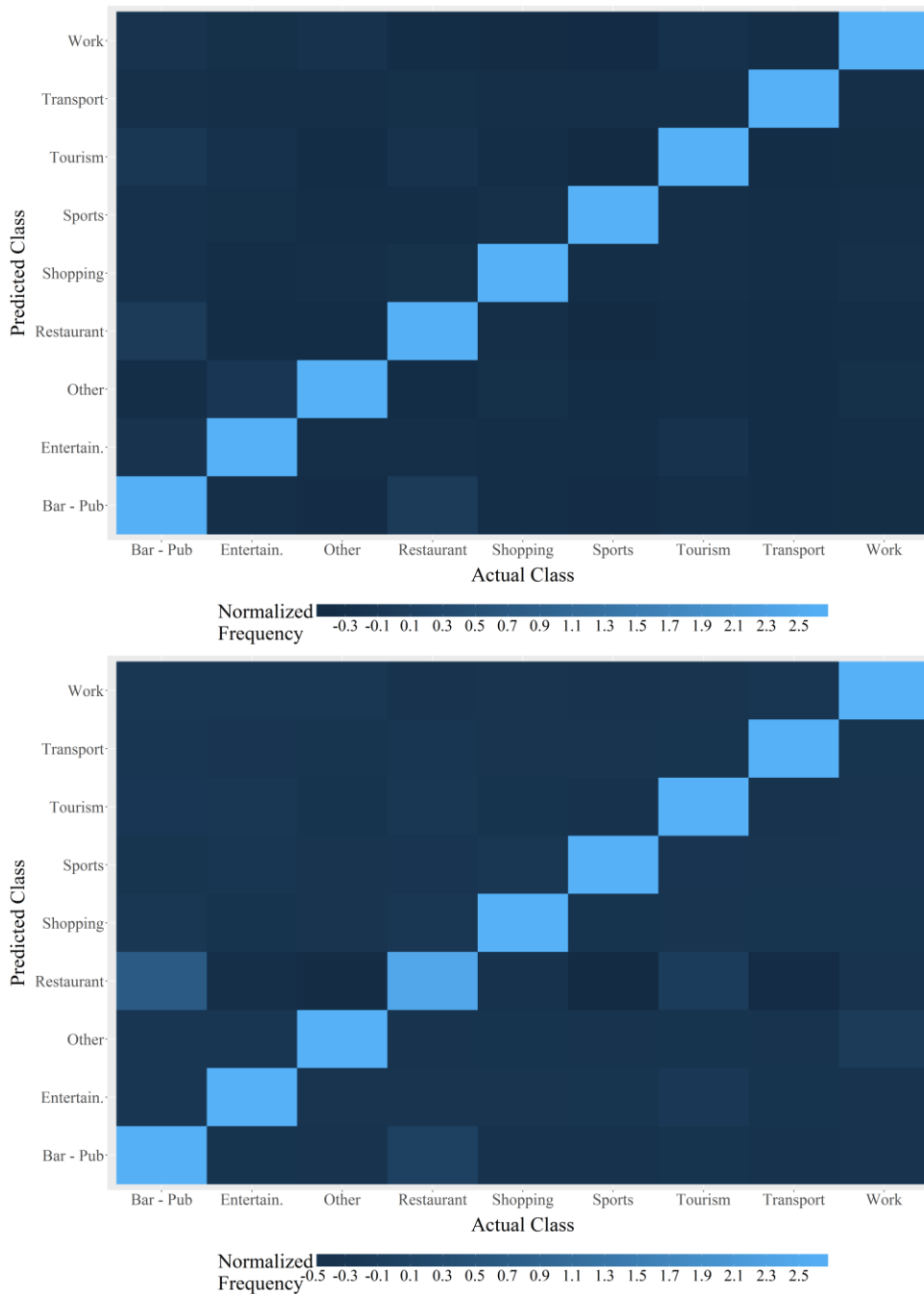
**FIGURE 7 Temporal distribution of Foursquare labeled activities**

would greatly improve the data availability on which researchers build their analysis. A set of classification methods believed to better fit the methodological framework proposed were selected and tested: a) Support Vector Machine (SVM), b) Generalized Linear Model via Penalized Maximum Likelihood (GLMPML) and c) Maximum Entropy (MAXENT). Given computational issues faced when training the classifiers for the whole dataset, the test of the best algorithm took place 10 times on a subset of data that included the major activities identified (Bar–Pub, Entertainment, Restaurant, Sports, Shopping, Other, Work, Tourism, Transport) and for 20,000 randomly selected cases. Each time, the dataset was split in training (85% – 17,000 entries) and testing set (15%, 3000 entries). The overall accuracy results from the classification are presented in the Table 1 and Figure 8. Table 1 includes a summary of the results from the examined classification methods on the metrics of precision, recall and the F–score. Results in Table 1 suggest that both the Maximum Entropy and the Generalized Linear Model via Penalized Maximum Likelihood are considered capable to be used for the identification of activities from Tweets.

**TABLE 1 Classification Performance**

	Average Precision	Precision Standard Deviation	Average Recall	Recall Standard Deviation	Average F–Score	F–Score Standard Deviation
<b>SVM</b>	0.6374	0.0104	0.5808	0.0132	0.6006	0.0117
<b>MAXENT</b>	0.8060	0.0112	0.7773	0.0081	0.7881	0.0080
<b>GLM</b>	0.8392	0.0078	0.6752	0.0068	0.7249	0.0064

1 Furthermore the precision of both algorithm indicate that the results would provide a sufficient  
 2 accuracy, especially taking into account the short number of characters allowed in each tweet.  
 3 Figure 8 on the other hand provides an overview of the classification in terms of confusion for the  
 4 GLM and MAXEXT case. As it is clearly indicated for both classification methods, the pair that  
 5 is mostly confused is the Bar–Pub and restaurant pair, which is logical when taking into account  
 6 that both location categories might involve same activities performed by users (as commonly met  
 7 for example drinking wine).  
 8



**FIGURE 8: Confusion Matrix for the Performance of Classification for the two selected performing classification algorithms: Max Entropy (top) GLM (bottom)**

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47

## CONCLUSIONS

In this research, we present a methodological framework for the characterization of transport-related activities for data originating from Twitter, using inferred data from Foursquare that defines a combined dataset. The resulted combined data is then used to derive classification methodologies for activity characterization using Twitter text. The methods to infer location characterizations (and, as a consequence, activities) and to enrich data, based on the user-centric activity data enrichment, have been applied for a large dataset from London. The results of the analysis indicate the capabilities of the proposed methods to derive activity patterns for individuals and to further use data from Social Media in transportation research and more specifically in activity-based models. Furthermore, the classification performance indicates that, for a sufficient number of cases, the data enrichment can indicate activities performed for all tweets recorded, based on the posted text, which –given the use of endogenous posting data– suggests that the generalization of the proposed methodology in other languages is possible.

Although the results of the analysis seem promising, there are some drawbacks that might be raised in a wide application of the methodology in other study areas and especially upon the final text classification for the inference of activities. To begin with, the methodologies are –of course– bounded to the wide-spread use of Foursquare in the area examined, which might not be the case in some countries or regions. Additionally, the classification has been performed only on posts that include a Foursquare URL or have been clustered in commonly visited venues by user locations. Within this data, only a small fraction was attributed to non-activities related posts (such as drinking wine and running denoted as activities, as opposed to –for example– posting or commenting on news). It is believed that for a proper representation of the activity space of individuals, based on text classification, non-activity related tweets should be detected (either using manual labeling or classification algorithms) to allow for a more complete representation of the reasons that users post on Social Media. Finally, although the examined algorithms seem to perform well, the performance of additional classification algorithms and the sensitivity of their parameters should be examined.

## LITERATURE

- 1 A. M. Kaplan and M. Haenlein, ‘Users of the world, unite! The challenges and opportunities of Social Media’, *Business Horizons*, vol. 53, no. 1, pp. 59–68, Jan. 2010.
- 2 G. Cao, S. Wang, M. Hwang, A. Padmanabhan, Z. Zhang, and K. Soltani, ‘A scalable framework for spatiotemporal analysis of location-based social media data’, *Computers, Environment and Urban Systems*, vol. 51, pp. 70–82, 2015.
- 3 E. Chaniotakis, C. Antoniou, and E. Mitsakis, ‘Data for Leisure Travel Demand from Social Networking Services’, in *4th hEART Symposium*, 2015.
- 4 J. H. Lee, S. Gao, and K. G. Goulias, ‘Comparing the Origin-Destination Matrices from Travel Demand Model and Social Media Data’, in *Transportation Research Board 95th Annual Meeting*, 2016, no. 16–0069.
- 5 A. Gal-Tzur, S. M. Grant-Muller, T. Kuflik, E. Minkov, S. Nocera, and I. Shoor, ‘The potential of social media in delivering transport policy goals’, *Transport Policy*, vol. 32, pp. 115–123, 2014.
- 6 S. M. Grant-Muller, A. Gal-Tzur, E. Minkov, S. Nocera, T. Kuflik, and I. Shoor, ‘Enhancing transport data collection through social media sources: methods, challenges and opportunities for textual data’, *IET Intelligent Transportation Systems*, 2014.
- 7 E. Cascetta, *Transportation systems analysis: models and applications*. 2009.

- 1 8 J. H. Lee, A. W. Davis, and K. G. Goulias, 'Activity Space Estimation with Longitudinal  
2 Observations of Social Media Data', in *Transportation Research Board 95th Annual  
3 Meeting*, 2016.
- 4 9 A. Noulas, C. Mascolo, and E. Frias-Martinez, 'Exploiting foursquare and cellular data to  
5 infer user activity in urban environments', in *2013 IEEE 14th International Conference on  
6 Mobile Data Management*, 2013, vol. 1, pp. 167–176.
- 7 10 S. Hasan and S. V Ukkusuri, 'Urban activity pattern classification using topic models from  
8 online geo-location data', *Transportation Research Part C Emerging Technologies*, vol. 44,  
9 pp. 363–381, 2014.
- 10 11 E. Chaniotakis and C. Antoniou, 'Use of Geotagged Social Media in Urban Settings:  
11 Empirical Evidence on Its Potential from Twitter', in *IEEE Conference on Intelligent  
12 Transportation Systems, Proceedings, ITSC*, 2015, vol. 2015–Octob, pp. 214–219.
- 13 12 Y. Gu, Z. (Sean) Qian, and F. Chen, 'From Twitter to detector: Real-time traffic incident  
14 detection using social media data', *Transportation Research Part C Emerging Technologies*,  
15 vol. 67, pp. 321–342, 2016.
- 16 13 S. Bregman, *Uses of social media in public transportation*, vol. 99. Transportation Research  
17 Board, 2012.
- 18 14 S. Jiang, A. Alves, F. Rodrigues, J. Ferreira Jr., and F. C. Pereira, 'Mining point-of-interest  
19 data from social networks for urban land use classification and disaggregation', *Computers,  
20 Environment and Urban Systems*, 2015.
- 21 15 A. Kumar, M. Jiang, and Y. Fang, 'Where not to go?: detecting road hazards using twitter',  
22 *Proceedings of the 37th international ACM conference*, vol. 2609550, pp. 1223–1226, 2014.
- 23 16 F. C. Pereira, F. Rodrigues, E. Polisciuc, and M. Ben-Akiva, 'Why so many people?  
24 Explaining Nonhabitual Transport Overcrowding With Internet Data', *IEEE Transactions  
25 on Intelligent Transportation Systems*, vol. 16, no. 3, pp. 1370–1379, 2015.
- 26 17 B. Hawelka, I. Sitko, E. Beinat, S. Sobolevsky, P. Kazakopoulos, and C. Ratti, 'Geo-located  
27 Twitter as proxy for global mobility patterns', *Cartography and Geographic Information  
28 Science*, vol. 41, no. 3, pp. 1–12, 2014.
- 29 18 Y.-S. Cho, G. Ver Steeg, and A. Galstyan, 'Where and Why Users“ Check In”.', 2014.
- 30 19 J. Li, Q. Qin, J. Han, L.-A. Tang, and K. H. Lei, 'Mining Trajectory Data and Geotagged  
31 Data in Social Media for Road Map Inference', *Transactions in GIS*, vol. 19, no. 1, pp. 1–  
32 18, Feb. 2015.
- 33 20 B. Jiang and Y. Miao, 'The evolution of natural cities from the perspective of location-based  
34 social media', *The Professional Geographer*, vol. 67, no. 2, pp. 295–306, 2015.
- 35 21 F. Yang, P. J. Jin, X. Wan, R. Li, and B. Ran, 'Dynamic origin-destination travel demand  
36 estimation using location based social networking data', in *Transportation Research Board  
37 93rd Annual Meeting*, 2014, no. 14–5509.
- 38 22 P. Jin, M. Cebelak, F. Yang, J. Zhang, C. Walton, and B. Ran, 'Location-Based Social  
39 Networking Data: Exploration into Use of Doubly Constrained Gravity Model for Origin-  
40 Destination Estimation', *Transportation Research Record: Journal of the Transportation  
41 Research Board*, no. 2430, pp. 72–82, 2014.
- 42 23 R. C. Team and others, 'R: A language and environment for statistical computing', 2013.
- 43 24 M. Ester, H.-P. Kriegel, J. Sander, and X. Xu, 'A density-based algorithm for discovering  
44 clusters in large spatial databases with noise.', in *Kdd*, 1996, vol. 96, no. 34, pp. 226–231.
- 45 25 M. Schaefer and T. Woodyer, 'Assessing absolute and relative accuracy of recreation-grade  
46 and mobile phone GNSS devices: a method for informing device choice', *Area*, vol. 47, no.  
47 2, pp. 185–196, 2015.