

An online map matching algorithm based on second-order hidden Markov model

Xiao Fu^{1*}, Jiaxu Zhang², Yue Zhang³

¹ Jiangsu Key Laboratory of Urban ITS, Jiangsu Province Collaborative Innovation Center of Modern Urban Traffic Technologies, School of Transportation, Southeast University, Nanjing, China

² State Key Laboratory of Information Engineering in Surveying, Mapping and Remote Sensing, Wuhan University, Wuhan, China

³ School of Geographic Sciences, East China Normal University, Shanghai, China

Correspondence should be addressed to Xiao Fu; fuxiao@seu.edu.cn

Abstract

Map matching is a key pre-process of trajectory data which recently have become a major data source for various transport applications and location-based services. In this paper, an online map matching algorithm based on second-order hidden Markov model (HMM) is proposed for processing trajectory data in complex urban road networks such as parallel road segments and various road intersections. Several factors such as driver's travel preference, network topology, road level and vehicle heading are well considered. An extended Viterbi algorithm and a self-adaptive sliding window mechanism are adopted to solve the map matching problem efficiently. To demonstrate the effectiveness of the proposed algorithm, a case study is carried out using a massive taxi trajectory dataset in Nanjing, China. Case study results show that the accuracy of the proposed algorithm outperforms the baseline algorithm built on the first-order HMM in various testing experiments.

Introduction

With the development of positioning and wireless communication technologies, floating car data (e.g. trajectories of taxis) have become major data source for many applications such as location-based services, intelligent transportation systems and transport policy appraisals [1-5]. The errors of positioning data collected by global positioning system (GPS) equipment on floating vehicles are inevitable and could come from satellite, transmission process and receiver [6]. Map matching is the process of matching GPS data with errors onto the road network in order to eliminate the impact of errors and maximize the effectiveness of data. In practical applications, map matching algorithm plays a vital role. For example, travel time prediction based on floating car data, which needs to match GPS points to the corresponding road segment accurately. Therefore, the map matching algorithm is the basis for the large-scale application of floating car data.

Existing map matching algorithms can be divided into four categories based on the technique they adopted [7]: geometric technique [8-10], topological technique [11-14], probability statistics technique [15] and integrations of multiple technologies [16-18]. The algorithms with

Table 1: Comparison of HMM-based map matching studies.

Category	Model	Sliding window	Factors of observation probability	Factors of transition probability	Order of HMM
Offline	Interactive-Voting Matching [24]	×	×	Speed constraint	First-order
	HMM Matching [23]	×	×	Distance difference	
	Unsupervised HMM [27]	×	Antenna location	Speed constraint	
	Quick Matching [25]	×	Speed constraint	×	
	Multistage Matching [26]	×	Vehicle heading	Heading change difference	
	Driver Path Preference Based HMM [28]	×	×	Distance difference; Speed constraint; Driver's travel preference	
Online	SnapNet [30]	Fixed size	Vehicle heading; Road level	Distance difference; Network topology (same road priority)	Second-order
	Spatial and Temporal Matching [29]	Fixed size	×	Speed constraint	
	Route Choice HMM [31]	Variable size	×	Distance difference; Free-flow travel time	
	This Study	Self-adaptive size	Vehicle heading; Road level; Driver's travel preference	Distance difference; Network topology (same/adjacent road priority)	

geometric technique utilize geometric information of GPS point and road network (e.g. distance, angle and shape) without considering the topology of road network. These algorithms show high efficiency of map matching, but the accuracy is low when matching low-precision GPS data to complex road network. As regards topological technique, both geometric factors and road topology are considered. To some extent, topological technique improves the matching accuracy, but is still vulnerable to the influence of low-frequency sampling interval and large sampling noise. The probability statistics technique sets an ellipse or rectangle confidence area for each GPS point, thus we can obtain the probability according to the distance between the GPS point and the position in confidence area. Optimal matching paths are determined according to values of the probability. Compared to geometric technique and topological technique, the probability statistics technique is relatively more complex and difficult to implement, and shows low time efficiency. By combining geometric, topological and probability factors, advanced techniques, such as Kalman filter [19], Bayesian filter [20], fuzzy logic model [21], multi-hypothesis tree [18] and hidden Markov model (HMM) [22], can effectively improve the map matching accuracy and achieve online incremental matching.

Of the advanced techniques, HMM has become popular in map matching studies. HMM is a prevailing paradigm of network-based dynamics modelling, which well suits the process of finding the most suitable matching point (i.e. hidden state) to each GPS point (i.e. observed state) on the road network in map matching problem. Existing map matching algorithms based on HMM can be categorized into two categories [20]: offline algorithms and online algorithms (refer to Table 1).

Offline HMM map matching algorithms are applied using historical data, batching the whole input trajectory to find the optimal matching path in the road network [23-28]. Whole trajectories enable offline algorithms to take account of the relationship between the front and

the back points to achieve higher accuracy. Offline algorithms show robustness to the reduction of sampling rate, but the computation efficiency is low. Online algorithms estimate the current segment immediately after obtaining GPS data, and this kind of algorithm can be used for providing online services such as real-time navigation and trajectory monitoring. Because of the unavailability of future points, online algorithms are more complicated and require higher computation demand for real-time applications. Most studies utilize the sliding window mechanism with fixed window size to realize online matching [29, 30]. As the number of GPS points increases, the points in sliding window change dynamically. However, under the condition of low data quality or complex road network, small window leads to a significant decrease in matching accuracy while large window brings a significant decrease in computation efficiency. A few online map matching algorithms adopt variable sliding window but it requires a lot of extra computation [31]. Considering these, in this study, we proposed self-adaptive sliding window to realize online map matching based on HMM, which promises accuracy and efficiency at the same time.

HMM build on the stochastic processes of observation and state transition. In the map matching context, two probabilities are important, called observation probability and transition probability. Observation probability is usually obtained by the Gaussian distribution of great-circle distance between GPS points and candidate points. In the literature, several factors have been considered in calculating observation probability. For instance, unsupervised HMM [27] considers the location of Antenna when matching mobile phone data. Other studies, e.g. Quick Matching [25], Multistage Matching [26] and SnapNet [30] consider more factors including the speed constraint, road level and vehicle heading. As regards transition probability calculation, to consider temporal relationship of different points, some factors such as speed constraint and free-flow travel time are considered in several studies [24, 27-29, 31]. To consider spatial relationship, some factors are included such as the difference between great-circle distance and route distance [23, 28, 30, 31], difference between vehicle's heading change and road segments' heading change [26] and same road priority [30]. Based on the analysis of advantages of each algorithm, this study is a pioneering endeavour devoted to comprehensively considering various factors in online map matching, i.e. road level, driver's travel preference, vehicle heading and network topology (same/adjacent road priority).

To the best of our knowledge, almost all map matching algorithms based on HMM adopt first-order HMM. The basic hypothesis of first-order HMM is that the observation probability is only related to the current state while the transition probability is only related to the previous state. Because the moving of a vehicle is a continuous process, there is a complex space-time relationship between the current state and the previous states. There is no doubt that first-order HMM over-simplifies several practical systems. Recently, Salnikov et al. [32] explored possibilities to enrich the system description and exploited empirical pathway information by means of second-order Markov models. Experiments show that the higher order model is more effective than the first-order model in dealing with space-time continuum. Therefore, a need is likely to exist for solving map matching problem using higher-order (e.g. second-order) HMM to achieve better map matching results.

Along the line of previous online studies, this study proposes a new map matching algorithm based on the HMM technique. The proposed algorithm extends the previous studies in following aspects. Firstly, the proposed novel map matching algorithm is on the basis of second-order HMM, which can better consider the space-time relationship among different states. It can be effectively applied to complex urban road network with parallel segments using low-frequency sampling GPS data. Secondly, the proposed algorithm comprehensively considers

driver's travel preference towards road segments, road level, vehicle heading and network topology when calculating the probability matrix of second-order HMM in order to improve the matching accuracy. Thirdly, the proposed algorithm introduces a self-adaptive sliding window mechanism. Compared to the conventional fixed window size mechanism, the introduced mechanism using a self-adaptive window size can significantly improve the map matching accuracy while has a reasonable computational performance.

In summary, the contributions of this work are three folds:

- An online map matching algorithm based on second-order hidden Markov model (HMM) is proposed, which can better consider the spatial-temporal relationship among different states and large perception fields.
- The proposed algorithm comprehensively considers driver's travel preference, road level, vehicle heading and network topology when calculating the probability matrix of second-order HMM to improve the matching accuracy.
- Experiments on real-world dataset show that with the help of the self-adaptive sliding window mechanism and an extended Viterbi algorithm, our second-order HMM based model can reach a high accuracy while ensuring efficiency.

The rest of this paper is organized as follows. In the next section, we state the problem of map matching. After the problem statement, an online map matching algorithm is proposed based on second-order HMM. A case study is carried out using a large taxi trajectory dataset in Nanjing, China to test the validity of the algorithm under various road conditions. Finally, we conclude this study and discuss directions for further research.

Problem Statement

Vehicle trajectory data are a series of GPS points recorded in chronological order. Each GPS point indicates longitude and latitude, vehicle speed, timestamp, etc. Because the errors of data collected by GPS equipment are inevitable, map matching is a key process before using the vehicle trajectory data. It is a process of matching GPS data onto the road segments and obtaining the continuous and specific locations of vehicles on the road. The concepts used in this study are listed as follows:

GPS point: A GPS point g_t is a record indicating the longitude, latitude, time stamp and velocity of the vehicle.

GPS trajectory: A GPS trajectory T is a series of GPS points. A T is showed as: $g_1 \rightarrow g_2 \rightarrow \dots \rightarrow g_n$.

Road network: Road network $G(V, E)$ is a directed graph where V is the set of vertexes and E is the set of edges.

Road segment: A road segment e is a directed edge in road network with length, road level, start vertex and end vertex.

Candidate point: The candidate point c_t^n is the n th candidate point matched with GPS point g_t on the road network.

Route: A route R is a sequence of road segments that matched best to a GPS trajectory T , each road segment belongs to the edge set E of road network $G(V, E)$. R is showed as: $e_1 \rightarrow e_2 \rightarrow \dots \rightarrow e_n$.

With the above concepts, the map matching problem solved in this study can be defined as:
Find the candidate points $c_t^1, c_t^2, \dots, c_t^n$ on each road segment e corresponding to GPS point g_t .
Select the most likely candidate points sequence for GPS trajectories T , and connect the matched road segments on network G to get route R .

Second-order HMM map matching

Data Pre-processing

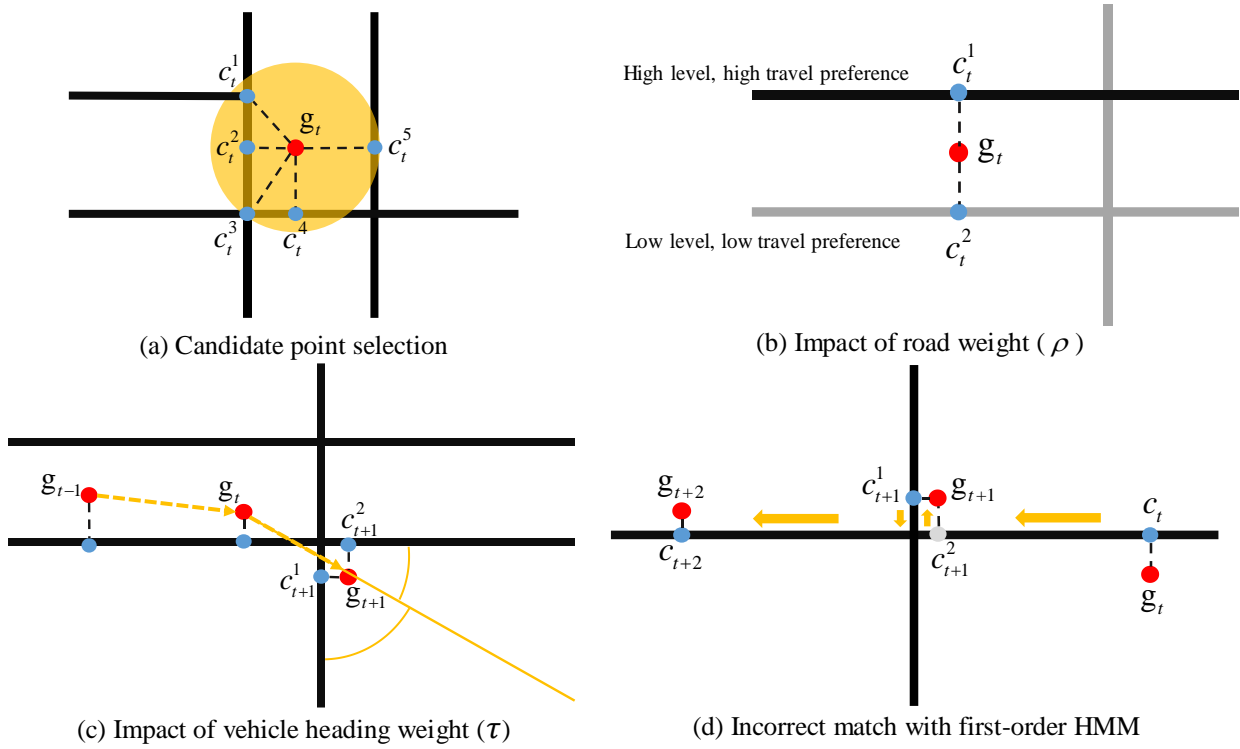


Figure 1: Illustration of merits of the proposed map matching method.

Generally, there are a lot of "redundancy" and "incompleteness" in floating vehicle GPS data, which may be caused by devices or road environments (e.g. stopping in or passing through tunnels). In order to ensure the efficiency and accuracy of map matching, we first need to pre-process the GPS data, including the removal of redundant data and the interpolation of missing data.

For the currently received data point g_t , calculate the great-circle distance [23] of g_t and g_{t-1} (denoted as $D_{t-1,t}$), if $D_{t-1,t}$ is less than a pre-defined lower bound, the current point g_t is omitted and not matched. If $D_{t-1,t}$ is greater than an upper bound, the two points will be interpolated linearly.

With the data pre-processing, the redundant GPS data points can be effectively eliminated to avoid unnecessary matching. At the same time, interpolation of two points with too large intervals helps to process low-frequency GPS data.

Candidate Point Selection

For the currently received data point g_t , we search for its candidate points (refer to Figure 1(a)) with the following steps:

Step1: Using the R-tree index, the road segments within a pre-defined error circle or nearest to the point g_t are selected as road segment candidates [13, 17].

Step2: Project g_t vertically on the candidate road segments, and the projection point c_t^i is a candidate point for g_t . If the projection point falls outside the segment, choose the closer vertex of the segment as c_t^i . As shown in Figure 1(a), the candidate points for g_t are $c_t^1, c_t^2, \dots, c_t^5$. The distances from g_t to the candidate points are denoted as $d_t^1, d_t^2, \dots, d_t^5$ respectively.

Observation Probability

In first-order HMM, the observation probability is used to measure the probability of getting some kinds of observed value in a hidden state [34]. The map matching algorithms based on HMM usually regard the GPS point g_t as the observation value of state t , and the actual position of g_t as the hidden value of state t . The observation probability is modelled using a Gaussian distribution for GPS trajectories. The first-order HMM observation probability in this paper is obtained as

$$P(g_t | c_t^i) = \frac{1}{\sqrt{2\pi}\sigma_t} e^{-0.5 \left(\frac{\tau \cdot \rho \cdot d_t^i}{\sigma_t} \right)^2}, \quad (1)$$

Where $P(g_t | c_t^i)$ is the observation probability of the candidate point c_t^i on g_t . d_t^i is the great-circle distance between g_t and the candidate point c_t^i . σ_t is the standard deviation of a Gaussian random variable that corresponds to the average great-circle distance between g_t and its candidate points. τ is a weight given on vehicle heading, which is related to the road direction angle α_{road} and the trajectory direction angle α_{GPS} :

$$\tau = \nu + \frac{e^{\left| \alpha_{road} - \alpha_{GPS} \right|}}{e^{2/\pi}}. \quad (2)$$

In Equation (2), the road direction angle α_{road} is the direction angle of the two vertexes of a segment. The trajectory direction angle α_{GPS} indicates the direction angle of the last GPS point and the current GPS point. Because of the bidirectional property of the road, there are

two results of $|\alpha_{road} - \alpha_{GPS}|$, and the smaller value of the two results should be used. ν is a parameter which can be estimated with real data.

ρ is a weight reflecting the effect of road including road level (denoted as $rlevel$) and driver's travel preference for the road segment (denoted as $plevel$):

$$\rho = 1 - \mu(rlevel + plevel), \quad (3)$$

where μ is a parameter to be estimated. In this study, $rlevel$ is within $[0, 5]$. A high $rlevel$ indicates a high level of road. The value of $plevel$ is also ranging from 0 to 5. Considering driver's travel experience as a sigmoid curve [33], $plevel$ can be derived as

$$plevel = \frac{5}{1 + e^{-\varpi + \varpi'}}, \quad (4)$$

where ϖ is the actual number of times drivers pass the road segment in a certain time period, and ϖ' is a pre-defined expected number.

In this way, the observation probability can be obtained. By using vehicle heading weight τ and road weight ρ , we can consider road level, driver's travel preference and heading of the floating vehicle at that time, which are significant in online map matching with limited information. Take Figure 1(b) as an example to illustrate the merit of road weight ρ . The current GPS point g_t is located in the middle of two parallel road segments. The distances from g_t to c_t^1 and c_t^2 are the same. In conventional map matching methods, c_t^1 or c_t^2 is selected randomly as the real position of vehicle. However, if road level and travel preference are taken into account using our proposed method, we can consider c_t^1 as the real position of vehicle. It can be seen that without subsequent GPS points, we must make full use of the information provided by existing GPS points and road network in order to improve the matching accuracy.

Figure 1(c) shows the merits of incorporating vehicle heading weight τ . The GPS point g_{t+1} is located near the intersection, which is close to the candidate point c_{t+1}^1 and c_{t+1}^2 , and the distance d_{t+1}^1 is the same as d_{t+1}^2 . Connecting g_t and g_{t+1} , the vehicle heading weight between the connecting line and the two segments is τ_1 and τ_2 . Considering the impact of vehicle heading weight, c_{t+1}^2 has a greater probability of observation, and we can suppose that c_{t+1}^2 is the real position of the vehicle at time $t+1$.

Transition Probability

In first-order HMM, the transition probability measures the transition from one hidden state to another [34]. The map matching algorithm based on HMM uses the transition probability to measure the probability of moving from a candidate point c_{t-1}^i at time $t-1$ to a candidate point

226 c_t^j at time t [30]. The formula for calculating the transition probability of the first-order HMM
 227 in this paper is given as Equation (5):

$$228 \quad P(c_t^j | c_{t-1}^i) = \begin{cases} p_{same} \frac{1}{\beta} e^{-\frac{s_t}{\beta}}, & c_t^j \text{ and } c_{t-1}^i \text{ are on the} \\ & \text{same/adjacent road segments} \\ (1-p_{same}) \frac{1}{\beta} e^{-\frac{s_t}{\beta}}, & \text{otherwise} \end{cases} \quad (5)$$

229 where $P(c_t^j | c_{t-1}^i)$ is the transition probability from candidate point c_{t-1}^i to candidate point
 230 c_t^j . p_{same} (>0.5) is a parameter. With Equation (5), we can get the transition probability with
 231 explicit consideration of network topology (i.e. considering if c_{t-1}^i and c_t^j are on the same or
 232 adjacent road segments). In this way, the topological relation of road segments is taken into
 233 account. β is the mean of s_t . s_t is the difference between the great-circle distance from g_{t-1}
 234 to g_t (denoted as $dist(g_{t-1}, g_t)$) and the route length from c_{t-1}^i to c_t^j (denoted as
 235 $routeDist(c_{t-1}^i, c_t^j)$):

$$236 \quad s_t = \left| dist(g_{t-1}, g_t) - routeDist(c_{t-1}^i, c_t^j) \right|. \quad (6)$$

237 Self-adaptive Sliding Window and Second-order Probability

238 Existing first-order HMM online map matching algorithms usually only focus on one single
 239 GPS point, considering its local geometric relation and road topology, which results in the
 240 precision of online map matching algorithm far behind the second-order map matching
 241 algorithm.

242 Figure 1(d) shows an example that the conventional first-order HMM online map matching
 243 results in incorrect match. Obviously, from GPS point g_t to g_{t+2} , the vehicle does not turn and
 244 the correct matching path should be $c_t \rightarrow c_{t+1}^2 \rightarrow c_{t+2}$. However, in the process of first-order
 245 HMM online incremental matching, an incorrect matching result is $c_t \rightarrow c_{t+1}^1 \rightarrow c_{t+2}$. The
 246 reason for this error is that the first-order HMM only considers the observation probability of
 247 a single point and the transition probability between two points. However, the measurement of
 248 transition probability should be on a larger scale. The real location of the current GPS point is
 249 not just related to the previous point, but to multiple previous points.

250 Higher-order HMM is an extension of first-order HMM [35]. The basic assumption of higher-
 251 order HMM is that the current state is not only related to one previous state, but to multiple
 252 previous states. In some cases, the second-order HMM is more consistent with the real situation,
 253 such as natural language processing, speech recognition and so on [36, 37]. For the map
 254 matching problem, because the vehicle movement is continuous, the real position of the current
 255 point is not only related to the previous point, but also to the trajectory formed by two or more
 256 points. Therefore, higher-order HMM is somewhat more suitable for map matching than
 257 traditional first-order HMM. Analogous to human eyes observing things, we should first pay

attention to the characteristics of things as a whole. For example, in the Figure 1(d), the connection from g_t to g_{t+2} is approximately a straight line, so the GPS point g_{t+1} is more likely to be matched to c_{t+1}^2 than c_{t+1}^1 . To overcome the matching errors which may resulting from first-order HMM and improve the accuracy of online map matching, in this study, we extend the first-order HMM map matching to a second-order one. Compared to the first-order HMM, the difficulties in using second-order HMM lie in the design of the probability matrix and how to improve the computational efficiency.

In the applications such as real-time navigation and travel time estimation, online map matching is necessary. The existing HMM map matching algorithms usually use sliding window to realize online matching. Denote the sliding window size as w (i.e. number of GPS points). If the window overflows after the current point g_t entering the window, the first point in the window g_{t-w} is removed, and the matching result of g_{t-w} point will be finally determined. As the new point continues to join, matching results within the window may be changed continuously. The introduction of sliding window makes online map matching possible, but it is difficult to determine the window size w . If w is too large, the matching speed will be too slow to meet the real-time performance requirement. If w is too small, the matching accuracy will be compromised. To solve this problem, a self-adaptive sliding window is proposed in this study.

In this study, we consider different sizes of self-adaptive sliding window. By calculating the average value of GPS points positioning error in the current window, sliding windows of different sizes are automatically selected to adapt to the current GPS positioning error, which can improve the accuracy of the online map matching as much as possible. The average value of GPS points positioning error (denoted as E_{ave}) can be obtained as

$$E_{ave} = \frac{\sum_{n=t-w+1}^{n=t} dist(g_n, c_n)}{w}, \quad (7)$$

where c_n is the candidate point which is matched to g_n .

Observation probability of second-order HMM $P(g_{t-1}, g_t | c_{t-1}^i, c_t^j)$ can be obtained from the first-order HMM:

$$P(g_{t-1}, g_t | c_{t-1}^i, c_t^j) = P(c_t^j | c_{t-1}^i) \cdot P(g_{t-1} | c_{t-1}^i) \cdot P(g_t | c_t^j). \quad (8)$$

Define the second-order HMM state transition probability (denoted as $P(c_t^i | c_{t-2}^j, c_{t-1}^k)$) as

$$P(c_t^i | c_{t-2}^j, c_{t-1}^k) = \frac{1}{\lambda} e^{-\frac{k_t}{\lambda}}, \quad (9)$$

where λ is the mean of k_t . k_t is the difference between the great-circle distance from g_{t-1} to g_{t+1} and the route length from c_{t-1}^i to c_{t+1}^j :

$$k_t = \left| \sum_{n=t-2}^{n=t-1} \text{dist}(g_n^i, g_{n+1}^j) - \sum_{n=t-2}^{n=t-1} \text{routeDist}(c_n^i, c_{n+1}^j) \right|. \quad (10)$$

The second-order transition probability describes the state transition between three consecutive candidate points, that is, the actual position of the current GPS point is related to the previous two points. In this way, the strong assumption of first-order HMM is relaxed and the accuracy of map matching is improved. In fact, we can continue to extend the proposed method to third-order HMM and define appropriate observation and transition probabilities to improve accuracy. However, three-order HMM will make the calculation process more complicated, which is not conducive to online map matching.

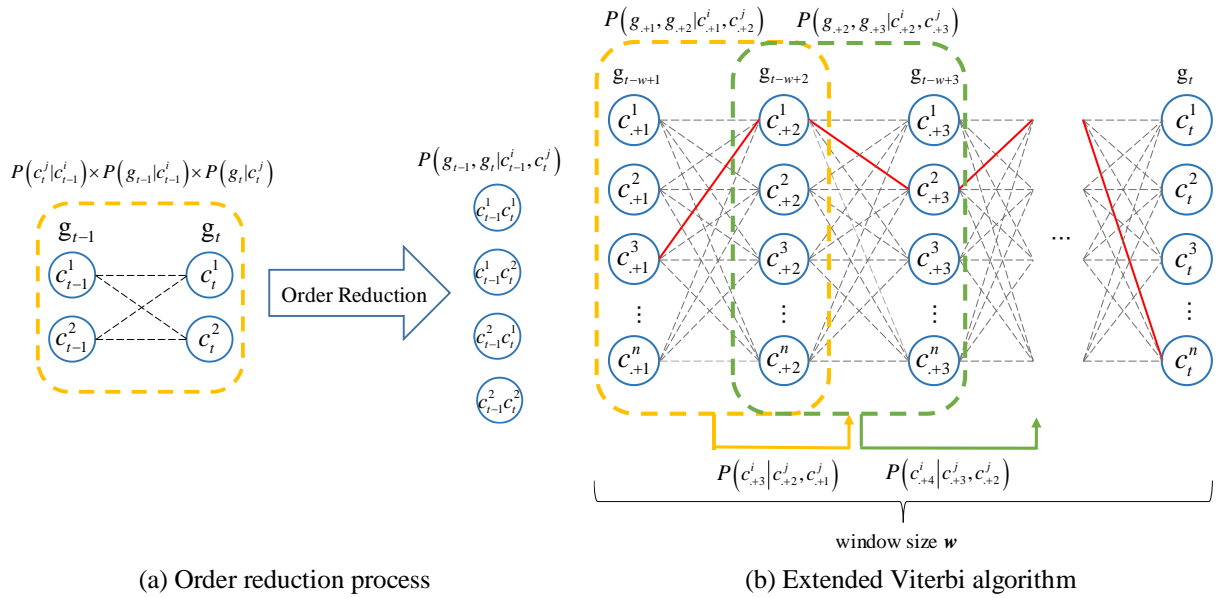


Figure 2: Illustration of the extended Viterbi algorithm.

Extended Viterbi Algorithm

In the previous sections, we introduce the second-order HMM to solve the map matching problem. Although we use sliding window to reduce the computational complexity of matching a single GPS point, the algorithm complexity of traversing the second-order HMM is still $O(n^w)$. Traversal search seriously affects the online performance of the matching algorithm. Thus, some dynamic programming algorithms should be used to reduce the complexity.

The objective function of second-order HMM dynamic programming is defined as

$$\max \prod_{n=t-w+3}^{n=t} \left(P(c_n^i | c_{n-2}^j, c_{n-1}^k) \times P(g_{n-2}, g_{n-1} | c_{n-2}^j, c_{n-1}^k) \right). \quad (11)$$

Viterbi algorithm is an efficient dynamic programming algorithm, which can effectively avoid repeated searches of path and quickly achieve the optimal solution. It is widely used to solve

308 the first-order HMM. To solving the second-order HMM with a complexity of $O(n^2)$, we
309 extend the traditional Viterbi algorithm [38] using an order reduction process as follows:

310 **Step 1: Order Reduction**

311 In the second-order HMM, $P(g_{t-1}, g_t | c_{t-1}^i, c_t^j)$ is regarded as the observation probability,
312 which is equivalent to the observation probability of a single candidate point in the first-order
313 HMM. Equation (8) shows that the observation probability of the second-order HMM is the
314 product of the observation probability of two consecutive candidates in the first-order HMM
315 and the state transition probability. Thus, the order of the second-order HMM can be reduced
316 by using Equation (8) (refer to Figure 2(a)). If the second-order HMM has two layers, each
317 layer has m and n nodes respectively, the second-order HMM can be reduced to one layer with
318 $m \times n$ nodes.

319 **Step 2: Recursive Tracing**

320 After Step1, we can use the traditional Viterbi algorithm for iterative calculation to solve the
321 second-order HMM in the following process (refer to Figure 2(b)):

322 a. Starting from the first layer's nodes, the observation probability of each layer's nodes after
323 reduction and the transition probability between adjacent two layers' nodes are calculated.

324 b. Calculate the maximum total probability of each node from the second layer to the last layer.
325 Save maximum total probability and precursor node of each node.

326 c. Select the node with the highest total probability in the last layer, and go back to its precursor
327 node until the first layer.

328 With the above steps, we can find the optimal matching path $(c_{t-w+1}^i, c_{t-w+2}^j, \dots, c_t^k)$ in the
329 sliding window.

330 **Case study**

331 In this section, we make sensitivity analyses of the parameters involved in the algorithm, and
332 use real data to show the merits of the proposed second-order HMM map matching algorithm.

333 **Data Preparation and Evaluation Metric**

334 We used the road network data of Qinhuai District in Nanjing, China, including 6901 sections
335 and 4647 nodes. Taxi GPS data with 30 seconds sampling interval collected in September 2016
336 were used, including 500 trajectories for 20 taxis. We manually match these trajectories to the
337 road network as the ground truth. In order to verify the effectiveness of the algorithm under
338 extreme conditions and reflect the advantages of the proposed algorithm, we re-sampled the
339 original data and added the random noise of Gaussian distribution. The re-sampling intervals
340 are 60s to 300s. The Gaussian noises with a standard deviation of 10m to 80m (convert to
341 degrees) were added to longitude and latitude.

Evaluation metric is defined as follows. First, we find the common matching sequence X (the sequence that matched correctly) between the matched output route M and the real trajectory T . Based on this sequence, the precision and the recall of the map-matching result (denoted as pcs and rc respectively) can be calculated as

$$pcs = \frac{X}{M}, \quad (12)$$

$$rc = \frac{X}{T}. \quad (13)$$

pcs is defined as the ratio of the length of matched sequence X and the total length of the matched trajectory M . rc is defined as the ratio between the length of matched sequence X and the total length of the real trajectory T . In this study, F_1 -score, which is widely used to evaluate the performance of classification models and prediction models [39], is adopted in this study to evaluate the proposed model:

$$F_1\text{-score} = \frac{2 \cdot pcs \cdot rc}{pcs + rc}. \quad (14)$$

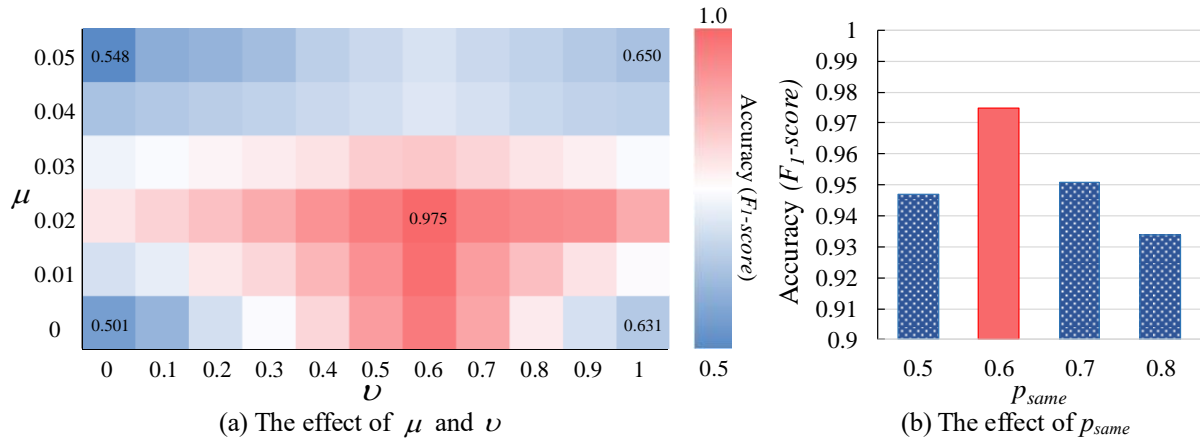


Figure 3: Sensitivity analyses of model parameters.

Table 2. Optimal parameter values.

Parameter	Optimal Value	Range	Impact
μ	0.02	0-0.05	Control the effect of road weight on results
ν	0.6	0-1	Control the effect of vehicle heading weight on results
p_{same}	0.6	0.5-0.8	Control the effect of same/adjacent road priority on results

Results

Effects of different parameters on map matching accuracy are investigated in this study. In the proposed model, there are three parameters to be estimated, i.e. μ, ν and p_{same} . According to previous studies, the approximate range of the three parameters can be obtained. Figure 3 shows the impact of different parameter values on F_1 -score and Table 2 shows the optimal parameter values. It can be seen that, when the road weight μ is around 0.02, the vehicle

heading weight ν is around 0.6, and the same/adjacent road priority p_{same} is around 0.6, their impact on the final performance becomes optimal and stable.

Figure 4(a) shows the effect of window size w on the accuracy of map matching. It can be seen that when $w = 3$, the value of F_1 -score increases significantly. The reason is that when the size of sliding window is larger than 3, the second-order HMM comes into play. Under different standard deviations of noise (SDNs), when the sliding window size increases from 3 to 10, the matching accuracy gradually remains unchanged. However, as the sliding window size increases, the computation time of matching a single GPS point increases rapidly. Thus, the optimal self-adaptive sliding window sizes are 3, 4 and 5.

Figure 4(b) shows the effects of the sample interval and the random SDN on accuracy of map matching. With the increases of sampling interval and SDN, the F_1 -score decreases. It can be seen from Figure 4(b) that when the sampling interval is between 30s to 90s and the SDN ranges from 0 to 30 m, the F_1 -score is kept above 0.9.

With the map matching algorithm proposed in this paper, various factors (i.e. road level, driver's travel preference, vehicle heading and network topology) are considered. Figure 5 shows some map matching results in complex urban road network environment. From Figure 5(a), it can be seen that the first-order HMM map matching algorithm may bring mismatch when it deals with parallel road segments. Under the constraints of topological relations, second-order HMM algorithm gives a greater transition probability to the segment which is adjacent to the previous segment to effectively reduce errors. When GPS points are near road intersections, first-order HMM algorithm may match GPS points to the segment crossed the current road. Second-order HMM and sliding window can help solve this problem. The second-order transition probability can effectively avoid the detour of matching trajectory at the intersection and improve the accuracy of map matching. Figure 5(b) shows an overview of map

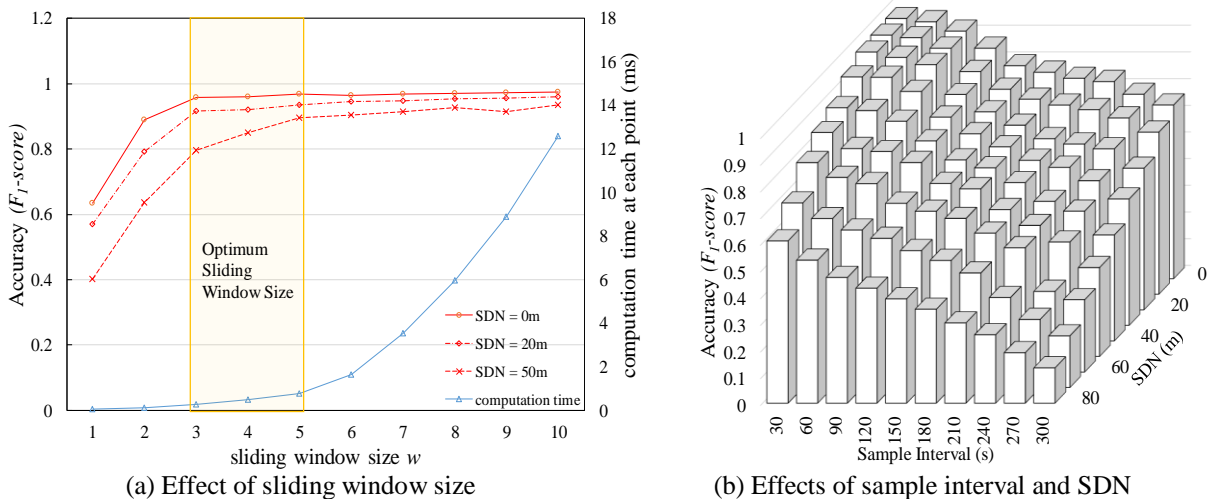
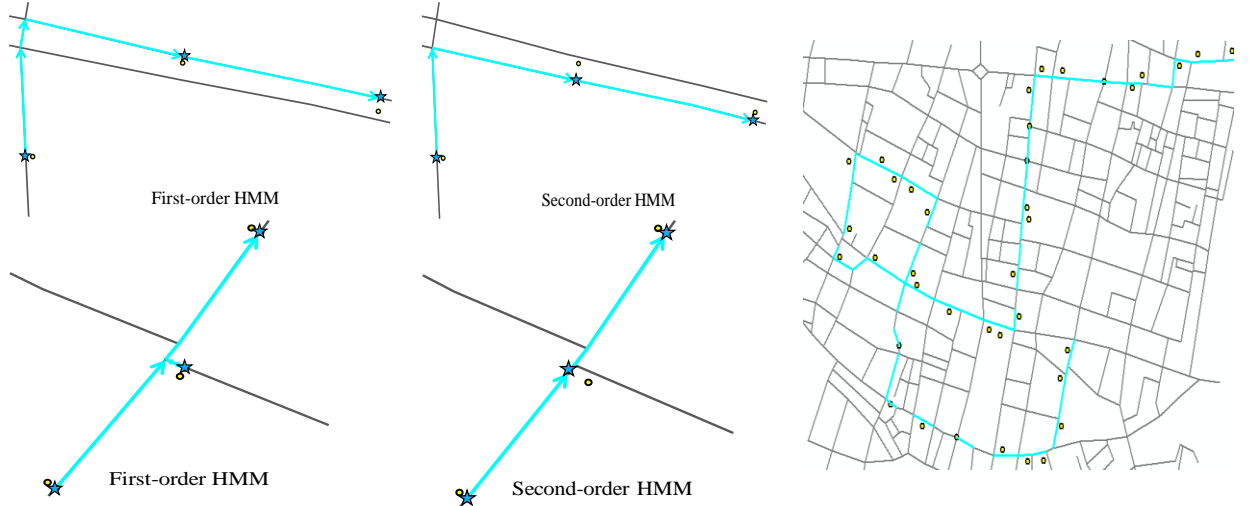
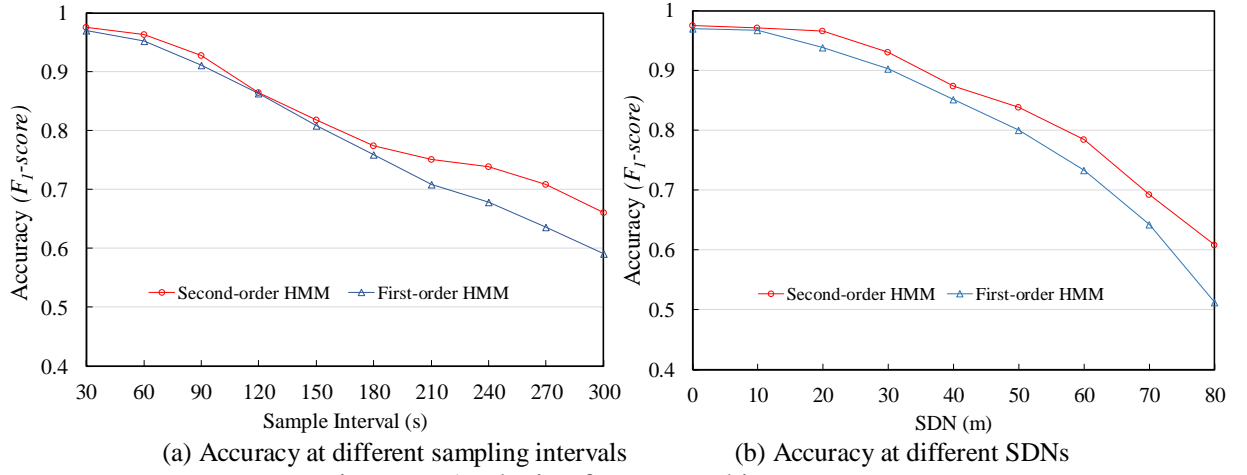


Figure 4: Effects of sliding window size and sampling on map matching accuracy.



(a) Map matching results at parallel segments and road intersections (b) An overview of map matching result
Figure 5: Demonstration of several map matching cases.



(a) Accuracy at different sampling intervals (b) Accuracy at different SDNs
Figure 6: Analysis of map matching accuracy.

matching result in central area of Nanjing, where the road network is dense and complex. The proposed algorithm is found well performed on parallel segments and intersections. This because the second-order HMM model has a wider field of view, and our method considers a variety of factors, which is helpful for map matching in complex conditions.

Figure 6(a) compares the accuracy of the proposed second-order HMM map matching algorithm with the accuracy of our baseline (first-order HMM map matching algorithm) at different sample intervals without adding random noise. It can be seen that the F_1 -score of the proposed algorithm is higher than that of the first-order HMM. With the increase of sampling interval, the advantages of the proposed algorithm become obvious. Taking the 300 seconds sampling interval as an example, the distance between two GPS points is about 2500 meters considering the average speed 30km/h on urban roads. In this situation, the position correlation between two consecutive GPS points is very low. The traditional first-order HMM algorithm only considers the transition probability between two points, so the error tends to be very large. Our proposed algorithm integrates several factors such as road level and driver's travel preference, and second-order transition probability can match GPS trajectory on a larger scale, so it shows higher accuracy (F_1 -score is about 0.67).

Table 3. Comparison of the accuracy (F_1 -score) with some state-of-the-art methods.

Method	Accuracy
HMM-DPP [28]	0.910
SnapNet [30]	0.909
This study	0.975

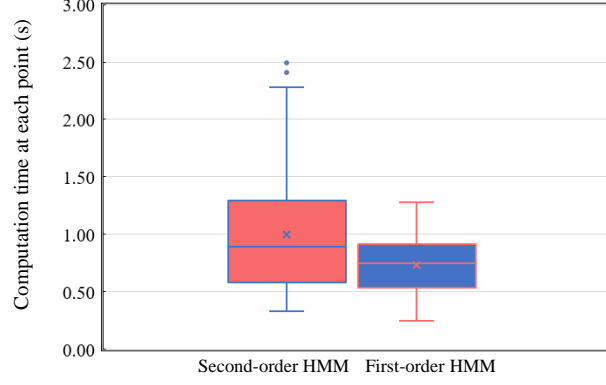


Figure 7: Comparison of map matching efficiency between first-order and second-order HMM.

Figure 6(b) compares the accuracy of the proposed second-order HMM map matching algorithm with our baseline (first-order HMM algorithm) at different SDNs with 30s sample interval. The map matching accuracy of the proposed algorithm is always higher than that of first-order algorithm. The reason is that the conventional first-order HMM algorithm only considers the difference between great-circle distance and route distance when calculating the observation probability of candidate points. When the positioning error of GPS point increases and the road network is dense, matching errors are numerous. In practice, the GPS positioning error is significant in city centre with dense high-rise buildings. As the proposed second-order HMM algorithm excelled conventional algorithms in accuracy (0.6 compared to 0.5 when SDN equals 80m), the proposed algorithm can be adopted to achieve high accuracy of map matching in the whole city.

When comparing with the state-of-the-art methods that most relevant to our proposed method on the condition of raw GPS data, the results in Table 3 show that our second-order HMM method performs well in accuracy. Figure 7 compares the efficiency of the proposed second-order HMM map matching algorithm with conventional first-order HMM algorithm. For the first-order HMM algorithm, the sliding window size is set to 5. It can be seen from Figure 7 that the computation time at each point using the second-order HMM algorithm is slightly longer than that using the first-order HMM algorithm, and the average computation time is less than 1s. In the process of self-adaption of sliding window size, a small number of outliers appear. For example, using the second-order HMM algorithm, there are a few points whose computation time is longer than 2 seconds. However, in this example, the overall matching efficiency is close to the first-order HMM map matching, which can meet the requirements of online map matching. Moreover, compared to the first-order HMM, the second-order HMM can better consider the spatial-temporal relationship among different states and larger perception fields, which can get remarkable accuracy under complex conditions.

Conclusions

Accurate and efficient matching of GPS data onto road network is the basis and prerequisite for conducting traffic flow analysis and providing location-based service. An online map matching algorithm based on second-order HMM is presented in this paper. Various factors (i.e. road level, driver's travel preference, vehicle heading and network topology) are explicitly considered in the algorithm, which effectively improve the accuracy of map matching in complex urban road network environment. An extended Viterbi algorithm is adopted to solve

the map matching problem efficiently. A self-adaptive sliding window mechanism is proposed to adjust window size on a real time basis and ensures high accuracy.

We tested the proposed algorithm using real road network and massive taxi GPS data collected in Nanjing, China. The proposed map matching approach was found outperform state-of-the-art algorithms built on the first-order HMM in various testing environment. Sliding window with self-adaptive size is shown to be an effective method for online incremental map matching. Some typical types of mismatching can be avoided in complex urban road network environment such as parallel road segments and various road intersections. The map matching accuracy of the proposed algorithm is demonstrated higher than that of conventional first-order HMM algorithm. The efficiency of the proposed algorithm is close to the first-order HMM map matching algorithm, which can meet the requirements of online map matching. Therefore, the proposed algorithm is applicable in real-time navigation, trajectory monitoring, traffic flow analysis and other related fields.

To solve the map matching problem, there are some other solutions such as considering driving direction and turning behaviour. The consideration of users with heterogeneous activity/travel behaviour is suggested as another interesting extension of the proposed method, potentially improving the accuracy of map matching [31, 40]. In the case study, the proposed algorithm is tested using a single processor. How to incorporate the parallel computing technologies into the proposed algorithm with a large number of trajectories needs further investigation [41]. Besides, the comparison of the advantages and disadvantages of the second-order HMM based method and other advanced map matching algorithms can also be the focus of future research.

Acknowledgments

The work described in this paper was jointly supported by the National Key Research and Development Program of China (Grant No. 2018YFB1600900), the National Natural Science Foundation of China (71601045) and “Zhishan” Scholars Programs of Southeast University.

References

- [1] Wong, W. and S. C. Wong. Network Topological Effects on the Macroscopic Bureau of Public Roads Function. *Transportmetrica*, 2015. 12:272–296.
- [2] Bao, J., P. Liu, X. Qin, and H. Zhou. Understanding the Effects of Trip Patterns on Spatially Aggregated Crashes with Large-scale Taxi GPS Data. *Accident Analysis & Prevention*, 2018. 120: 281–294.
- [3] Huang, W., W. Jia, J. Guo, B. M. Williams, G. Shi, Y. Wei, and J. Cao. Real-time Prediction of Seasonal Heteroscedasticity in Vehicular Traffic Flow Series. *IEEE Transactions on Intelligent Transportation Systems*, 2018. 19: 3170–3180.
- [4] He, Z., Y. Lv, L. Lu, and W. Guan. Constructing spatiotemporal speed contour diagrams: using rectangular or non-rectangular parallelogram cells? *Transportmetrica B: Transport Dynamics*, 2019. 7: 44–60.
- [5] Chen, B. Y., H. Yuan, Q. Q. Li, S. L. Shaw, W. H. K. Lam, and X. Chen. Spatiotemporal data model for network time geographic analysis in the era of big data. *International Journal of Geographical Information Science*, 2016, 30(6): 1041–1071.
- [6] Bierlaire, M., J. Chen, and J. Newman. A Probabilistic Map Matching Method for Smartphone GPS Data. *Transportation Research Part C*, 2013. 26: 78–98.

- [7] Luo, Q., J. Auld, and V. Sokolov. Addressing Some Issues of Map-matching for Large-scale, High-frequency GPS Data Sets. Presented at 95th Annual Meeting of the Transportation Research Board, Washington, D.C., 2016.
- [8] White, C. E., D. Bernstein, and A. L. Kornhauser. Some map matching algorithms for personal navigation assistants. *Transportation Research Part C*, 2000. 8: 91–108.
- [9] Taylor, G., G. Blewitt, D. Steup, S. Corbett, and A. Car. Road Reduction Filtering for GPS-GIS Navigation. *Transactions in GIS*, 2001. 5: 193–207.
- [10] Schweizer J., S. Bernardi, and F. Rupi. Map-matching algorithm applied to bicycle global positioning system traces in Bologna. *IET Intelligent Transport Systems*, 2016. 10(4):244–250.
- [11] Velaga, N. R., M. A. Quddus, and A. L. Bristow. Improving the Performance of a Topological Map-Matching Algorithm through Error Detection and Correction. Presented at 91th Annual Meeting of the Transportation Research Board, Washington, D.C., 2012.
- [12] Chen, B. Y., H. Yuan, Q. Li, W. H. K. Lam, S. L. Shaw, and K. Yan. Map-Matching Algorithm for Large-scale Low-frequency Floating Car Data. *International Journal of Geographical Information Science*, 2014. 28(1): 22–38.
- [13] Tradišauskas, N., J. Juhl, H. Lahrmann, and C.S. Jensen. Map matching for intelligent speed adaptation. *IET Intelligent Transport Systems*, 2009. 3(1): 57–66.
- [14] He, Z., X. She, L. Zhuang, and P. Nie. On-line map-matching framework for floating car data with low sampling rate in urban road networks. *IET Intelligent Transport Systems*, 2013. 7(4): 404–414.
- [15] Zheng, L., X. Liu, B. Yi. Dynamic Weighted Real-time Map Matching Algorithm Considering Spatio-Temporal Property. *Journal of Computer Applications*, 2017. 37(8): 2381–2386.
- [16] Liu, X. L., K. Liu, M. X. Li, and F. Lu. A ST-CRF Map-Matching Method for Low-Frequency Floating Car Data. *IEEE Transactions on Intelligent Transportation Systems*, 2017. 18(5): 1241–1254.
- [17] Gong, Y. J., E. Chen, X. Zhang, L. M. Ni, and J. Zhang. AntMapper: An Ant Colony-based Map Matching Approach for Trajectory-based Applications. *IEEE Transactions on Intelligent Transportation Systems*, 2018. 19(2): 390–401.
- [18] Zhang, K., S. Liu, Y. Dong, D. Wang, Y. Zhang, and L. Miao. Vehicle positioning system with multihypothesis map matching and robust feedback. *IET Intelligent Transport Systems*, 2017. 11(10): 649–658.
- [19] Takenga, C., T. Peng, and K. Kyamakya. Post-processing of Fingerprint Localization using Kalman Filter and Map-Matching Techniques. Presented at The 9th International Conference on Advanced Communication Technology, 2007.
- [20] Taguchi, S., S. Koide, and T. Yoshimura. Online Map Matching with Route Prediction. *IEEE Transactions on Intelligent Transportation Systems*, 2018. 1–10.
- [21] Ren, M. and H. A. Karimi. A Fuzzy Logic Map Matching for Wheelchair Navigation. *GPS Solutions*, 2012. 16(3): 273–282.
- [22] Chen, Z. and R. C. Qiu. Prediction of Channel State for Cognitive Radio using Higher-Order Hidden Markov Model. *Proceedings of the IEEE SoutheastCon 2010 (SoutheastCon)*, 2010. DOI: 10.1109/SECON.2010.5453870.
- [23] Newson, P. and J. Krumm. Hidden Markov Map Matching through Noise and Sparseness. *Proceedings of ACM SIGSPATIAL*, 2009. 336–343.

516 [24] Yuan, J., Y. Zheng, C. Zhang, X. Xie, and G. Sun. An Interactive-Voting Based Map Matching
517 Algorithm. *Proceedings of Eleventh International Conference on Mobile Data Management*, 2010. DOI:
518 10.1109/MDM.2010.14.

519 [25] Song, R., W. Lu, W. Sun, W. Huang, and C. Chen. Quick Map Matching using Multi-core CPUs.
520 *Proceedings of the ACM-GIS*, 2012. 605–608. DOI:10.1145/2424321.2424428.

521 [26] Atia, M. M., A. R. Hilal, C. Stellings, E. Hartwell, J. Toonstra, W. B. Miners, and O.A. Basir. A Low-
522 cost Lane-Determination System using GNSS/IMU Fusion and HMM-based Multistage Map Matching.
523 *IEEE Transactions on Intelligent Transportation Systems*, 2017. 1–11.

524 [27] Bonnetain, L., A. Furno, J. Krug, and N-E. E. Faouz. Can We Map-Match Individual Cellular Network
525 Signaling Trajectories in Urban Environments? Data-Driven Study. *Transportation Research Record:
526 Journal of the Transportation Research Board*, 2019. <https://doi.org/10.1177/0361198119847472>.

527 [28] Song, C., X. Yan, N. Stephen, A.A. Khan. Hidden Markov model and driver path preference for floating
528 car trajectory map matching. *IET Intelligent Transport Systems*, 2018. 12(10): 1433–1441

529 [29] Lou, Y., C. Zhang, Y. Zheng, W. Wang, and Y. Huang. Map-Matching for Low-Sampling-Rate GPS
530 Trajectories. *Proceedings of the ACM-GIS*, 2009. 352–361.

531 [30] Mohamed, R., H. Aly, and M. Youssef. Accurate Real-time Map Matching for Challenging
532 Environments. *IEEE Transactions on Intelligent Transportation Systems*, 2017. 18(4): 847–857.

533 [31] Jagadeesh, G. R., and T. Srikanthan. Online Map-Matching of Noisy and Sparse Location Data with
534 Hidden Markov and Route Choice Models. *IEEE Transactions on Intelligent Transportation Systems*,
535 2017. 1–12.

536 [32] Salnikov, V., M. T. Schaub, and R. Lambiotte. Using Higher-Order Markov Models to Reveal Flow-
537 based Communities in Networks. *Scientific Reports*, 2016: 6.

538 [33] Leibowitz, N., B. Baum, G. Enden, and A. Karniel. The Exponential Learning Equation as a Function of
539 Successful Trials Results in Sigmoid Performance, *Journal of Mathematical Psychology*, 2010. 54(3):
540 338–340.

541 [34] Eddy, S. R. What is a Hidden Markov Model? *Nature Biotechnology*, 2004. 22: 1315–1316.

542 [35] Ye, F. and Y. Wang. Research Advancement of High-Order Hidden Markov Model. *Advances in
543 Mathematics*, 2014. 2(2): 845–848.

544 [36] Lee, T., F. Zheng, W. Wu, and D. Chen. The Hidden Markov Model of Co-articulation and its
545 Application to the Continuous Speech Recognition. *Journal of Electronics (China)*, 2000. 17(3): 242–
546 247.

547 [37] Gales, M. J. F., S. Watanabe, and E. Fosler-Lussier. Structured discriminative Models for Speech
548 Recognition: An Overview. *IEEE Signal Processing Magazine*, 2012. 29(6): 70–81.

549 [38] He, Y. Extended Viterbi algorithm for second order hidden Markov process. *Pattern Recognition*, 1988.
550 9th International Conference on. IEEE.

551 [39] Voorhees, E. M. Variation in Relevance Judgments and the Measurement of Retrieval Effectiveness.
552 *Information Processing & Management*, 2000. 36(5): 697–716.

553 [40] Fu, X. and W. H. K. Lam. Modelling Joint Activity-Travel Pattern Scheduling Problem in Multi-modal
554 Transit Networks. *Transportation*, 2018. 45: 23–49.

555 [41] Li, Q., T. Zhang, and Y. Yu. Using Cloud Computing to Process Intensive Floating Car Data for Urban
556 Traffic Surveillance. *International Journal of Geographical Information Science*, 2011. 25: 1303–1322.