# Vehicular traffic flow intensity detection and prediction through mobile data usage

Maurice Saliba[1], Charlie Abela[1], and Colin Layfield[2]

[1] Department of Artificial Intelligence,
Faculty of ICT, University of Malta, Malta,
`mauricesaliba@gmail.com`, `charlie.abela@um.edu.mt`
[2] Department of Computer Information Systems,
Faculty of ICT, University of Malta, Malta,
`colin.layfield@um.edu.mt`

**Abstract.** A novel approach, consisting of an ensemble of data-mining and machine learning techniques, is proposed to prove that it is possible to extract and predict vehicular traffic patterns from mobile usage data. An anonymized mobile phone usage dataset from a telecommunications provider in Malta was used to generate an origin-destination (OD) matrix that defines the top two locations towards which each user travels to through clustering. The OD matrix was used to infer user trips over fastest routes between these top two locations across time. We then applied spatial binning techniques to deduce the aggregate distribution of traffic load on the traffic network. A predictive model based on an artificial neural network was trained with grid nodes' traffic levels in a time series to predict traffic level for specific nodes.
Our findings are promising and show that the built models are more effective to measure and predict traffic flow demand for specific locations rather than the actual traffic flow rate. The proposed solution needs improvement by adding a dynamic traffic assignment to the whole algorithm. This would give more accurate results, especially for traffic flow points that tend to be congested, by capturing user route selection changes and get more precise localization of delay causes.

**Keywords:** Mobility patterns · Vehicular traffic congestion · Artificial intelligence · Machine learning.

## 1   Introduction

The dynamics of traffic flow are determined by the travel needs of the masses. The daily commutes of every individual impacts those of others. The interaction on a large scale of all the vehicles in a time series is difficult to model and to predict in a robust and responsive manner [19]. Traffic sensors, cameras and induction loops are all sources of information that can be used to both detect high traffic intensity or even forecast it beforehand. However, traffic intensity measurement with these methods is physically limited. Camera feeds and

inductive-loop detectors cannot be installed in every road of the transport infrastructure. Devices carried by travellers, or embedded in vehicles, are more practical to build smart solutions for traffic management [19].

Mobile traces can be processed and used to offer location based services that have a wide application spectrum that go beyond solving mobility issues [10, 4, 8, 9]. This formidable data source, however, poses a challenge. Location data, which usually comes in large amounts, has to be harvested, ingested efficiently and ideally processed in real time for the required final purpose which is value added location based services.

The range of applications and branches of research abound on remote collection of mobile users' geolocation information. To name a few, applications include: traffic patterns and prediction modelling, crowd management, hotspot detection, lost device recovery, emergency rescue, use for investigative authorities, location-based recommendation and advertising systems, contextualized information, social interaction based application, epidemiology etc.

Calabrese et al [4] emphasized that studies on human mobility patterns would be vital for improved and sustainable urban planning and could boost the environment's well being given that transportation in 2004 already accounted for 22% of primary energy use.

Steenbruggen et al [17] discuss how mobile geolocation data can be used to differentiate weekday traffic patterns from those of weekends. Another specific type of prediction based on mobile usage discussed in [10] is jam detection. Macroscopic monitoring and analysis of vehicle mobility through mobile traces is a wide area of study that has ramifications in many areas of research [17].

In this paper we focus on measuring traffic flow and the prediction of traffic flow changes over time for a selection of locations by using mobile data usage. A combination of data mining and machine learning techniques are used to devise a data processing pipeline. Through this pipeline it is possible to:

1. process raw event data records containing cell tower locations, date and time based on which we carry out preliminary descriptive statistical analysis;
2. zoom into the main areas of activity of users by using unsupervised machine learning techniques that cluster the most dense groups of geolocation data points;
3. determine routes between these main activity areas and collect spatial-grid aggregated data from daily trips done along these routes from thousands of users;
4. use the transformed data that is representative of traffic flow in various locations to train and validate a predictive model using artificial neural networks [4, 18, 10, 2];
5. feed visualization tools that enable insightful visual inspection of traffic patterns projected on maps.

A selection of methods that are encountered in literature are applied and evaluated. The real challenges arise in the quest for a high spatio-temporal resolution when modelling traffic, given that mobile usage records' geolocation dataset is

sparse and tracks the position of users with a considerable margin of displacement error [10, 8]. In section 2, background and related work, we will expound on techniques used for general human mobility modelling from mobile call detail records (CDRs) found in literature. In section 3 the methodology to form the data processing pipeline is described. Sections 4 and 5 follow with results and evaluation and a conclusion with possible future work respectively.

## 2 Background and Related work

In this section we will go over mainstream techniques and approaches that make use of mobile data for traffic flow detection and prediction.

### 2.1 Mobile location data sources

Our research revolves around geolocation attributes of mobile data usage. Mobile device location traces have their limitations when used for vehicular traffic analyses. In contrast to surveys, they lack demographics [4, 6] and the market share of a mobile service provider that made the dataset available for scientific research might not be really representative of the commuting patterns [3]. Many studies highlighted the importance of removing bias when pre-processing such datasets before any further processing is done [12, 18]. Passive data, gathered in the form of CDRs, are not suited to extract different modes of travelling, route assignment and the classification of detailed activity types [6].

Mobile device location data is not only limited to data that originates from cellular networks. Global Positioning System (GPS) is the most current reliable source of geolocation because of their higher resolution with a lower margin of error [9]. Using GPS data for a mobility study is more challenging because it needs the continuous consent of users to get such data and drains the battery quickly especially because of long signal acquisition time [20, 1].

Research generally focuses on voice CDRs to trace mobility. Mobile data usage was rarely found in literature to be used to detect vehicular traffic or predict it because of unavailability of such datasets [10, 3].

### 2.2 Origin and destination matrices computation

A recurrent topic in traffic flow analyses is the study of how to deduce origin and destination (OD) locations for travelling vehicles [12]. Numerous research work focused on tackling the problem posed by traffic congestion detection by first deducing the OD matrix [18, 12, 2–4, 6].

ODs are used to extract main activity hubs. Gonzalez et al [8] state that 40% of the time users are at their two preferred locations. Therefore most trips can be mostly explained as being between several locations since users tend to be highly inclined to be regular in spatial and temporal terms. All this leads to the safe assumption that the majority of trips are between home and work. In literature it is commonly found that locations that were likely to be recorded in

OD matrices were home and work [3, 6, 16]. In [3] home location is detected by checking which 500 metres square cell has the most activity during the night for every specific user. Colak et al [6] and Calabrese et al [3] also label zones such as home and work and try to find purpose behind other types of trips.

### 2.3   Trip generation from OD matrices

Route selection is necessary to link origins and destinations from OD matrices to generate OD trips. In Iqbal et al[12] the route is determined by a function of least travel time path. In Toole et al[18], Open Street Maps[3] (OSM), which is an open source map editing framework, is used to infer routing. Some studies assign trips to a user when there are consecutive calls in the same day and the calls are done from different locations. Typically, two consecutive 'stays' that are not more than 1 day apart would constitute a trip [6, 18].

## 3   Methodology

In this section we describe step by step how we extracted traffic flow information and built predictive models from mobile data usage.

### 3.1   Dataset used

Experimentation was done on an anonymized cell generated CDRs' dataset that was provided by GO Plc Malta[4], which is one of the main Malta telecommunication service providers. The dataset was recorded in October 2016 and had approximately 100 million mobile data usage records. Main data fields of interest in the data structure were an anonymous identifier, timestamp, volume, duration and geo-coordinates.

### 3.2   Traffic flow detection by trip generation assigned traffic

The method we adopted to detect traffic on the road network involved first the generation of an OD matrix that contained main stay locations for users in a time series. A trip was generated between each main location for each user as it will be explained in section 3.4. The trip information includes turn by turn directions with longitude and latitude coordinates. Traffic load was assigned to junctions and turns depending on the time retrieved from OSM data. The major challenge here proved to be the traffic assignment, given that there is an interaction of a lot of vehicles at a given point in time with a complex structure of roads and possible unexpected events such as weather, accidents and road closures.

---

[3] https://www.openstreetmap.org
[4] https://go.com.mt

### 3.3   Main activity hubs extraction through DBSCAN clustering

One of the main steps of the proposed solution was to derive main areas of activity from the mobile data usage of subscribers. This was achieved using clustering. Clustering was also used to remove noise in the form of sudden displacements through frequent oscillations between cell towers by finding a centroid of activity. DBSCAN (density based spatial clustering of applications) clustering was chosen over k-means for reasons similar to those explained in [5, 11]. DBSCAN does not need to set the number of clusters at the outset. Moreover it finds clusters of non-spherical nature and leaves noisy elements out of the computed clusters [5]. The algorithm is more sensitive to density rather than to aggregate distances of surrounding points.



**Fig. 1.** DBSCAN clustering to find main user activity hubs. (Sample illustration)

The chosen values for DBSCAN hyper-parameters were 500m for radius (based on average distance of 350m between cell towers in urban areas), minimum required points was set to 3 and euclidean distance was chosen as the distance metric. Figure 1 shows plotting of a sample of mobile activity clusters delineated by rectangular boundaries.

### 3.4   OD Matrix trip Generation

We decided to focus on two main areas of activity per user as the basis of our OD matrix generation, namely home and work location. This was based on results reported in the literature review (refer to section 2).

The top two mobile data usage activity clusters per user where retrieved from the resulting users' clusters computed through the DBSCAN. The user's CDRs

that had geographical coordinates located in the two main activity cluster areas were then filtered into a new dataset through the spatial joining technique [7]. This process resulted in a dataset containing all data usage records that had a location in either of the top two clusters for any user in the time series.

### 3.5   Trip generation, route choice and traffic assignment

We inferred the routes between origin and destination from the OSM using a method discussed in Toole et al [18]. The fastest route was assigned for each entry in the OD matrix together with duration information from the trip. The routing engine *Open Source Routing Machine* (OSRM)[5] was used for this purpose. Choosing the fastest route by default is a limitation of this research and must be considered as a source of bias.



**Fig. 2.** Traffic flow count snapshot at 7:00 a.m. October 2016. Illustrated through CartoDB temporal mapping. Tool allows to visually investigate traffic by sliding the date-time setting. Circle size is proportional to traffic flow count. Large circles are depicted in notorious traffic hotspots in Malta.

---

[5] http://project-osrm.org/docs/v5.15.2/api/#route-service

The difference between the actual trip duration retrieved from observed departures and arrivals per user, and the OSRM derived trip duration, was considered to be the global trip delay. After computing delays for each trip per user, aggregated statistics were collated to describe typical delays at different hours for both weekdays and weekends. Trip delay differences are evident even between Saturdays and Sundays but they were highly similar for weekdays.

Traffic was assigned to the road network depending on manoeuvres' steps with geolocation given by the OSRM. These steps have the time information when user travelled through the geolocation. This information was used to distribute traffic count on the road network. Figure 2 illustrates traffic flow count data at a given point in time.

### 3.6   Prediction using a Multilayer Perceptron Classifier (MLPC)

The next step in the data processing pipeline consisted in predicting traffic from a stipulated time ahead for a given location point. This prediction had to be based on data that was harvested some time before. All of the original datasets had records with timestamps set in the past, so we simulated prediction of traffic flow by trying to forecast traffic at a certain point in time which is ahead of a given timestamp. Evaluation was then carried out with a variable number of first principle component analysis (PCA) components, prediction multi-steps ahead and possible classes that describe level of traffic. The training and testing inputs for the MLPC model were traffic counts in the grid and the output was the level of traffic at a certain location in the relative future.

The multi-step time series prediction model was trained and tested with a variable amount of steps ahead. Each step was already defined to be 5 minutes long. The experimentation was performed with 3, 6, 12 and 288 steps ahead that transalte to 15 minutes, 30 minutes, 1 hour and 1 day.

## 4   Results and evaluation

In our evaluation process we evaluated four main experimental procedures:

1. Average trip count per hour for weekdays and weekdays;
2. Average global trip delay per hour for weekdays and weekends;
3. Traffic flow count in a selection of locations;
4. Traffic count prediction for a selection of locations.

### 4.1   Average trip count per hour

Through a linear regression we showed that the trip distribution derived from mobile usage CDRs' generated OD Matrix has a significant correlation with the trip distribution as reported in a National Household Travel Survey (NHTS) conducted in 2010 [14] (see figure 3). The same linear regression model was used to scale up the trip distribution in 2010 to the one registered in 2016 for this study. A correlation statistical analysis gave the result of a Pearson correlation coefficient of 0.94 with a p-value of $1.13628e - 11^{*}$.
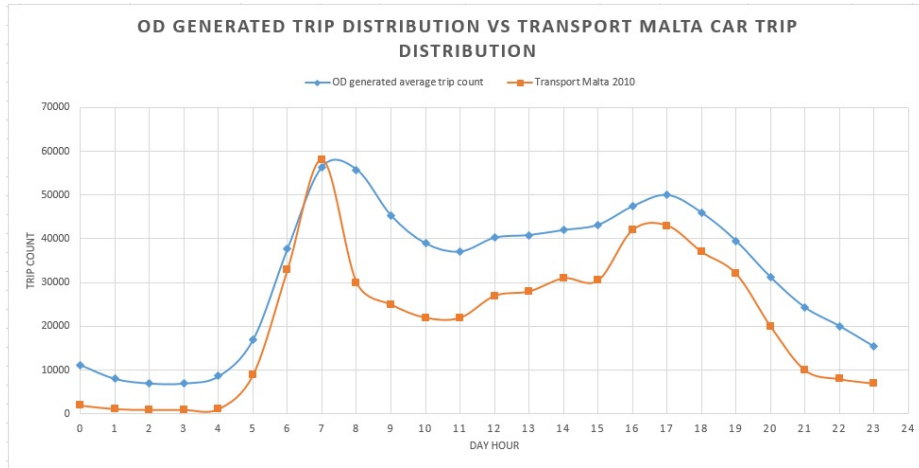
**Fig. 3.** Comparison between OD average trip distribution over a month and Transport Malta 2010 survey results ([14])

### 4.2   Trip average delay per hour

Seven whole days of Google Distance matrix API (DMAPI) data from June 2018 was scraped by retrieving duration information for every quarter of an hour. Estimated trip delay was calculated by subtracting the estimated trip duration from the trip duration in traffic. Average trip delay was then computed per hour for four different Malta routes that link cities, namely Mosta to Marsa, Mellieha to Swieqi (MS), Birkirkara to Sliema (BS) and Valletta (the capital city) to Mgarr (VM). The overall trip delay average was calculated for these routes.

Correlation results showed that there is a strong linear relationship between the routes' trip delay pattern which were investigated with the DMAPI. Correlation between DMAPI and OD-OSRM trip delay estimation was less but still considerable. Between DMAPI average overall trip delay and OD-OSRM non shifted expected trip delay data, there was a correlation of 0.69 ($p < 0.001$). Correlation was computed between data retrieved in June 2018 for DMAPI and data retrieved in October 2016 for OD-OSRM. In October traffic in Malta is much heavier than in June because schools start in this period. It is important to note that during summer, government employees work half days and schools are closed.

### 4.3   Traffic flow count

The ground truth to evaluate traffic flow count experimentation came from work done by Nigel Pace in his dissertation submitted in 2017 [15]. Directional traffic flow counts were manually gathered from web camera streams recorded from four locations. These were gathered from Kappara and Marsa roadways for traffic which was both northbound and southbound. The Marsa roadway is referred to

as the Marsa-Hamrun bypass, which is the road leading to and from the Santa Venera tunnels. The Kappara roadways get and feed traffic to the old Kappara roundabout which today has been replaced by a flyover. The dates for the data collection were from Monday 17th October to Friday 21st October 2016. Data for the day of Tuesday 18th October was missing from the dataset. The traffic flow count consisted of an average traffic flow count per minute taken over intervals of 15 minutes. This resulted into 11 samples per location for data gathered from 6.00 a.m to 8.45 a.m. for every day. A daily average for every quarter of an hour was then taken for both the actual data and the one generated with OD-OSRM.

| linear regression model OD-OSRM traffic counts vs video stream traffic counts | regression result output parameters | | | |
| --- | --- | --- | --- | --- |
| | Pearson's Coefficient | $R^2$ | p-value [*] | degrees of freedom |
| Kappara North | 0.68 | 0.46 | 4.13E-07 | 43 |
| Kappara South | 0.75 | 0.56 | 3.85E-09 | 43 |
| Marsa North | -0.46 | 0.22 | 0.0013 | 43 |
| Marsa South | -0.31 | 0.10 | 0.04 | 43 |

**Table 1.** Correlation statistics for linear regression models. OD-OSRM traffic flow count is the predictor variable and video stream traffic count is the dependent variable. [*]Results are all significant with $p < 0.05$.

A regression was established between OD-OSRM traffic count and video stream traffic counts and results are shown in table 1. There is strong correlation with Kappara traffic flows[6] but a weak negative one with Marsa[7] located traffic flows. This can be attributed to the fact that traffic tends to be slower in Marsa traffic flow points when compared with the Kappara traffic flow points. OD-OSRM measurements were based on trips that had been detected but if actual vehicular traffic slows down due to congestion the OD-OSRM traffic flow count does not reflect actual traffic counts. Therefore, two conclusions are derived from this. The first conclusion is that reliable regression models can be trained on actual traffic data for traffic flow road sections which do not experience heavy traffic slow down. Secondly the regression model mapped traffic flow counts gave a reliable account of what flow capacity is 'expected' to be serviced at any given point in time from a given road section in order that traffic flows smoothly.

## 4.4 Traffic flow count prediction for a selection of locations

Table 2 shows the evaluation results obtained for traffic flow count prediction. It is evident that models trained to predict for smaller time ahead intervals generally perform better than models that are trained with a lengthier prediction

---

[6] geocordinates: 35.904416, 14.487168
[7] geocordinates: 35.898072, 14.486804

time interval for the same location. Performance of prediction of four levels of traffic was done with level one being low or no traffic and level four having the highest level of traffic. The levels were mapped with a logarithmic function as a ratio to the highest level of traffic. All models proved to have highest recall and precision for class one traffic flow counts.

| Prediciton evaluation metrics for a given traffic flow point with variable prediction ahead | | Evaluation results' metrics | | | |
|---|---|---|---|---|---|
| Traffic flow section | Prediciton time interval ahead | Accuracy | weighted Precision | weighted Recall | F1-Score |
| Kappara south bound | 15 minutes | 0.67 | 0.66 | 0.67 | 0.66 |
| | 30 minutes | 0.67 | 0.66 | 0.67 | 0.66 |
| | 60 minutes | 0.64 | 0.64 | 0.64 | 0.64 |
| | 1 day | 0.61 | 0.60 | 0.61 | 0.61 |
| Imsida skate park | 15 minutes | 0.70 | 0.70 | 0.70 | 0.70 |
| | 30 minutes | 0.68 | 0.70 | 0.68 | 0.69 |
| | 60 minutes | 0.68 | 0.69 | 0.68 | 0.68 |
| | 1 day | 0.62 | 0.62 | 0.62 | 0.62 |
| Hamrun - Valletta | 15 minutes | 0.91 | 0.89 | 0.91 | 0.90 |
| | 30 minutes | 0.91 | 0.89 | 0.91 | 0.90 |
| | 60 minutes | 0.90 | 0.89 | 0.90 | 0.90 |
| | 1 day | 0.90 | 0.88 | 0.90 | 0.89 |
| Marsa roadway leading to Aldo Moro | 15 minutes | 0.67 | 0.68 | 0.67 | 0.67 |
| | 30 minutes | 0.66 | 0.66 | 0.66 | 0.66 |
| | 60 minutes | 0.64 | 0.65 | 0.64 | 0.65 |
| | 1 day | 0.58 | 0.58 | 0.58 | 0.58 |

**Table 2.** Classification evaluation metrics for 4 traffic flow road sections with 4 label classification and PCA set to extract 324 first components. Testing was done with 4 sizes of prediction time window ahead for each prediction location.

It appears from the table that the best overall classification metric scores were attained for Hamrun-Valletta roadway. However on examination of a confusion matrix for classification results per label we noted that the model performed very badly for high traffic flow count classes. There were no results for class four and for classes two and three the precision and recall metrics are very low. In fact, when computing the F1-score for classes two and three, both result were found to be low at 0.14 and 0.0 respectively. Lv et al stated that ANN models trained with low traffic counts do not perform well. In this same work evaluation relative error is greater when traffic flow is small. Results are only being quoted when traffic flow measurement amounts to 450 vehicles or more for a 15 minute time window [13].

In contrast predictive results for Marsa road that leads to Aldo Moro street are less promising than those for Hamrun-Valletta arterial road. Still, the predictive efficacy results are very good, especially when examined in the perspective

of confusion matrices. Class four cases, which are classified as class one or class two cases are very few and, even if almost half of class four test values were predicted as class three, in practice, this would still make the model useful.

## 5   Conclusions and future work

This research posed questions on whether it is feasible to get vehicular traffic descriptive and predictive analytics from mobile usage data. We showed how from top users' activity locations it is also possible to achieve accurate results in getting global trip counts and trip delays. Trip data was then used to actually map traffic flow demand on the road grid. However, it was found that traffic flow mapping gave more accurate results when the level of traffic congestion was low.

Finally, an MLPC was found to be really efficient in predicting traffic flow for a set of locations. The confidence level given by the prediction results is high and if traffic flow input used to train the predictive model is accurate the method we devised could be used for practical scenarios to forecast traffic in real-time.

## References

1. Ahas, R., Tiru, M., Saluveer, E., Demunter, C.: Mobile telephones and mobile positioning data as source for statistics : Estonian experiences. Presentation for NTTS (2011)
2. Alexander, L., Jiang, S., Murga, M., González, M.C.: Origin-destination trips by purpose and time of day inferred from mobile phone data. Transportation Research Part C: Emerging Technologies **58**, 240–250 (2015). https://doi.org/10.1016/j.trc.2015.02.018, http://dx.doi.org/10.1016/j.trc.2015.02.018
3. Calabrese, C., Giusy, F., Lorenzo, D., Liu, L., Ratti, C., Calabrese, F., Lorenzo, G.D.: Estimating Origin-Destination flows using opportunistically collected mobile phone location data from one million users in Boston Metropolitan Area Terms of Use Estimating Origin-Destination flows using opportunistically collected mobile phone location da. IEEE Pervasive Computing **10**(4), 36–44 (2011). https://doi.org/10.1109/mprv.2011.41
4. Calabrese, F., Diao, M., Di Lorenzo, G., Ferreira, J., Ratti, C.: Understanding individual mobility patterns from urban sensing data: A mobile phone trace example. Transportation Research Part C: Emerging Technologies **26**, 301–313 (2013). https://doi.org/10.1016/j.trc.2012.09.009, http://dx.doi.org/10.1016/j.trc.2012.09.009
5. Chakraborty NKNagwani Lopamudra Dey, S.: Performance Comparison of Incremental K-means and Incremental DBSCAN Algorithms. International Journal of Computer Applications **27**(11), 975–8887 (2011)
6. Çolak, S., Alexander, L.P., Alvim, B.G., Mehndiratta, S.R., González, M.C.: Analyzing Cell Phone Location Data for Urban Travel. Transportation Research Record: Journal of the Transportation Research Board **2526**, 126–135 (2015). https://doi.org/10.3141/2526-14, http://trrjournalonline.trb.org/doi/10.3141/2526-14

7. Eldawy, A., Mokbel, M.F.: Spatialhadoop: A mapreduce framework for spatial data. In: Data Engineering (ICDE), 2015 IEEE 31st International Conference on. pp. 1352–1363. IEEE (2015)

8. Gonzalez, M.C., Hidalgo, C.A., Barabasi, A.L.: Understanding individual human mobility patterns. Nature **453**(7196), 779–782 (2008). https://doi.org/10.1038/nature06958

9. Hoteit, S., Chen, G., Viana, A., Fiore, M.: Filling the gaps: On the completion of sparse call detail records for mobility analysis. In: Proceedings of the Eleventh ACM Workshop on Challenged Networks. pp. 45–50. ACM (2016)

10. Hoteit, S., Secci, S., Sobolevsky, S., Ratti, C., Pujolle, G.: Estimating human trajectories and hotspots through mobile phone data. Computer Networks **64**, 296–307 (2014). https://doi.org/10.1016/j.comnet.2014.02.011, http://dx.doi.org/10.1016/j.comnet.2014.02.011

11. Huang, F., Zhu, Q., Zhou, J., Tao, J., Zhou, X., Jin, D., Tan, X., Wang, L.: Research on the parallelization of the dbscan clustering algorithm for spatial data mining based on the spark platform. Remote Sensing **9**(12), 1301 (2017)

12. Iqbal, M.S., Choudhury, C.F., Wang, P., González, M.C.: Development of origin-destination matrices using mobile phone call data. Transportation Research Part C: Emerging Technologies **40**, 63–74 (2014). https://doi.org/10.1016/j.trc.2014.01.002, http://dx.doi.org/10.1016/j.trc.2014.01.002

13. Lv, Y., Duan, Y., Kang, W., Li, Z., Wang, F.Y.: Traffic flow prediction with big data: a deep learning approach. IEEE Transactions on Intelligent Transportation Systems **16**(2), 865–873 (2015)

14. Malta, T.: National household travel survey 2010. transport malta (2011)

15. Pace, N.: Investigating the Potential of Big Data in the Management of Traffic in Malta (2017)

16. Ranjan, G., Zang, H., Zhang, Z.L., Bolot, J.: Are call detail records biased for sampling human mobility? ACM SIGMOBILE Mobile Computing and Communications Review **16**(3), 33–44 (2012)

17. Steenbruggen, J., Tranos, E., Nijkamp, P.: Data from mobile phone operators: A tool for smarter cities? Telecommunications Policy **39**(3-4), 335–346 (2015). https://doi.org/10.1016/j.telpol.2014.04.001, http://dx.doi.org/10.1016/j.telpol.2014.04.001

18. Toole, J.L., Colak, S., Sturt, B., Alexander, L.P., Evsukoff, A., González, M.C.: The path most traveled: Travel demand estimation using big data resources. Transportation Research Part C: Emerging Technologies **58**, 162–177 (2015). https://doi.org/10.1016/j.trc.2015.04.022, http://dx.doi.org/10.1016/j.trc.2015.04.022

19. Vlahogianni, E.I., Karlaftis, M.G., Golias, J.C.: Short-term traffic forecasting: Where we are and where were going. Transportation Research Part C: Emerging Technologies **43**, 3–19 (2014). https://doi.org/https://doi.org/10.1016/j.trc.2014.01.005, http://www.sciencedirect.com/science/article/pii/S0968090X14000096

20. Wang, M.h., Schrock, S.D.: Feasibility of Using Cellular Telephone Data to Determine the Truckshed of Intermodal Facilities. Tech. rep., University of Nebraska - Lincoln (2012)