

ParPredict: A partially-ordered sequential rules based framework for mobility prediction

Amirat Hanane

Received: date / Accepted: date

Abstract Predicting the future movement of mobile users has emerged as an important technology topic in many applications related to intelligent transportation systems (ITS) and Location-based services (LBS). Numerous prediction models were proposed relying on probabilistic models (e.g. Markov Chain) or data mining techniques (e.g. neural network, sequential patterns mining). Mining sequential patterns and rules is one of the data mining techniques used. Mining sequential rules from sequence databases is an active research topic that is broadly applied for many real-world scenarios. In this paper, we propose to adapt a novel kind of sequential rules called partially order sequential rules for route prediction problem. We aim to further compare this kind with standard sequential rule for the task of mobility prediction. An experimental evaluation conducted on real and synthetic datasets show that the proposed model outperforms a state-of-the-art sequential model in terms of accuracy and prediction coverage.

Keywords Route prediction · ITS · LBS · Partially-ordered · Sequential rules mining

1 Introduction

Nowadays, predicting future movements of vehicles is of great importance for many applications including improving the quality of intelligent transportation systems by providing real-time traffic data and thus allows predicting congestion levels and trip durations. In other words, the prior knowledge of the future movement of vehicles enables estimating future congestion levels and upcoming traffic hazards. Another application is optimizing hybrid fuel consumption. Researchers from Nissan have shown that a hybrid fuel economy

Département informatique
Université Kasdi Merbah, Ouargla, Algeria
E-mail: hanane.amirat@gmail.com

by up to 7.8% could be achieved if the route is known in advance [4]. As for the context of Location-Based services, route prediction can be used for targeted advertising to deliver advertising messages to customers who are likely approaching areas of interest.

Route prediction consists of finding the future road segment for a mobile user. It is mainly based on the assumption that the driving behavior presents a spatial and temporal regularity. For instance, people always tend to take the same routes from home to work at almost a specific time in the morning. Hence, a significant number of a person's trips are repeated, and thus, it is possible to predict that he or she will take that route.

Route prediction has attracted a lot of attention from the research community in the last decade. Most route prediction approaches usually consist of using machine learning techniques such as neural networks or statistical models like Markov model, probabilistic tree, and Bayesian inference. There are two main reasons behind the choice of these techniques for route prediction problem. First, prediction is considered one of the major tasks of data mining. This fact justifies the wide use of data mining prediction techniques to predict future routes. Second, prediction is not precise most of the time where prediction is always performed with a certain probability and confidence, but never reach 100% sure of the prediction. Thus, statistical models have constituted suitable tools in such problem.

Although these techniques and models have shown good performance in prediction, they suffer from two key limitations. First, accurate predictions using data mining techniques mainly depend on defining a set of configuration parameters such as appropriate architecture, function, and input weights such as the in case of neural networks. Second, most of the proposed statistical and probabilistic based approaches assume the markovian hypothesis that the next road segment only depends on the previous segment (in the case of first order Markov model) or requires building a model that is exponentially large if more than one element has to be considered (i.e. k-order Markov models). These assumptions unfortunately do not hold most of the time making route prediction unrealistic or unpractical in many cases.

To accommodate these limitations, we propose, in this paper, to design a route prediction approach based on sequential rules mining. Predicting using sequential rules occurring in a set of training sequences has the advantage of being unsupervised, scalable, noise-tolerant to generate accurate predictions. Moreover, identifying sequential patterns are only order-dependent, meaning additional sets of items could occur between the identified sets without consequence as long as the sets are in that order. The technique is also more robust to noise in terms of still being able to identify the underlying pattern. For example, let us suppose a frequent sequential pattern $P=ABC$ with a frequency $=3$. If there exists an item d which has occurred as noise $ABdC$, the frequency of the pattern P will increase by one resulting a new value ($4 = 3+1$) without an effect of noisy data. Thus, this technique may be useful for uncertain and noisy data. This fact has motivated us the use of sequential patterns for uncertain data such as GPS collected and used for route prediction.

In the context of route prediction, sequential rules mining has been applied in [10] as a core for their prediction model. One drawback regarding this approach is that the idea has been well presented, however; the performance of their proposal was not experimented to study the effectiveness of the approach. Moreover, this approach is generally too specific and restrictive in terms of the order required for visiting locations (they are order-dependent). Hence, a small variation in the order of locations visited by a user may lead to different predictions or the inability to make predictions. For instance, a set of persons (e.g. students) may tend to visit the same places in a city but the sequential order may be slightly different for different persons. Thus, prediction using standard sequential rules may fail if the order of locations visited by a person does not match any of the extracted sequential patterns.

In this paper, we improve on the approach in [10] by adapting a new kind of sequential rules named partially-ordered sequential rules (POSR) [7] for the task of mobility prediction. This kind of sequential rules has proved to greatly improve prediction accuracy while requiring a smaller training set accuracy for the task of webpage recommendation for large clickstream datasets in [5] comparing to the standard sequential rule. The distinctive characteristic of POSR compared to standard sequential rules is that it maintains the occurrence order of all the set of items (i.e. location) in the antecedent of the rule that comes before all the set of items in the consequent while it tolerates the order variation among the set of items in the antecedent and consequent of the rule. Consequently, POSR rules are considered to be more general than standard sequential rules where several standard sequential rules could be represented by a single partially ordered sequential rule.

The remainder of the paper is organized as follows. Section 2 gives some definitions related to our prediction problem and presents a literature review of many existing proposals. Section 3 provides formal definitions required for sequential rule mining and presents partially-ordered sequential rules we are using in this paper. Section 4 describes the architecture of the prediction framework. The evaluation of the proposed model and the experimental results are reported in Section 5. Finally, we conclude our work and give some future works in Section 6.

2 Preliminaries

In this section, the definitions of the key terms required for the comprehension of our approach are provided and the problem addressed in this paper is formulated.

- **Definition 1** (*Road segment*). A road segment with a unique identifier rs_i is represented by a unidirectional edge between two nodes (junctions) [14].
- **Definition 2** (*path, mobility or movement pattern*). A mobility pattern is an abstraction of vehicle locations. A mobility pattern $P = \langle rs_1, rs_2, \dots, rs_n \rangle$

Table 1 A sample of mobility sequences of vehicles.

Vehicle ID	Mobility Sequence
V_1	$rs_0rs_2, rs_3, rs_4, rs_6$
V_2	rs_1, rs_2, rs_3, rs_5
V_3	rs_2, rs_6, rs_7
V_4	rs_2, rs_3, rs_4, rs_6

can be defined as the sequence of road segments traversed by a vehicle during its trip in a specific geographic area within a given map. For instance, Table 1 depicted a sample of mobility pattern of four vehicles $V:V_1, V_2, V_3, V_4$. Each vehicle has its own trip expressed by road segment identifier rs_i . For example, the vehicle V_2 has traversed the road segment rs_1 followed successively by rs_2 , rs_3 and rs_5 .

- **Definition 3** (*Route prediction problem*). Given such a mobility pattern $P = \langle rs_1, rs_2, \dots, rs_n \rangle$, the future route prediction problem is to predict the road segment route that the vehicle will arrive at next.

3 Related work

To predict the future route for a mobile user, several techniques such as data mining, statistical models, and trip matching have been proposed in the literature.

3.1 Statistical models

Most of the proposals relying on the use of statistical models (first-order Markov [11], variable-order Markov [16], Hidden Markov [13], etc.) are based on Markovian assumption. For instance, and to perform personal route prediction, the researchers in [15] have adopted a first-order Markov chain model to build a probability transition matrix containing the probability associated with each link. To deal with fast-growing personal mobility data and challenges in real-time in-vehicle application, data reduction algorithms on the probability matrix have been proposed. The main drawback of this approach is that the dataset used to validate the work contains just one driver which is insufficient to represent the accuracy approved in their work. The authors in [12] have proposed three prediction models (1) statistical model based on mining frequent itemset (2) n-order Markov model where $n \geq 4$ (3) and a Pattern Matching Model adapting Markov model by easing constraints of the number of previous items to be considered in n-order Markov. The best accuracy has been reported by the second model with 70%. As previously mentioned, the main drawback of Markov models is their assumption that the next route for a driver only depends on the actual current route which is not always the case.

However, considering more than one element requires building an exponentially large model (i.e. case of k-order Markov).

3.2 Data mining

For route prediction problem, data mining techniques have been applied as the core component of the prediction approach such as neural network, sequential pattern mining in [9]. Two architectures of neural network have been employed feed-forward [17] and recurrent bidirectional [2]. In [2], the authors have designed a system to predict the destination of a taxi based on the beginning of its trajectory and associated meta-information such as the departure time, the driver identifier, and client information. A recurrent bidirectional neural network has been applied to encode the representation of the taxi's prefix with associated metadata, whereas the mean shift clustering technique has been used to get clusters of destinations of all the training trajectories. Yet, accurate prediction basing on a neural network requires defining a set of input weights, reward function, and suitable architectures that is not a trivial task.

Besides, mining sequential rules from sequence databases is an important and active research topic with broad applications, such as customer shopping transaction analysis, mining weblogs. Many approaches have been proposed to predict the next route of a person using sequential pattern mining. For instance, the researchers in [3] have applied a mining algorithm called CRPM (Continuous Route Pattern Mining). The latter is based on the well-known sequential pattern mining algorithm PrefixSpan [8] to extract route patterns from historical movement data and prediction relies mainly on a pattern tree built from these patterns. In [3], the authors have attempted to predict simultaneously the intended destination and the future route of a person. From real GPS data, the authors have proposed to cluster important places the user may depart from or go to using FBM (Forward-Backward Matching) clustering algorithm. Trajectories are abstracted and then extract the movement patterns using an extended CRPM algorithm. Important factors that must to be considered and fixed when mining sequential patterns and rules are (1) threshold value needed for pattern extracting and (2) the confidence associated with generated rules.

The closer work to our proposal is presented in [9] where the authors have proposed several communication schemes to collect historical vehicular paths. After collecting all the paths and determining the minimum support threshold (*minsup*), the most frequently traveled vehicular paths can be extracted and used as vehicular movement patterns. Relying on a minimum confidence threshold (*minconf*), a set of sequential rules are generated from the extracted patterns and used to forecast a vehicle's future route.

4 ParPredict: a partially-ordered sequential rules for mobility prediction

To address the problem of route prediction, this article proposes to adapt a newly proposed algorithm for mining partially order sequential rules. In this section, we first define some key terms related to POSR mining followed by the architecture of our ParPredict framework.

4.1 Definitions

- **Definition 1** (*sequence database*). A sequence database SDB is a set of tuples of the form $\langle id, S_i \rangle$, where S is a sequence identified by an identifier id . For instance, Table 1 represents a sample of a sequence database of the form $\langle V_i, rs \rangle$ comprising the sequence rs of the set of road segments traversed by the set V of vehicles.
- **Definition 2** (*mobility sequence and mobility subsequence*). Let $RS = \{rs_1, rs_2, \dots, rs_m\}$ be a finite alphabet of items (symbols here also called road segments). A sequence $S = \langle rs_1, rs_2, \dots, rs_n \rangle$ is an ordered set of road segments $rs_k \in RS$. A sequence $S' = \langle rs_j, rs_{j+1}, \dots, rs_k \rangle$ is said to be a subsequence of S if and only if there exists $j \leq k \leq n$. For instance, consider the sequence of database of table 1, the sequence $\langle rs_1, rs_2 \rangle$ is a subsequence of three sequences of vehicles V_1, V_2, V_4 .
- **Definition 3** (*Sequential pattern and frequent sequential mobility pattern*). A sequential pattern is defined as a subsequence of many sequences in the sequence database. A *frequent sequential pattern* is a sequential pattern X that is repeated a number of time over the total number of sequences that exceeds a predefined threshold $minsup$ ($support(X) \geq minsup$). For instance, the sequences $S_1 : \langle rs_1, rs_2 \rangle$ and $S_2 : \langle rs_3, rs_4 \rangle$ are two sequential mobility patterns, however, only $\langle rs_1, rs_2 \rangle$ is frequent if $minsup=3$.
- **Definition 4** (*Sequential rule*). A sequential rule $SR : X \Rightarrow Y$ is the relationship between two bility patterns $X, Y \subseteq RS$ such that $X \cap Y = \emptyset$ and $X, Y \neq \emptyset$. X is called the antecedent whereas Y is the consequent of SR .
- **Definition 5** (*Support*). The support of a pattern X is defined as the number of sequences from SDB where the items of X occurs, divided by the number of sequences in SDB.
- **Definition 6** (*Confidence*). Given $SR : X \Rightarrow Y$ a sequential rule, the confidence of SR is computed as:

$$Conf(SR) = \frac{Support(X \cup Y)}{Support(Y)}$$

- **Definition 7** (*Partially-ordered sequential rule*). Let $PR : P \Rightarrow L$ be a sequential rule. PR is said to be partially ordered if the items of P occur in a sequence (in any order), the items in L will occur afterward in the same

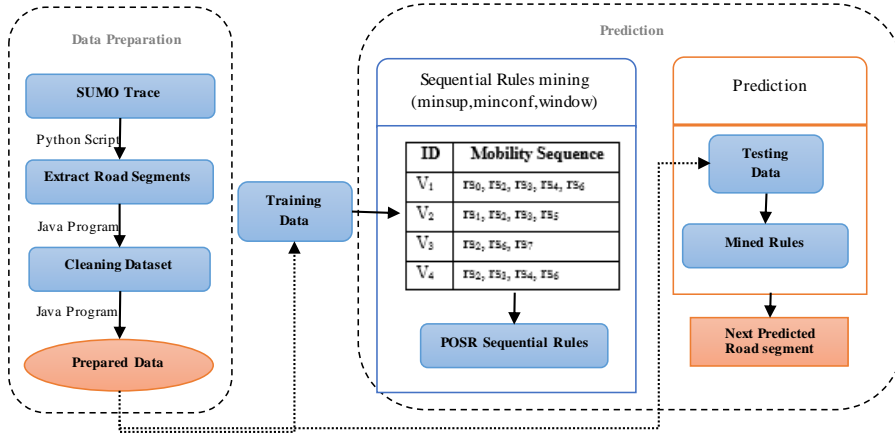


Fig. 1 Architecture of ParPredict.

sequence (in any order). In other words, the requirements of a sequential ordering inside the antecedent X and inside the consequent Y of the rule are ignored but the requirement of a sequential relationship between the antecedent and consequent of a rule is preserved

4.2 Framework Architecture

As depicted in Figure 1, ParPredict consists of two key steps: pre-processing (or data preparation) and prediction.

4.2.1 Pre-processing (data preparation)

In this pre-processing step, *ParPredict* generates the movement sequence database to be employed in prediction. In our case, two types of mobility datasets were used: synthetic and real GPS mobility data.

- **Synthetic dataset.** Using a SUMO (Simulation of Urban MObility) scenario that contains the ".rou" file that defines the realistic route description and the ".cfg" configuration file, ParPredict first generates the vehicles' mobility trace. It then cleans the trace using python programs to remove useless information in order to get and only preserve the path of each vehicle. The output of this step is a set of route segment ID traversed by each vehicle (i.e. mobility patterns).

1. **Generate the mobility dataset.** By executing a python script (networkdump.py) to the SUMO configuration file (.cfg), an XML output file containing detailed information about each vehicle (edge ID, lane ID, etc.) for each timestamp is generated.

road segment rs_4 as the susceptible road segment that a vehicle will pass through.

5 Experimental Evaluation

We have conducted a comparative study between partially and standard sequential rules. This section explains the experimental settings and discusses the obtained results.

Experimental setting This section describes the experimental study we have conducted to evaluate the performance of *ParPredict* framework and compare it with the system proposed in [9]. For sake of simplicity, we call the latter *StandPredict*. The evaluation experiments were performed on a computer equipped with a Core i7-4500 CPU, 16GB of RAM, and 1TB of Hard Disk. Java implementations the RuleGen [18] and TRuleGrowth [5] sequential rules mining algorithms were used. To well measure the performance of our proposal, the synthetic trace called *Lust* and the realist taxi driver datasets were used. *Lust* is a large-scale vehicular mobility model that allows generating vehicle driving traces based on real traffic volume counts and Luxembourg map. This mobility trace is available as a SUMO scenario and it is prepared as presented in section 4.2.1. The taxi cabs dataset comprises the GPS coordinates of approximately 500 taxis collected over 30 days in the San Francisco Bay Area. To generate mobility sequences from this dataset, the steps presented in [1] were followed.

Each dataset is randomly divided into a training set and a testing set by setting a training ratio parameter. The training set is used to generate sequential rules that are used thereafter to perform prediction for each testing sequence. In all the conducted experiments, the training ratio is set to 0.7 (70%).

5.1 Evaluation metrics

To measure the performance of *ParPredict*, the two following metrics widely employed for prediction proposes, are used.

- **Overall Accuracy.** It can be defined as the number of future routes successfully predicted to the total number of sequences in the testing set.

$$OverallAccuracy = \frac{(\text{Number of successful predictions})}{(\text{Number of Testing sequences})}$$

- **Coverage.** The coverage determines the number of sequences where a prediction has been performed by the number of testing sequences. This measure indicates whether a matching rule is found for a mobility sequence or not.

$$Coverage = \frac{(\text{Number of matching rules})}{(\text{Number of Testing sequences})}$$



Fig. 3 Scalability of ParPredict.

5.2 Experiments and results

5.2.1 Experiment 1: scalability

In this experiment, we aim to study the scalability of *ParPredict* by increasing the number of mobility sequences using TaxiCab and Lust datasets for *minsup* and *minconf* thresholds set both to 10%. As depicted in Figure 3, it could be found that prediction performance is decreased when the number of training mobility sequences is increased. This is reasonable because as more location data is involved, the less frequent are the mobility patterns as more users are included in data. This finding is well observable with the TaxiCab dataset. Results also show that *ParPredict* outperforms *StandPredict* for all cases.

5.2.2 Experiment 2: effect of varying *minsup*

Figure 4 exhibits the results obtained in the second experiment that compares the performance of *ParPredict* with *StandPredict* on the TaxiCab dataset for different *minsup* values. For both sequential rules mining models, as *minsup* is increased, results become worse. This may be due to the fact that by increasing *minsup*, less frequent patterns are found but *ParPredict* gives good results as POSR tolerate order variations in mobility patterns.

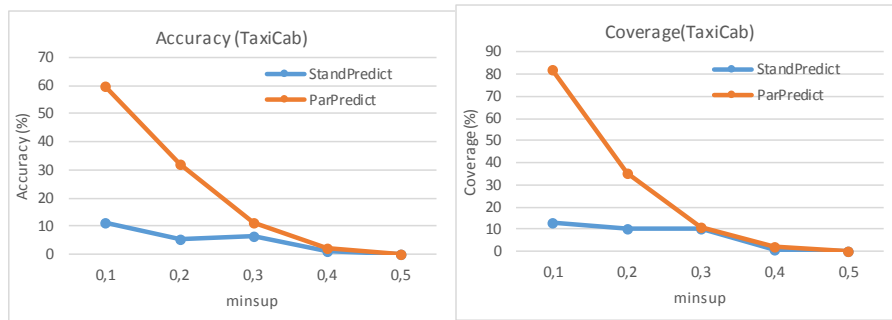


Fig. 4 Impact of varying minsup threshold on prediction.

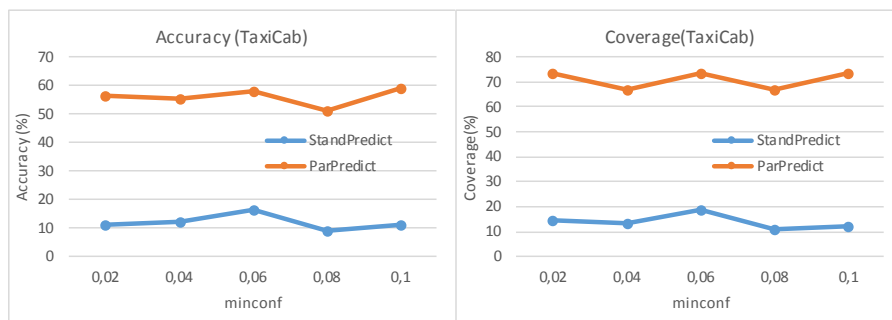


Fig. 5 Impact of varying minconf threshold on prediction.

5.2.3 Experiment 3: effect of varying minconf

Figure 5 checks the influence of varying *minconf* threshold on performance. From the results, it can be noticed that increasing *minconf* slightly improves the performance for all sequential rules mining predictors.

To summarize and from all the above obtained results, it can be concluded that *ParPredict* has presented good results comparing to standards sequential rules mining based models in terms of all evaluation metrics. Besides, due to the limit of pages number, many experiments were not included in the paper and they will be a subject of an extended version of this work such as studying the number of rules generated, the spatial and temporal complexity of both models, etc.

6 Conclusion

This paper has proposed a mobility predictor system called *ParPredict*. The latter is a sequential rules-based framework. It initially mines partially ordered

sequential rules then employs these rules to predict new locations to drivers. Unlike standard sequential-based predictors, the proposed framework does not require a strict ordering of locations in the learned rules that lead to more accurate predictions. The experimental study was carried out using two mobility datasets. Results have shown that *PaPredict* has good performance compared to a state-of-the-art model. For future work, we plan to (1) extend *ParPredict* so it considers the temporal factor of mobility, and (2) remove redundancy in the mined rules by discovering a compact set of rules.

References

1. Amirat, H., Lagraa, N., Fournier-Viger, P., Ouinten, Y.: Nextroute: a lossless model for accurate mobility prediction. *J. Ambient Intell. Humaniz. Comput.* **11**(7), 2661–2681 (2020). DOI 10.1007/s12652-019-01327-w. URL <https://doi.org/10.1007/s12652-019-01327-w>
2. de Brébisson, A., Simon, É., Auvolet, A., Vincent, P., Bengio, Y.: Artificial Neural Networks Applied to Taxi Destination Prediction (2015). URL <http://arxiv.org/abs/1508.00021>
3. Chen, L., Lv, M., Chen, G.: A system for destination and future route prediction based on trajectory mining. *Pervasive and Mobile Computing* **6**(6), 657–676 (2010). DOI 10.1016/j.pmcj.2010.08.004. URL <http://dx.doi.org/10.1016/j.pmcj.2010.08.004>
4. Deguchi, Y., Kuroda, K., Shouji, M., Kawabe, T.: HEV Charge / Discharge Control System Based on Navigation Information. In: *Convergence International Congress & Exposition On Transportation Electronics*, vol. 1 (2004)
5. Fournier-Viger, P., Gueniche, T., Tseng, V.S.: Using partially-ordered sequential rules to generate more accurate sequence prediction. In: *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 7713 LNAI, pp. 431–442. Springer Berlin Heidelberg (2012)
6. Fournier-Viger, P., Lin, J.C.W., Gomariz, A., Gueniche, T., Soltani, A., Deng, Z., Lam, H.T.: The SPMF Open-Source Data Mining Library Version 2. pp. 36–40. Springer, Cham (2016). DOI 10.1007/978-3-319-46131-1_8. URL http://link.springer.com/10.1007/978-3-319-46131-1_8
7. Fournier-Viger, P., Wu, C.W., Tseng, V.S., Cao, L., Nkambou, R.: Mining Partially-Ordered Sequential Rules Common to Multiple Sequences. *IEEE Transactions on Knowledge and Data Engineering* **27**(8), 2203–2216 (2015). DOI 10.1109/TKDE.2015.2405509. URL <http://ieeexplore.ieee.org/document/7045582/>
8. Jian Pei, J., Jiawei Han, J., Mortazavi-Asl, B., Jianyong Wang, J., Pinto, H., Qiming Chen, Q., Dayal, U., Mei-Chun Hsu, M.C.: Mining sequential patterns by pattern-growth: the PrefixSpan approach. *IEEE Transactions on Knowledge and Data Engineering* **16**(11), 1424–1440 (2004). DOI 10.1109/TKDE.2004.77. URL <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=1339268>
9. Merah, A.F., Samarah, S., Boukerche, A.: Vehicular movement patterns: A prediction-based route discovery technique for VANETs. *IEEE International Conference on Communications* pp. 5291–5295 (2012). DOI 10.1109/ICC.2012.6364141
10. Merah, A.F., Samarah, S., Boukerche, A., Mammeri, A.: A sequential patterns data mining approach towards vehicular route prediction in VANETs. *Mobile Networks and Applications* **18**(6), 788–802 (2013). DOI 10.1007/s11036-013-0459-6
11. Neto, F.D.N., Baptista, C.D.S., Campelo, C.E.C.: Prediction of Destinations and Routes in Urban Trips with Automated Identification of Place Types and Stay Points pp. 80–91 (2015)
12. Petróczi, A.I., Gáspár-Papanek, C.: Route prediction on tracking data to location-based services. *Lecture Notes in Computer Science* **5733 LNCS**, 69–77 (2009). DOI 10.1007/978-3-642-03700-9₈

13. Qiu, D., Papotti, P., Blanco, L.: Future locations prediction with uncertain data. *Lecture Notes in Computer Science* **8188 LNAI(PART 1)**, 417–432 (2013). DOI 10.1007/978-3-642-40988-2_27
14. SUMO: Simulation of Urban MObility download — SourceForge.net. URL <https://sourceforge.net/projects/sumo/>
15. Wang, X., Ma, Y., Di, J., Murphey, Y.L., Qiu, S., Kristinsson, J., Meyer, J., Tseng, F., Feldkamp, T.: Building efficient probability transition matrix using machine learning from big data for personalized route prediction. *Procedia Computer Science* **53(1)**, 284–291 (2015). DOI 10.1016/j.procs.2015.07.305. URL <http://dx.doi.org/10.1016/j.procs.2015.07.305>
16. Xue, G., Li, Z., Zhu, H., Liu, Y.: Traffic-known urban vehicular route prediction based on partial mobility patterns. *Proceedings of the International Conference on Parallel and Distributed Systems - ICPADS* pp. 369–375 (2009). DOI 10.1109/ICPADS.2009.129
17. Ye, Q., Chen, L., Chen, G.: Predict personal continuous route. *IEEE Conference on Intelligent Transportation Systems, Proceedings, ITSC* pp. 587–592 (2008). DOI 10.1109/ITSC.2008.4732585
18. Zaki, M.: Scalable algorithms for association mining. *IEEE Transactions on Knowledge and Data Engineering* **12(3)**, 372–390 (2000). DOI 10.1109/69.846291. URL <http://ieeexplore.ieee.org/document/846291/>