

# Trajectory Data Mining: An Overview

YU ZHENG, Microsoft Research

The advances in location-acquisition and mobile computing techniques have generated massive spatial trajectory data, which represent the mobility of a diversity of moving objects, such as people, vehicles, and animals. Many techniques have been proposed for processing, managing, and mining trajectory data in the past decade, fostering a broad range of applications. In this article, we conduct a systematic survey on the major research into *trajectory data mining*, providing a panorama of the field as well as the scope of its research topics. Following a road map from the derivation of trajectory data, to trajectory data preprocessing, to trajectory data management, and to a variety of mining tasks (such as trajectory pattern mining, outlier detection, and trajectory classification), the survey explores the connections, correlations, and differences among these existing techniques. This survey also introduces the methods that transform trajectories into other data formats, such as graphs, matrices, and tensors, to which more data mining and machine learning techniques can be applied. Finally, some public trajectory datasets are presented. This survey can help shape the field of *trajectory data mining*, providing a quick understanding of this field to the community.

Categories and Subject Descriptors: H.2.8 [Database Management]: Database Applications—Data mining, spatial databases and GIS; I.2.6 [Artificial Intelligence]: Learning—Knowledge acquisition

General Terms: Algorithms, Measurement, Experimentation

Additional Key Words and Phrases: Spatiotemporal data mining, trajectory data mining, trajectory compression, trajectory indexing and retrieval, trajectory pattern mining, trajectory outlier detection, trajectory uncertainty, trajectory classification, urban computing

## ACM Reference Format:

Yu Zheng. 2015. Trajectory data mining: An overview. ACM Trans. Intell. Syst. Technol. 6, 3, Article 29 (May 2015), 41 pages.

DOI: <http://dx.doi.org/10.1145/2743025>

## 1. INTRODUCTION

A *spatial trajectory* is a trace generated by a moving object in geographical spaces, usually represented by a series of chronologically ordered points, for example,  $p_1 \rightarrow p_2 \rightarrow \dots \rightarrow p_n$ , where each point consists of a geospatial coordinate set and a time stamp such as  $p = (x, y, t)$ .

The advance in location-acquisition technologies has generated a myriad of spatial trajectories representing the mobility of various moving objects, such as people, vehicles, and animals. Such trajectories offer us unprecedented information to understand moving objects and locations, fostering a broad range of applications in location-based social networks [Zheng 2011], intelligent transportation systems, and urban computing [Zheng et al. 2014b]. The prevalence of these applications in turn calls for systematic research on new computing technologies for discovering knowledge from trajectory data. Under the circumstances, *trajectory data mining* has become an increasingly

---

Authors' addresses: Y. Zheng, Microsoft Research, Building 2, No. 5 Danling Street, Haidian District, Beijing 100080, China; email: yuzheng@microsoft.com.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212) 869-0481, or [permissions@acm.org](mailto:permissions@acm.org).

© 2015 ACM 2157-6904/2015/05-ART29 \$15.00

DOI: <http://dx.doi.org/10.1145/2743025>

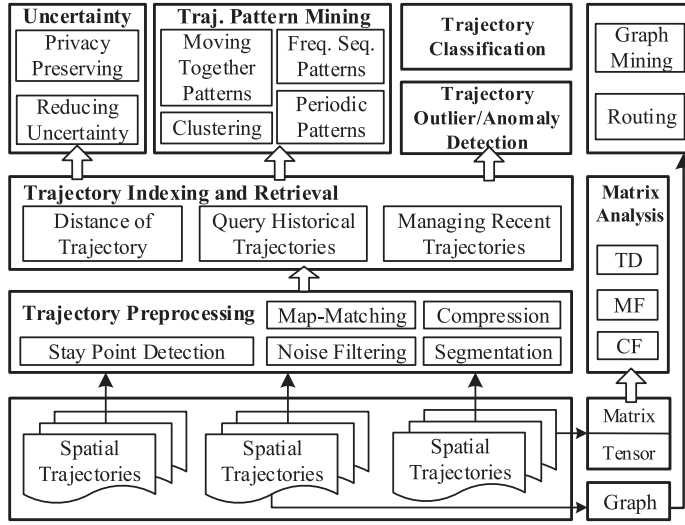


Fig. 1. Paradigm of trajectory data mining.

important research theme, attracting the attention from numerous areas, including computer science, sociology, and geography.

Intensive and extensive individual research has been done in the field of *trajectory data mining*. However, we are lack of a systematic review that can well shape the field and position existing research. Facing a huge volume of publications, the community is still not very clear about the connections, correlations and difference among these existing techniques. To this end, we conduct a comprehensive survey that thoroughly explores the field of trajectory data mining, according to the paradigm shown in Figure 1:

First, in Section 2, we classify the sources generating trajectory data into four groups, listing a few key applications that trajectory data can enable in each group.

Secondly, before using trajectory data, we need to deal with a number of issues, such as noise filtering, segmentation, and map matching. This stage is called *trajectory preprocessing*, which is a fundamental step of many trajectory data mining tasks. The goal of *noise filtering* is to remove from a trajectory some noise points that may be caused by the poor signal of location positioning systems (e.g., when traveling in a city canyon). *Trajectory compression* is to compress the size of a trajectory (for the purpose of reducing overhead in communication, processing, and data storage) while maintaining the utility of the trajectory. A *stay point detection* algorithm identifies the location where a moving object has stayed for a while within a certain distance threshold. A stay point could stand for a restaurant or a shopping mall that a user has been to, carrying more semantic meanings than other points in a trajectory. *Trajectory segmentation* divides a trajectory into fragments by time interval, spatial shape, or semantic meanings, for a further process like clustering and classification. *Map matching* aims to project each point of a trajectory onto a corresponding road segment where the point was truly generated. We detail trajectory pre-processing in Section 3.

Thirdly, many online applications require instantly mining trajectory data (e.g., detecting traffic anomalies), calling for effective data management algorithms that can quickly retrieve particular trajectories satisfying certain criteria (such as spatiotemporal constraints) from a big trajectory corpus. There are usually two major types of queries: the nearest neighbors and range queries. The former is also associated with a

distance metric, for example, the distance between two trajectories. Additionally, there are two types (historical and recent) of trajectories, which need different managing methods. We will introduce *trajectory indexing and retrieval* in Section 4.

Fourthly, based on the first two steps, we can then conduct mining tasks, like trajectory pattern mining, trajectory uncertainty, outlier detection, and classification.

- Trajectory Uncertainty*: Objects move continuously while their locations can only be updated at discrete times, leaving the location of a moving object between two updates uncertain. To enhance the utility of trajectories, a series of research tried to model and reduce the uncertainty of trajectories. On the contrary, a branch of research aims to protect a user's privacy when the user discloses her trajectories. We review uncertainty of trajectory in Section 5.
- Trajectory Pattern Mining*: The huge volume of spatial trajectories enables opportunities for analyzing the mobility patterns of moving objects, which can be represented by an individual trajectory containing a certain pattern or a group of trajectories sharing similar patterns. In Section 6, we survey the literature that is concerned with four categories of patterns: moving together patterns, trajectory clustering, periodic patterns, and frequent sequential patterns.
- Trajectory Classification*: Using supervised learning approaches, we can classify trajectories or segments of a trajectory into some categories, which can be activities (like hiking and dining) or different transportation modes, such as walking and driving. We show examples of trajectory classification in Section 7.
- Trajectory Outlier Detection*: Different from trajectory patterns that frequently occur in trajectory data, trajectory outliers (a.k.a. anomalies) can be items (a trajectory or a segment of trajectory) that are significantly different from other items in terms of some similarity metric. It can also be events or observations (represented by a collection of trajectories) that do not conform to an expected pattern (e.g., traffic congestion caused by a car accident). Section 8 introduces outlier/anomaly detection from trajectory data.

Finally, besides studying trajectories in its original form, we can transform trajectories into other formats, such as graph, matrix, and tensor (see the right part of Figure 1). The new representations of trajectories expand and diversify the approaches for trajectory data mining, leveraging existing mining techniques (e.g., graph mining, Collaborative Filtering (CF), Matrix Factorization (MF), and Tensor Decomposition (TD)). In Section 9, we present representative examples of the transformation.

The contribution of this article is fourfold. First, the article presents a framework for trajectory data mining, defining the scope and road map for this field. The framework provides a panorama with which people can quickly understand and step into this field. Second, individual research works are well positioned, categorized, and connected in each layer of this framework. Professionals can easily locate the methods they need to solve a problem, or find the unsolved problems. Third, this article proposes a vision to transfer trajectories into other formats, to which a diversity of existing mining techniques can be applied. This expands the original scope of trajectory data mining, advancing the methodologies and applications of this field. Fourth, we collect a list of sources from which people can obtain various public trajectory datasets for research. We also introduce the conferences and journals that are concerned with the research on trajectory data.

## 2. TRAJECTORY DATA

In this section, we classify the derivation of trajectories into four major categories, briefly introducing a few application scenarios in each category. Trajectory data representing human mobility can help build a better social network [Bao et al. 2015; Zheng

2011; Zheng et al. 2012b] and travel recommendation [Zheng and Xie 2011b; Zheng et al. 2011c; Zheng et al. 2009b].

- (1) **Mobility of people:** People have been recording their real-world movements in the form of spatial trajectories, passively and actively, for a long time.
  - Active Recording:* Travelers log their travel routes with GPS trajectories for the purpose of memorizing a journey and sharing experiences with friends. Bicyclers and joggers record their trails for sports analysis. In Flickr, a series of geotagged photos can formulate a spatial trajectory as each photo has a location tag and a time stamp corresponding to where and when the photo was taken. Likewise, the “check-ins” of a user in a location-based social network can be regarded as a trajectory, when sorted chronologically.
  - Passive Recording:* A user carrying a mobile phone unintentionally generates many spatial trajectories represented by a sequence of cell tower IDs with corresponding transition times. Additionally, transaction records of a credit card also indicate the spatial trajectory of the cardholder, as each transaction contains a time stamp and a merchant ID denoting the location where the transaction occurred.
- (2) **Mobility of transportation vehicles:** A large number of GPS-equipped vehicles (such as taxis, buses, vessels, and aircrafts) have appeared in our daily life. For instance, many taxis in major cities have been equipped with a GPS sensor, which enables them to report a time-stamped location with a certain frequency. Such reports formulate a large amount of spatial trajectories that can be used for resource allocation [Yuan et al. 2011b, 2013b], traffic analysis [Wang et al. 2014; Yuan et al. 2013a], and improving transportation networks [Zheng et al. 2011a].
- (3) **Mobility of animals:** Biologists have been collecting the moving trajectories of animals like tigers and birds, for the purpose of studying animals’ migratory traces, behavior, and living situations [Lee et al. 2007; Li et al. 2010c].
- (4) **Mobility of natural phenomena:** Meteorologists, environmentalists, climatologists, and oceanographers are busy collecting the trajectories of some natural phenomena, such as hurricanes, tornados, and ocean currents. These trajectories capture the change of the environment and climate, helping scientists deal with natural disasters and protect the natural environment we live in.

### 3. TRAJECTORY DATA PREPROCESSING

This section introduces a fourfold of basic techniques that we need to process a trajectory before starting a mining task, consisting of noise filtering, stay point detection, trajectory compression, and trajectory segmentation.

#### 3.1. Noise Filtering

Spatial trajectories are never perfectly accurate, due to sensor noise and other factors, such as receiving poor positioning signals in urban canyons. Sometimes, the error is acceptable (e.g., a few GPS points of a vehicle fall out of the road the vehicle was actually driven), which can be fixed by map-matching algorithms (introduced in Section 3.5). In other situations, as shown in Figure 2, the error of a noise point like  $p_5$  is too big (e.g., several hundred meters away from its true location) to derive useful information, such as travel speed. So, we need to filter such noise points from trajectories before starting a mining task. Though this problem has not been completely solved, existing methods fall into three major categories.

*Mean (or Median) Filter.* For a measured point  $z_i$ , the estimate of the (unknown) true value is the mean (or median) of  $z_i$  and its  $n-1$  predecessors in time. The mean (median) filter can be thought of as a sliding window covering  $n$  temporally adjacent values of

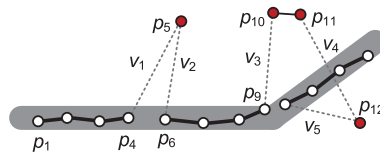


Fig. 2. Noise points in a trajectory.

$z_i$ . In the example shown in Figure 2,  $p_5 \cdot z = \sum_{i=1}^5 p_i \cdot z / 5$ , if we use a mean filter with a sliding window size of 5. The median filter is more robust than the mean filter when handling extreme errors. The mean (median) filters are practical for handling individual noise points like  $p_5$  in a trajectory with a dense representation. However, when dealing with multiple consecutive noise points, for example,  $p_{10}$ ,  $p_{11}$ , and  $p_{12}$ , a larger size of sliding window is needed. This results in a bigger error between the calculated mean (or median) value and a point's true position. When the sampling rate of trajectory is very low (i.e., the distance between two consecutive points could be longer than several hundred meters), the mean and median filters are not good choices anymore.

*Kalman and Particle Filters.* The trajectory estimated from the Kalman filter is a trade-off between the measurements and a motion model. Besides giving estimates that obey the laws of physics, the Kalman filter gives principled estimates of higher order motion states like speed. While the Kalman filter gains efficiency by assuming linear models plus Gaussian noise, the particle filter relaxes these assumptions for a more general, but less efficient, algorithm. A tutorial-like introduction to using the Kalman and particle filters to fix noisy trajectory points can be found in Lee and Krumm [2011].

The initialization step of the particle filtering is to generate  $P$  particles  $\mathbf{x}_i^{(j)}$ ,  $j = 1, 2, \dots, P$  from the initial distribution. For example, these particles would have zero velocity and be clustered around the initial location measurement with a Gaussian distribution. The second step is “importance sampling,” which uses the dynamic model  $P(\mathbf{x}_i | \mathbf{x}_{i-1})$  to probabilistically simulate how the particles change over one timestep. The third step computes “importance weights” for all the particles using the measurement model  $\omega_i^{(j)} = P(z_i | \hat{\mathbf{x}}_i^{(j)})$ . Larger importance weights correspond to particles that are better supported by the measurement. The important weights are then normalized so they sum to one. The last step in the loop is the “selection step” when a new set of  $P$  particles  $\mathbf{x}_i^{(j)}$  is selected from the  $\hat{\mathbf{x}}_i^{(j)}$  proportional to the normalized importance weights  $\omega_i^{(j)}$ . Finally, we can compute a weight sum by  $\hat{\mathbf{x}}_i = \sum_{j=1}^P \omega_i^{(j)} \hat{\mathbf{x}}_i^{(j)}$ .

The Kalman and particle filters, model both the measurement noise and the dynamics of the trajectory. However, they depend on the measurement of an initial location. If the first point in a trajectory is noisy, the effectiveness of the two filters drops significantly.

*Heuristics-Based Outlier Detection.* While the previously mentioned filters replace a noise measurement in a trajectory with an estimated value, the third category of methods removes noise points directly from a trajectory by using outlier detection algorithms. The noise filtering method, which has been used in T-Drive [Yuan et al. 2010a, 2011a, 2013a] and GeoLife [Zheng et al. 2009a; Zheng et al. 2010] projects, first calculates the travel speed of each point in a trajectory based on the time interval and distance between a point and its successor (we call this a segment). The segments, such as  $p_4 \rightarrow p_5$ ,  $p_5 \rightarrow p_6$ , and  $p_9 \rightarrow p_{10}$  (illustrated by the dotted lines in Figure 2), with a speed larger than a threshold (e.g., 300km/h), are cut off. Given that the number of noise points is much smaller than common points, the separated points like  $p_5$  and  $p_{10}$  can be regarded as outliers. Some distance-based outlier detection can easily find



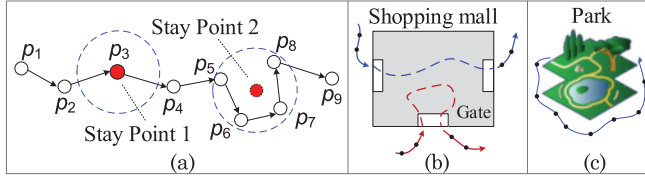


Fig. 3. Stay points in a trajectory.

the number of  $p_5$ 's neighbors within a distance  $d$  is smaller than  $p$  proportion of the points in the entire trajectory. Likewise,  $p_{10}$ ,  $p_{11}$ , and  $p_{12}$  can be filtered. While such algorithms can handle the initial error in a trajectory and data sparsity problems, setting the threshold  $d$  and  $p$  is still based on heuristics.

### 3.2. Stay Point Detection

Spatial points are not equally important in a trajectory. Some points denote locations where people have stayed for a while, such as shopping malls and tourist attractions, or gas stations where a vehicle was refueled. We call this kind of points “*Stay Points*.” As shown in Figure 3(a), there are two types of stay points occurring in a trajectory. One is a single point location, for example, *Stay Point 1*, where a user remains stationary for a while. This situation is very rare, because a user’s positioning device usually generates different readings even in the same location. The second type, like *Stay Points 2* shown in Figure 3(a), is more generally observed in trajectories, representing the places where people move around (e.g., as depicted in Figures 3(b) and 3(c)) or remain stationary but with positioning readings shifting around.

With such stay points, we can turn a trajectory from a series of time-stamped spatial points  $\mathbf{P}$  into a sequence of meaningful places  $\mathbf{S}$ ,

$$\mathbf{P} = p_1 \rightarrow p_2 \rightarrow \dots \rightarrow p_n, \Rightarrow \mathbf{S} = s_1 \xrightarrow{\Delta t_1} s_2 \xrightarrow{\Delta t_2} \dots \xrightarrow{\Delta t_{n-1}} s_n,$$

therefore facilitating a diversity of applications, such as travel recommendations [Zheng and Xie 2011b; Zheng et al. 2011c], destination prediction [Ye et al. 2009], taxi recommendation [Yuan et al. 2011b, 2013b], and gas consumption estimation [Zhang et al. 2013, 2015]. On the other hand, in some applications, for example, estimating the travel time of a path [Wang et al. 2014] and driving direction suggestion [Yuan et al. 2013a], such stay points should be removed from a trajectory during the preprocessing.

Li et al. [2008] first proposed the stay point detection algorithm. This algorithm first checks if the distance between an anchor point (e.g.,  $p_5$ ) and its successors is in a trajectory larger than a given threshold (e.g., 100 m). It then measures the time span between the anchor point and the last successor (i.e.,  $p_8$ ) that is within the distance threshold. If the time span is larger than a given threshold, a stay point (characterized by  $p_5$ ,  $p_6$ ,  $p_7$ , and  $p_8$ ) is detected; the algorithm starts detection the next stay point from  $p_9$ . Yuan et al. [2011b, 2013b] improved this stay point detection algorithm based on the idea of density clustering. After finding  $p_5$  to  $p_8$  is a candidate stay point (using  $p_5$  as an anchor point), their algorithm further checks the successor points from  $p_6$ . For instance, if the distance from  $p_9$  to  $p_6$  is smaller than the threshold,  $p_9$  will be added into the stay point.

### 3.3. Trajectory Compression

Basically, we can record a time-stamped geographical coordinate every second for a moving object. But, this costs a lot of battery power and the overhead for communication, computing, and data storage. In addition, many applications do not really need such a precision of location. To address this issue, two categories of trajectory

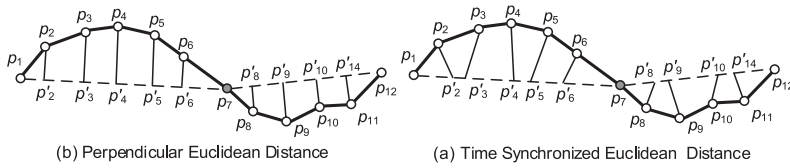


Fig. 4. Distance metric measuring the compression error.

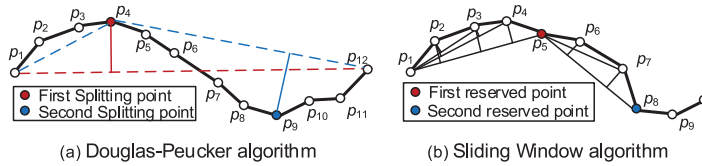


Fig. 5. Illustration of Douglas-Peucker algorithm.

compression strategies (based on the shape of a trajectory) have been proposed, aiming to reduce the size of a trajectory while not compromising much precision in its new data representation [Lee and Krumm 2011]. One is the offline compression (a.k.a. batch mode), which reduces the size of trajectory after the trajectory has been fully generated. The other is online compression, compressing a trajectory instantly as an object travels.

*Distance Metric.* Besides the two strategies, there are two distance metrics to measure the error of a compression: *perpendicular Euclidean distance* and *time synchronized Euclidean distance*. As illustrated in Figure 4, supposing we compress a trajectory with 12 points into a representation of three points (i.e.,  $p_1$ ,  $p_7$ , and  $p_{12}$ ), the two distance metrics are the summation of the lengths of the segments connecting  $p_i$  and  $p'_i$ , in Figures 4(a) and 4(b), respectively. The latter distance assumes a constant speed traveling between  $p_1$  and  $p_7$ , calculating the projection of each original point on  $\overline{p_1 p_7}$  by time intervals.

*Offline Compression.* Given a trajectory that consists of a full series of time-stamped points, a batched compression algorithm aims to generate an approximated trajectory by discarding some points with a negligible error from the original trajectory. This is similar to the line simplification problem, which has been studied in the computer graphics and cartography research communities [McMaster 1986].

A well-known algorithm, called Douglas-Peucker [Douglas and Peucker 1973], is used to approximate the original trajectory. As demonstrated in Figure 5(a), the idea of Douglas-Peucker is to replace the original trajectory by an approximate line segment, for example,  $\overline{p_1 p_{12}}$ . If the replacement does not meet the specified error requirement (perpendicular Euclidean distance is used in this example), it recursively partitions the original problem into two subproblems by selecting the point contributing the biggest error as the splitting point (e.g.,  $p_4$ ). This process continues until the error between the approximation and the original trajectory is below a specified error. The complexity of the original Douglas-Peucker algorithm is  $O(N^2)$ , where  $N$  is the number of points in a trajectory. Its improvement achieves  $O(N \log N)$  [Hershberger and Snoeyink 1992]. To ensure that the approximated trajectory is optimal, Bellman's algorithm [Bellman 1961] employs a dynamic programming technique with a complexity of  $O(N^3)$ .

*Online Data Reduction.* As many applications require one to transmit trajectory data in a timely fashion, a series of online trajectory compression techniques have been proposed to determine whether a newly acquired spatial point should be retained

in a trajectory. There are two major categories of online compression methods. One is the window-based algorithms, such as the *Sliding Window* algorithm [Keogh et al. 2001] and *Open Window* algorithm [Maratnia and de By 2004]. The other is based on the speed and direction of a moving object.

The idea of the *Sliding Window* algorithm is to fit the spatial points in a growing sliding window with a valid line segment and continue to grow the sliding window until the approximation error exceeds some error bound. As illustrated in Figure 5(b),  $p_5$  will be first reserved as the error for  $p_3$  exceeds the threshold. Then, the algorithm starts from  $p_5$  and reserve  $p_8$ . Other points are negligible. Different from the *Sliding Window* algorithm, the *Open Window* algorithm [Maratnia and de By 2004] applies the heuristic of the Douglas-Peucker algorithm to choose the point with the maximum error in the window (e.g.,  $p_3$  in Figure 5(b)) to approximate the trajectory segment. This point is then used as a new anchor point to approximate its successors.

Another category of algorithms consider speed and directions as key factors when doing online trajectory compression. For instance, Potamias et al. [2006] use a safe area, derived from the last two locations and a given threshold, to determine whether a newly acquired point contains important information. If the new data point is located within the safe area, then this location point is considered as redundant and thus can be discarded; otherwise, it is included in the approximated trajectory.

*Compression with Semantic Meaning.* A series of research [Richter et al. 2012; Chen et al. 2009] aims to keep the semantic meanings of a trajectory, when compressing the trajectory. For instance, in a location-based social network [Zheng 2011], some special points where a user stayed, took photos, or changed direction greatly, would be more significant than other points in presenting semantic meanings of a trajectory. Chen et al. [2009] proposed a Trajectory Simplification (TS) algorithm, which considers both the shape skeleton and the aforementioned special points. TS first divides a trajectory into walking and nonwalking segments using a trajectory segmentation algorithm [Zheng et al. 2008a] (see Section 3.4). A point is weighted by its heading change degree and the distance to its neighbors.

Another branch of research [Kellaris et al. 2009; Song et al. 2014] considers trajectory compression with the constraints of transportation networks. For example, we can reduce the redundant points on the same road segment. We can even discard all the newly acquired points after an anchor point, as long as the moving object is traveling on the shortest path from the anchor point to its current location. This branch of work usually needs the support of map-matching algorithms (refer to Section 3.5). In 2014, PRESS [Song et al. 2014] was proposed to separate the spatial representation of a trajectory from its temporal representation. PRESS consists of a hybrid spatial compression algorithm and an error bounded temporal compression algorithm, compressing the spatial and temporal information of trajectories, respectively. The spatial compression combines frequent sequential pattern mining techniques with Huffman Coding to reduce the size of a trajectory; that is, a frequently traveled path can be represented by a shorter code, therefore saving storage.

### 3.4. Trajectory Segmentation

In many scenarios, such as trajectories clustering and classification, we need to divide a trajectory into segments for a further process. The segmentation not only reduces the computational complexity but also enables us to mine richer knowledge, such as subtrajectory patterns, beyond what we can learn from an entire trajectory. In general, there are three types of segmentation methods.

The first category is based on *time interval*. For example, as illustrated in Figure 6(a), if the time interval between two consecutive sampling points is larger than a given



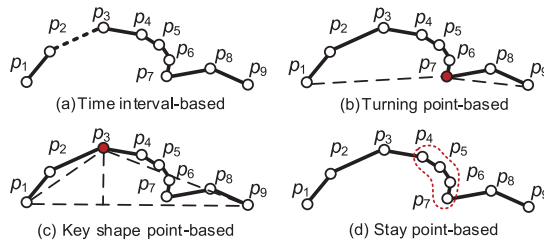


Fig. 6. Methods of trajectory segmentation.

threshold, a trajectory is divided into two parts at the two points, that is,  $p_1 \rightarrow p_2$  and  $p_3 \rightarrow \dots \rightarrow p_9$ . Sometimes, we can divide a trajectory into segments of the same time length.

The second category of methods is based on the *shape of a trajectory*. For example, as demonstrated in Figure 6(b), we can partition a trajectory by the turning points with heading direction changing over a threshold. Alternatively, we can employ the line simplification algorithms, such as the Douglas-Peucker algorithm, to identify the key points maintaining a trajectory's shape, as depicted in Figure 6(c). The trajectory is then partitioned into segments by these key points. Similarly, Lee et al. [2007] proposed to partition a trajectory by using the concept of Minimal Description Language (MDL), which is comprised of two components:  $L(H)$  and  $L(D|H)$ .  $L(H)$  is the length, in bits, of the description of the hypothesis  $H$ ; and  $L(D|H)$  is the length, in bits, of the description of the data when encoded with the help of the hypothesis. The best hypothesis  $H$  to explain  $D$  is the one that minimizes the sum of  $L(H)$  and  $L(D|H)$ . More specifically, they use  $L(H)$  to denote the total length of partitioned segments (like  $\overline{p_1 p_2}$  and  $\overline{p_1 p_9}$ ), while letting  $L(D|H)$  represent the total (perpendicular and angle) distance between the original trajectory and the new partitioned segments. Using an approximation algorithm, they find a list of characteristic points that minimize  $L(H) + L(D|H)$  from a trajectory. The trajectory is partitioned into segments by these characteristic points.

The third category of methods is based on the *semantic meanings* of points in a trajectory. As illustrated in Figure 6(d), a trajectory can be divided into segments, that is,  $p_1 \rightarrow p_2 \rightarrow p_3$  and  $p_8 \rightarrow p_9$ , based on the stay points it contains. Whether we should keep the stay points in the divided results depends on applications. For example, in a task of travel speed estimation, we should remove the stay points (from a taxi's trajectory) where a taxi was parked to wait for passengers [Yuan et al. 2013b]. On the contrary, to estimate the similarity between two users [Lee et al. 2008], we can only focus on the sequences of stay points, while skipping other raw trajectory points between two consecutive stay points.

Another semantic meaning-based trajectory segmentation is to divide a trajectory into segments of different transportation modes, such as driving, taking a bus, and walking. For example, Zheng et al. [2008a, 2008b, 2010c] proposed a walk-based segmentation method. The key insight is that people have to walk through the transition between two different transportation modes. Consequently, we can first distinguish *walk points* from *non-walk points* in a trajectory, based on a point's speed ( $p \cdot v$ ) and acceleration ( $p \cdot a$ ). The trajectory can then be divided into alternate *Walk Segments* and *non-Walk Segments*, as illustrated in Figure 7(a). In reality, however, as shown in Figure 7(b), a few points from *non-Walk Segments* may be detected as possible walk points, for example, when a bus moves slowly in traffic congestion. On the other hand, due to the locative error, a few points from walk segments might exceed the upper bound of travel speed ( $v_t$ ), therefore being recognized as non-walk points. To address this issue, a segment is merged into its backward segment, if the distance or time span

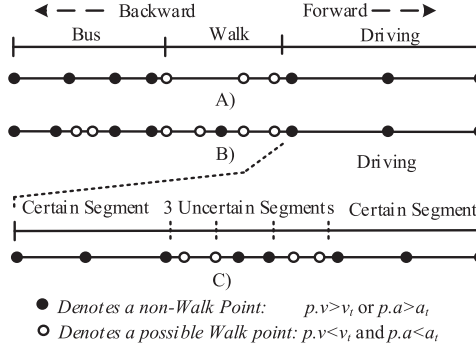


Fig. 7. Change point-based segmentation method.

of the segment is less than a threshold. After that, a segment is regarded as a *Certain Segment* if its length exceeds a threshold, as presented in Figure 7(c). Otherwise, it is deemed as an *Uncertain Segment*. As common users do not frequently change their transportation modes within a short distance, uncertain segments are merged into one non-walk segment if the number of consecutive uncertain segments exceeds a certain threshold (three in this example). Later, features are extracted from each segment to determine its exact mode.

### 3.5. Map Matching

Map matching is a process to convert a sequence of raw latitude/longitude coordinates to a sequence of road segments. Knowledge of which road a vehicle was/is on is important for assessing traffic flow, guiding the vehicle's navigation, predicting where the vehicle is going, and detecting the most frequent travel path between an origin and a destination, and so forth. Map matching is not an easy problem, given parallel roads, overpasses, and spurs [Krumm 2011]. There are two approaches to classify map-matching methods, based on the additional information used, or the range of sampling points considered in a trajectory.

According to the additional information used, map-matching algorithms can be categorized into four groups: *geometric* [Greenfeld 2002], *topological* [Chen et al. 2003; Yin and Wolfson 2004], *probabilistic* [Ochieng et al. 2004; Pink and Hummel 2008; Quddus et al. 2006], and other *advanced* techniques [Lou et al. 2009; Newson and Krumm 2009; Yuan et al. 2010b]. Geometric map-matching algorithms consider the shape of individual links in a road network, for example, matching a GPS point to the nearest road. Topological algorithms pay attention to the connectivity of a road network. Representative algorithms are those that use the Fréchet distance to measure the fit between a GPS sequence and candidate road sequence [Brakatsouls et al. 2005]. To deal with noisy and low-sampling rate trajectories, probabilistic algorithms [Ochieng et al. 2004; Pink and Hummel 2008; Quddus et al. 2006] make explicit provisions for GPS noise and consider multiple possible paths through the road network to find the best one. More advanced map-matching algorithms have emerged recently that embrace both the topology of the road network and the noise in trajectory data, exemplified by Lou et al. (2009), Newson and Krumm [2009], and Yuan et al. [2010b]. These algorithms find a sequence of road segments that simultaneously come close to the noisy trajectory data and form a reasonable route through the road network.

According to the range of sampling points considered, map-matching algorithms can be classified into two categories: *local/incremental* and *global* methods. The local/incremental algorithms [Civili et al. 2005; Chawathe 2007] follow a greedy strategy of

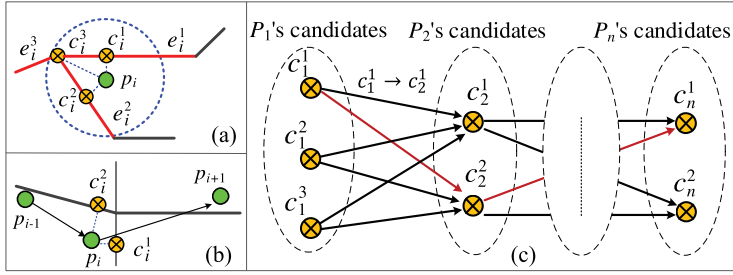


Fig. 8. An advanced map-matching algorithm.

sequentially extending the solution from an already matched portion. These methods try to find a local optimal point based on the distance and orientation similarity. Local/incremental methods run very efficiently, often adopted in online applications. However, when the sampling rate of a trajectory is low, the matching accuracy degrades. Instead, *Global* algorithms [Alt et al. 2003; Brakatsouls et al. 2005] aim to match an entire trajectory with a road network, for example, considering the predecessors and successors of a point. Global algorithms are more accurate, but less efficient, than *local* methods, usually applied to offline tasks (e.g., mining frequent trajectory patterns), where entire trajectories have already been generated.

Advanced algorithms [Lou et al. 2009; Newson and Krumm 2009; Yuan et al. 2010b] embrace local and global information (or geometric, topological, and probability) to deal with the mapping of a low-sampling-rate trajectory. As shown in Figure 8(a), the algorithm proposed in Lou et al. [2009] first finds the local candidate road segments that are within a circle distance to each point in a trajectory. For instance, road segments  $e_i^1$ ,  $e_i^2$ , and  $e_i^3$  are within the circle distance to  $p_i$ , and  $c_i^1$ ,  $c_i^2$ , and  $c_i^3$  are the candidate points on these road segments. The distance between  $p_i$  and a candidate point  $\text{dist}(c_i^j, p_i)$  indicates the probability  $N(c_i^j)$  that  $p_i$  can be matched to the candidate point. This *probability* can be regarded as the *local* and *geometric* information, which is modeled by a normal distribution:

$$N(c_i^j) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{\text{dist}(c_i^j, p_i)^2}{2\sigma^2}}.$$

The algorithm also considers the transition probability between the candidate points of each two consecutive trajectory points. For example, as depicted in Figure 8(b),  $c_i^2$  is more likely to be the true match of  $p_i$ , considering  $p_{i-1}$  and  $p_{i+1}$ . The transition probability between two candidate points is denoted by the ratio between their Euclidean distance and the road network distance. The transition is actually based on the *topologic* information of a road network. Finally, as shown in Figure 8(c), combining the local and transition probabilities, the map-matching algorithm finds a path (on a candidate graph) that maximizes the *global* probability of matching. The idea is similar to the hidden Markov model where emission and transition probabilities are considered to find the most possible sequence of status given a sequence of observations [Newson and Krumm 2009].

#### 4. TRAJECTORY DATA MANAGEMENT

Mining massive trajectories is very time consuming, as we need to access different samples of the trajectories or different parts of a trajectory many times. This calls for effective data management techniques that can quickly retrieve the trajectories (or parts of a trajectory) needed. Different from moving object databases that are concerned with the current location of a moving object, the trajectory data management

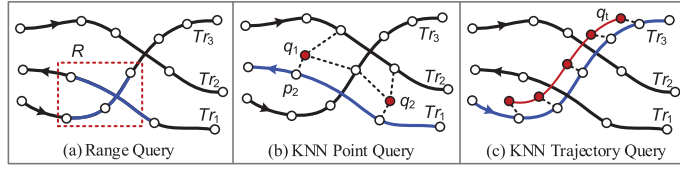


Fig. 9. Two categories of queries for trajectory data.

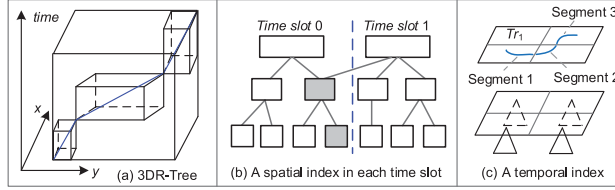


Fig. 10. Three approaches answering range queries.

introduced in this section deals with the traveling history of a moving object. A more comprehensive survey on trajectory data management can be found in Deng et al. [2011].

#### 4.1. Trajectory Indexing and Retrieval

There are two major types of queries: K-Nearest Neighbor (KNN) queries and Range queries, as depicted in Figure 9.

*Range queries* retrieve the trajectories falling into (or intersecting) a spatial (or spatiotemporal) range. For example, as shown in Figure 9(a), a range query can help us retrieve the trajectories of vehicles passing a given rectangular region  $R$  between 2pm and 4pm in the past month. The retrieved trajectories (or segments) can then be used to derive features, such as the travel speed and traffic flow, for data mining tasks like classification and prediction. There are three approaches to answering such kind of spatiotemporal range queries.

The first approach regards the time as the third dimension besides the 2D geographical space, building a 3D-Rtree based on trajectories, as depicted in Figure 10(a). A spatiotemporal range query is then formulated as a three-dimensional (3D) query box. So, answering such a query means finding the nodes on a 3D-Rtree within the 3D query box. The 3D-Rtree works well for indexing trajectories generated in the near recent (e.g., in the past few hours). When the time span of the trajectories to be indexed lasts for a long period (i.e., more segments of newly generated trajectories will be inserted into a 3D-Rtree index), however, the overlap among 3D boxes bounding segments of different trajectories occurs more often. This results in a frequent update of indexing structure and a significant increase of node accesses when retrieving a trajectory. Though ST-R-tree and TB-tree [Pfoser et al. 2000] have been proposed to address this issue, the overlap among different 3D boxes still keeps on increasing as time goes by.

The second approach divides a time period into multiple time intervals, building an individual spatial index like R-tree for the trajectories generated in each interval. The part of indexing structure that does not change over time is shared by two time slots. Representative indexing structures are multiple version R-tree, such as Rt-Tree [Xu et al. 1990], HR-Tree [Tao and Papadias 2001], and H+R-Tree [Tao and Papadias 2001], as illustrated in Figure 10(b). Given a spatiotemporal range query, such an index first finds the time slots falling in the temporal range, and then retrieves the trajectories intersecting the spatial query range from each spatial index of these time slots.

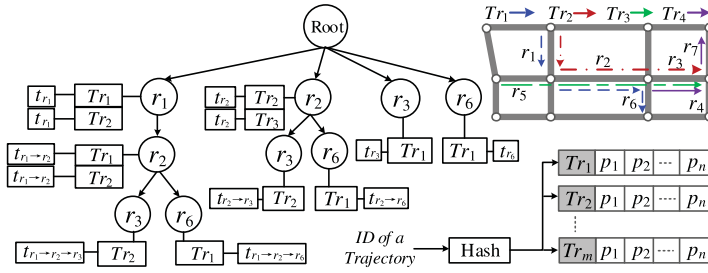


Fig. 11. Suffix-tree-like index for maintaining trajectories.

The third approach partitions a geographical space into grids, and then builds a temporal index for the trajectories falling in each grid. As shown in Figure 10(c), CSE-tree [Wang et al. 2008] divides a trajectory into several segments by the grids. Each segment falling in a grid is represented by a 2D point whose coordinates are the starting time and ending time of the segment. These points are then indexed by a hybrid B+tree. When retrieving trajectories satisfying a spatiotemporal query, CSE-tree first finds the grids intersecting the spatial range of the query, and then searches the hybrid B+tree of these grids for the segments of trajectories falling in the temporal range of the query. Finally, CSE-tree merges the IDs of trajectory segments (and their starting and ending times) retrieved from different grids.

*KNN queries* retrieve the top- $K$  trajectories with the minimum aggregate distance to a few points (entitled the KNN point query [Chen et al. 2010; Tao et al. 2002; Tang et al. 2011]) or a specific trajectory (entitled the KNN trajectory query [Yi et al. 1998; Agrawal et al. 1993]).

As depicted in Figure 9(b), an example of the KNN point query is to retrieve the trajectories of vehicles that are close to two given restaurants (e.g.,  $q_1$  and  $q_2$ ). Sometimes, the order between the query points is also considered [Chen et al. 2010], for example, finding the top- $k$  nearest trajectories first passing  $q_1$  and then  $q_2$ . Without the order,  $Tr_1$  is the nearest trajectory to the two points. However,  $Tr_2$  becomes the nearest after considering the order. The KNN point queries concern more about whether a trajectory provides a good connection to query locations rather than whether the trajectory is similar to the query in shape. Additionally, the number of query points is usually very small and can be far away from each other in applications. As a result, we cannot connect these query points sequentially to formulate a trajectory and then call the solution designed for the KNN trajectory query to solve it.

As illustrated in Figure 9(c), a KNN trajectory query can find the GPS logs of people traveling through a specific route. To answer such a query, the first step is to define a similarity/distance function between two trajectories. Then efficient query processing algorithms are designed to address the problem of searching over a large set of candidate trajectories. Sometimes, we need to retrieve the trajectories of vehicles traversing a specific path. There are two ways to achieve the goal.

One is to regard a path on a road network as a trajectory and use the KNN trajectory query to detect the trajectories that are close to the path. The other way is first to convert a trajectory into a sequence of road segments by using a map-matching algorithm. Some indexing structures are then built to manage the relationship between paths and the trajectories passing them. Figure 11 presents a suffix-tree-based indexing structure [Wang et al. 2014] that manages the four trajectories  $Tr_1$ ,  $Tr_2$ ,  $Tr_3$ , and  $Tr_4$  traversing a road network. Here, each node in the indexing tree stands for a road segment; each path on the tree corresponds to a route on the road network. Each node stores the IDs and travel times of the trajectories that traverse the path from the root



to the node. For example,  $t_{r_1 \rightarrow r_2 \rightarrow r_3}$  stands for the time for traveling path  $r_1 \rightarrow r_2 \rightarrow r_3$ . By searching for the tree, we can easily get the IDs of trajectories passing a path and retrieve the points of each trajectory through a hash table (as shown in the bottom-right part of Figure 11). The detailed content of each trajectory can be stored on disk if the memory is not big enough. Because the size of the index grows quickly as the number of trajectories increases, such index is only suitable for managing trajectories generated recently.

#### 4.2. Distance/Similarity of Trajectories

When answering KNN queries or clustering trajectories, we need to calculate the distance (alternatively, we can say similarity) between a trajectory and a few points, or the distance between two trajectories.

The distance between a point  $q$  and a trajectory  $A$  is usually measured by the distance from  $q$  to its nearest point in  $A$ , denoted as  $D(q, A) = \min_{p \in A} D(p, q)$ ; for example,  $q_1$  and  $p_2$  shown in Figure 9(b). An approach extending the distance from a single point  $q$  to multiple query points  $Q$  is  $D(Q, A) = \sum_{q \in Q} e^{D(q, A)}$ , or  $S(Q, A) = \sum_{q \in Q} e^{-D(q, A)}$ , written in a similarity fashion. The intuition of using the exponential function is to assign a larger contribution to a closer matched pair of points while giving much lower value to those faraway pairs. Chen et al. [2010] define the best connect distance, which can measure the distance between a trajectory and a few points with or without an order.

The distance between two trajectories is usually measured by some kind of aggregation of distances between trajectory points. Closest-pair distance uses the minimal distance between the points in two trajectories  $(A, B)$  to represent the similarity of trajectories, that is,  $CPD(A, B) = \min_{p \in A, p' \in B} D(p, p')$ . Assuming that two trajectories are of the same length, sum-of-pairs distance uses the sum of corresponding points from the two trajectories to denote the distance, that is,  $SPD(A, B) = \sum_{i=1}^n D(p_i, p'_i)$ . As the assumption may not hold in reality, Dynamic Time Wrapping (DTW) distance was proposed to allow “repeating” some points as many times as needed in order to get the best alignment [Agrawal et al. 1993]. As some noise points from a trajectory may cause a big distance between trajectories, the concept of the Longest Common Subsequence (LCSS) is employed to address this issue. The LCSS-based distance allows one to skip some noise points when calculating the distance of trajectories, using a threshold  $\delta$  to control how far in time we can go in order to match one point from a trajectory to a point in another trajectory. Another threshold  $\varepsilon$  is used to determine whether two points (from two different trajectories) are matched. Chen and Ng [2004] proposed the Edit Distance on Real Sequence (EDR) distance, which is similar to LCSS in using a threshold  $\varepsilon$  to determine a match, while assigning penalties to the gaps between two matched subtrajectories. Chen et al. [2005] also proposed the Edit Distance with Real Penalty (ERP) distance aiming to combine the merits of DTW and EDR, by using a constant reference point for computing distance. Note that DTW is not a metric, as it does not satisfy the triangle inequality. EDR is metric, and thus can be used to prune unnecessary trajectories.

Basically, LCSS and Edit Distance were proposed for matching strings. When used to match two trajectories, there is a threshold  $\varepsilon$  need to set; this is not easy. K-BCT [Chen et al. 2010] is a parameter-free similarity metric for trajectories, combining the merits of DTW and LCSS. During the matching process, K-BCT can repeat some trajectory points and skip unmatched trajectory points including outliers.

*The Distance between Two Trajectory Segments:* A distance measure for trajectory segments is based on the Minimum Bounding Rectangles (MBR) of segments [Jeung et al. 2011]. As demonstrated in Figure 12(a), the MBRs of two segments  $(L_1, L_2)$  are  $(B_1, B_2)$ , each of which is described by the coordinates of the low bound point  $(x_i, y_i)$

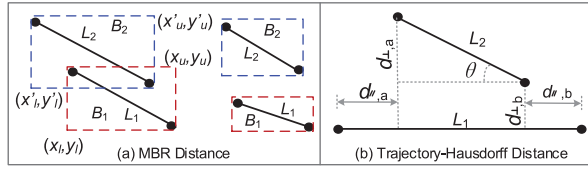


Fig. 12. Distance metrics for trajectory segments.

and upper bound point  $(x_u, y_u)$ . The MBR-based distance  $D_{min}(B_1, B_2)$  is defined as the minimum distance between any two points from  $(B_1, B_2)$ , calculated as

$$\sqrt{(\Delta([x_l, x_u], [x'_l, x'_u]))^2 + (\Delta([y_l, y_u], [y'_l, y'_u]))^2},$$

where the distance between two intervals is defined as

$$\Delta([x_l, x_u], [x'_l, x'_u]) = \begin{cases} 0 & [x_l, x_u] \cap [x'_l, x'_u] \neq \emptyset \\ x'_l - x_u & x'_l > x_u \\ x_l - x'_u & x_l > x'_u \end{cases}.$$

In the two examples shown in Figure 12(a), the distance between  $L_1$  and  $L_2$  is 0 and  $y'_l - y_u$ , respectively.

As depicted in Figure 12(b), Lee et al. [2007] proposed a distance function, entitled Trajectory-Hausdorff Distance ( $D_{Haus}$ ), which is a weighted sum of three terms: (1) The aggregate perpendicular distance ( $d_{\perp}$ ) that measures the separation between two trajectories, (2) the aggregate parallel distance ( $d_{\parallel}$ ) that captures the difference in length between two trajectories, and (3) the angular distance ( $d_{\theta}$ ) that reflects the orientation difference between two trajectories. Formally,

$$D_{Haus} = w_1 d_{\perp} + w_2 d_{\parallel} + w_3 d_{\theta},$$

where  $d_{\perp} = \frac{d_{\perp,a}^2 + d_{\perp,b}^2}{d_{\perp,a} + d_{\perp,b}}$ ,  $d_{\parallel} = \min(d_{\parallel,a}, d_{\parallel,b})$ ,  $d_{\theta} = ||L_2|| \cdot \sin\theta$ , and  $w_1, w_2$ , and  $w_3$  are weights depending on applications.

## 5. UNCERTAINTY IN A TRAJECTORY

As the location of a moving object is recorded at a certain time interval, the trajectory data we obtain is usually a sample of the object's true movement. On one hand, the movement of an object between two consecutive sampling points becomes unknown (or called uncertain). To this end, we expect to reduce the uncertainty of a trajectory. On the other hand, in some applications, to protect a user's privacy that could be leaked from her trajectories, we need to make a trajectory even more uncertain.

### 5.1. Reducing Uncertainty from Trajectory Data

Many trajectories have been recorded with a very low sampling rate, leading to an object's movement between sampling points uncertain; we call them uncertain trajectories. For instance, as shown in Figure 13(a), the GPS coordinates of a taxi ( $p_1, p_2, p_3$ ) were recorded every few minutes to reduce communication loads, resulting in multiple possible paths between two consecutive sampling points. As illustrated in Figure 13(b), people's check-in records in a location-based social networking service like FourSquare can be regarded as trajectories if we connect them chronologically. As people do not check in very often, the time interval (and distance) between two consecutive check-ins may be hours (and several kilometers). Consequently, we have no idea how a user traveled between two check-ins. As demonstrated in Figure 13(c), to save energy, the GPS

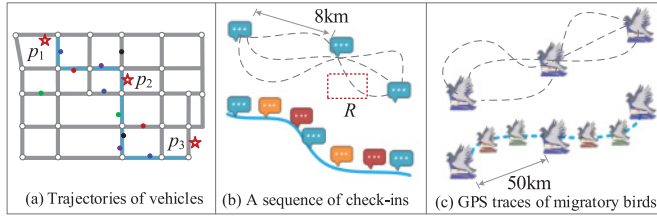


Fig. 13. Examples of uncertain trajectories.

logger installed on a migratory bird can only send a location record every half day. As a result, the path that a bird flew over two particular locations is quite uncertain.

**5.1.1. Modeling Uncertainty of a Trajectory for Queries.** Several models of uncertainty paired with appropriate query evaluation techniques [Pfoser and Jensen 1999; Cheng et al. 2008] have been proposed for moving object databases to answer queries, for example, “is it possible for an object to intersect a query window.” As illustrated in Figure 13(b), we do not know whether the trajectory formulated by the three blue check-ins should be retrieved or not by the range query  $R$ , without modeling the uncertainty of the trajectory. Many of these techniques aim at providing conservative bounds for the positions of uncertain objects between two sampling points. This is usually achieved by employing geometric objects, such as cylinders [Trajcevski et al. 2004, 2009] or beads [Trajcevski et al. 2010], as trajectory approximations. These models concern little about data mining, and therefore are not the focus of our article. Recent approaches use independent probability density functions at each point of time [Cheng et al. 2004], or stochastic processes [Qiao et al. 2010; Xu et al. 2013; Emrich et al. 2012; Niedermayer et al. 2014] (e.g., Markov chains), to better model the uncertain positions of an object and answer different queries.

**5.1.2. Path Inference from Uncertain Trajectories.** Different from the aforementioned models aiming at the retrieval of existing trajectories by different queries, a new series of techniques infers (or say “constructs”) the most likely  $k$  route(s) that a moving object could travel (i.e., the missing subtrajectory) between a few sample points based on a bunch of uncertain trajectories. The major insight is that trajectories sharing (or partially sharing) the same/similar routes can often supplement each other to make themselves more complete. In other words, it is possible to interpolate an uncertain trajectory by cross-referring other trajectories on (or partially on) the same/similar route, that is, “uncertain + uncertain  $\rightarrow$  certain.” For example, given the uncertain trajectories of many taxicabs (marked by different colored points in Figure 13(a)), we could infer that the blue path is the most likely route traversing ( $p_1, p_2, p_3$ ). Likewise, based on the check-in data of many users, as depicted in Figure 13(b), we could find the blue curve the most possible travel path between the three blue check-ins. Similarly, given the uncertain GPS traces of many birds, we can identify the path that birds fly over a few locations. Reducing the uncertainty of trajectories can support scientific studies and enable many applications, such as travel recommendation and traffic management. There are two categories of methods to complement an uncertain trajectory.

One is designed for the trajectories generated in a road network setting [Zheng et al. 2012a]. What set this category of methods apart from map-matching algorithms lies in two aspects. First, the methods for reducing the uncertainty of trajectories leverage the data from many other trajectories, while map-matching algorithms only use the geometric information from a single trajectory and the topological information of road networks. Second, the sampling rate of trajectories handled by the uncertainty methods

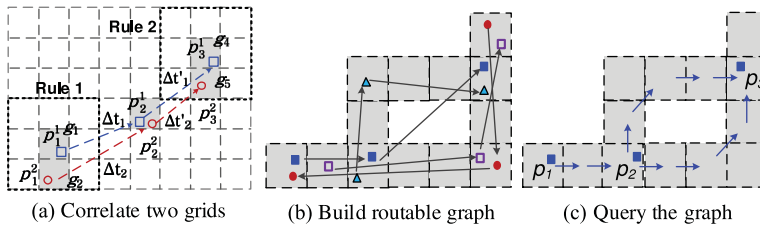


Fig. 14. The most likely route based on uncertain trajectories.

can be very low, for example, more than 10 minutes. This seems nearly impossible for a map-matching algorithm.

The other is for a free space, where moving objects (like flying birds or people hiking a mountain) do not follow paths in road networks [Wei et al. 2012], as illustrated in Figures 13(b) and 13(c). The major challenges are twofold. One is to determine those trajectories that may be relevant to a series of query points. The other is to construct a route that can approximate a bunch of relevant trajectories. As shown in Figure 14(a), the method proposed in Wei et al. [2012] first partitions a geospace into uniform grids (the size of a grid depends on the required inference accuracy), and then maps trajectories onto these grids. Some grids can be connected to form a region if the trajectories passing them satisfy one of the following two rules: (1) If the starting points ( $p_1^1, p_1^2$ ) of two trajectory segments are located in two grids ( $g_1, g_2$ ) that are geospatial neighbors, and the ending points ( $p_2^1, p_2^2$ ) of the two segments are located in the same grid, and the travel times ( $\Delta t_1, \Delta t_2$ ) of the two segments are similar, then the two grids ( $g_1, g_2$ ) can be connected. (2) If the starting points ( $p_1^1, p_1^2$ ) are located in the same grid, and ending points ( $p_3^1, p_3^2$ ) fall in the grids ( $g_4, g_5$ ) that are neighbors, travel times ( $\Delta t_1', \Delta t_2'$ ) of the two segments are similar, then grids ( $g_4, g_5$ ) can be connected.

After turning disjoint grids into connected region(s), as demonstrated in Figure 14(b), we can build a routable graph where a node is a grid. The direction and travel time between two adjacent grids in the graph is inferred based on the trajectories passing the two grids. Finally, as depicted in Figure 14(c), given three query points, we can find the most likely route on the graph based on a routing algorithm. To find a more detailed path, a regression can be performed over the trajectories passing the identified route.

Su et al. [2013] proposed an anchor-based calibration system that aligns trajectories to a set of fixed anchor points. The approach considers the spatial relationship between anchor points and trajectories. It also trains inference models from historical trajectories to improve the calibration.

## 5.2. Privacy of Trajectory Data

Instead of making a trajectory more certain, a series of techniques aim to protect a user from the privacy leak caused by the disclosure of the user's trajectories [Abul et al. 2008; Xue et al. 2013; Chow and Mokbel 2011]. This kind of technology tries to blur a user's location, while ensuring the quality of a service or the utility of the trajectory data. There are two major scenarios that we need to protect a user's trajectory data from the privacy leak.

One is in real-time continuous location-based services for example, tell me the traffic conditions that are 1 km around me. In this scenario, a user may not want to exactly disclose her current location when using a service. Different from the simple location privacy, the spatiotemporal correlation between consecutive samples in a trajectory may help infer the exact location of a user. Techniques trying to protect the privacy leak in this scenario include spatial cloaking [Mokbel et al. 2007], mix zones [Beresford

and Stajano 2003], path confusion [Hoh et al. 2010], Euler histogram based on short IDs [Xie et al. 2010], dummy trajectories [Kido et al. 2005], and so on.

The second is the publication of historical trajectories. Collecting many trajectories of an individual may allow attackers to infer her home and work places, therefore identifying who the individual is. Major techniques for protecting users' privacy in such scenario include clustering-based [Abul et al. 2008], generalization-based [Nergiz et al. 2009], suppression-based [Terrovitis and Mamoulis 2008], and grid-based [Gid'ofalvi et al. 2007] approaches. A comprehensive survey on trajectory privacy can be found in [Chow and Mokbel 2011].

## 6. TRAJECTORY PATTERN MINING

In this section, we study four major categories of patterns that can be discovered from a single trajectory or a group of trajectories. They are moving together patterns, trajectory clustering, sequential patterns, and periodic patterns.

### 6.1. Moving Together Patterns

This branch of research is to discover a group of objects that move together for a certain time period, such as *flock* [Gudmundsson and Kreveld 2006; Gudmundsson et al. 2004], *convoy* [Jeung et al. 2008a, 2008b], *swarm* [Li et al. 2010a], *traveling companion* [Tang et al. 2012a, 2012b], and *gathering* [Zheng et al. 2013; Zheng et al. 2014a]. These patterns can help the study of species' migration, military surveillance, and traffic event detection, and so on. These patterns can be differentiated between each other based on the following factors: the shape or density of a group, the number of objects in a group, and the duration of a pattern.

Specifically, a *flock* is a group of objects that travel together within a disk of some user-specified size for at least  $k$  consecutive time stamps. A major concern with flock is the predefined circular shape, which may not well describe the shape of a group in reality, and therefore may result in the so-called lossy-flock problem. To avoid rigid restrictions on the size and shape of a moving group, the *convoy* is proposed to capture generic trajectory pattern of any shape by employing the density-based clustering. Instead of using a disk, a convoy requires a group of objects to be density connected during  $k$  consecutive time points. While both flock and convoy have a strict requirement on consecutive time period, Li et al. [2010a] proposed a more general type of trajectory pattern, called *swarm*, which is a cluster of objects lasting for at least  $k$  (possibly nonconsecutive) time stamps. While convoy and swarm need to load entire trajectories into memory for a pattern mining, the *traveling companion* [Tang et al. 2012a] uses a data structure (called *traveling buddy*) to continuously find convoy/swarmlike patterns from trajectories that are being streamed into a system. So, the traveling companion patterns can be regarded as an online (and incremental) detection fashion of convoy and swarm.

To detect some incidents, such as celebrations and parades, in which objects join in and leave an event frequently, the *gathering* pattern [Zheng et al. 2013a, 2014a] further loses the constraints of the aforementioned patterns by allowing the membership of a group to evolve gradually. Each cluster of a gathering should contain at least  $m_p$  participators, which are the objects appearing in at least  $k_p$  clusters of this gathering. As the gathering pattern is used to detect events, it also requires the geometric property (like location and shape) of a detected pattern to be relatively stable.

Figure 15(a) illustrates these patterns. If the set requirement of time stamps is  $k = 2$ , a group  $\langle o_2, o_3, o_4 \rangle$  is a flock from  $t_1$  to  $t_3$ . Though  $o_5$  is a companion of the group, it cannot be included due to the fixed size of the disk employed by the flock definition. On the other hand, a convoy can include  $o_5$  into the group, since  $\langle o_2, o_3, o_4, o_5 \rangle$  is density-based connected from  $t_1$  to  $t_3$ . The five objects also form a swarm during the nonconsecutive



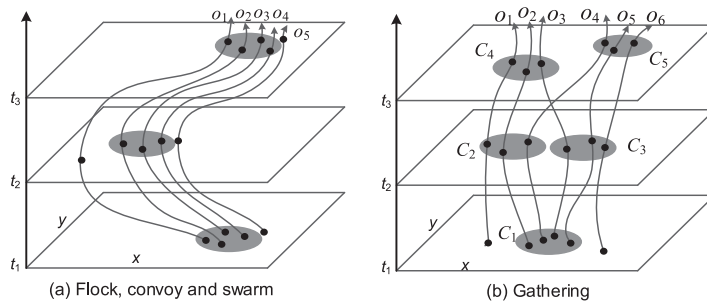


Fig. 15. Examples of moving together patterns.

time period  $t_1$  and  $t_3$ . As demonstrated in Figure 15(b), if we set  $k_p = 2$  and  $m_p = 3$ , then  $\langle C_1, C_2, C_4 \rangle$  is a gathering.  $\langle C_1, C_3, C_5 \rangle$  is not a gathering as  $C_5$  is too far away from  $C_2$  and  $C_3$ .

The aforementioned pattern mining algorithms usually use a density-based distance metric (in a Euclidean space) to find a cluster of moving objects. Jensen et al. [2007] extend the distance metric by considering semantic factors, such as heading directions and speed, of a moving object.

## 6.2. Trajectory Clustering

To find representative paths or common trends shared by different moving objects, we usually need to group similar trajectories into clusters. A general clustering approach is to represent a trajectory with a feature vector, denoting the similarity between two trajectories by the distance between their feature vectors. However, it is not easy to generate a feature vector with a uniform length for different trajectories, as different trajectories contain different and complex properties, such as length, shape, sampling rate, number of points, and their orders. In addition, it is difficult to encode the sequential and spatial properties of points in a trajectory into its feature vector.

Given the challenges mentioned previously, a series of technique works have been done. Since the distance metrics between trajectories have been introduced in Section 4.2, we hereafter focus on the clustering methods proposed for trajectories. Note that the clustering methods discussed in this section are dedicated for trajectories in free spaces (i.e., without a road network constraint). Though there are a few publications (e.g., Kharrat et al. [2008]) discussing the trajectory clustering in a road network setting, this problem can actually be solved by the combination of map-matching and graph clustering algorithms. That is, we can first use map-matching algorithms to project trajectories onto a road network and then employ graph clustering algorithms to find a subgraph (i.e., a collection of roads) on the road network.

Gaffney and Smyth [1999] and Cadez et al. [2000] proposed to group similar trajectories into clusters by using a regression mixture model and the Expectation-Maximization (EM) algorithm. This algorithm clusters trajectories with respect to the overall distance between two entire trajectories. However, moving objects rarely travel together for an entire path in the real world. To this end, Lee et al. [2007] proposed to partition trajectories into line segments and to build groups of close trajectory segments using the Trajectory-Hausdorff Distance, as illustrated in Figure 16(a). A representative path is later found for each cluster of segments. Since trajectory data are often received incrementally, Li et al. [2010b] further proposed an incremental clustering algorithm, aiming to reduce the computational cost and storage of received trajectories. Both Lee [2007] and Li [2010] adopted a *Micro-and-Macroclustering framework*, which was proposed by Aggarwal et al. [2003] to cluster data streams. That is, their methods

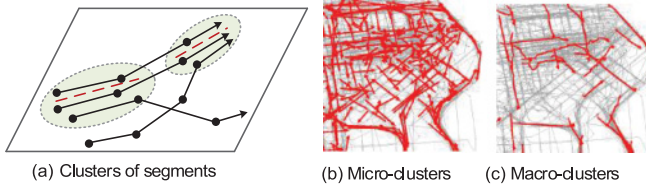


Fig. 16. Trajectory clustering based on partial segments [Li et al. 2010b].

first find microclusters of trajectory segments (as demonstrated in Figure 16(b)), and then group microclusters into macroclusters (as shown in Figure 16(c)). A major insight of Li's work [Li et al. 2010b] is that new data will only affect the local area where the new data were received rather than the faraway areas.

### 6.3. Mining Sequential Patterns from Trajectories

A branch of research is to find the *sequential patterns* from a single trajectory or multiple trajectories. Here, a sequential pattern means a certain number of moving objects traveling a common sequence of locations in a similar time interval. The locations in a travel sequence do not have to be consecutive. For instance, two trajectories  $A$  and  $B$ ,

$$A : l_1 \xrightarrow{1.5h} l_2 \xrightarrow{1h} l_7 \xrightarrow{1.2h} l_4, \quad B : l_1 \xrightarrow{1.2h} l_2 \xrightarrow{2h} l_4,$$

share a common sequence  $l_1 \rightarrow l_2 \rightarrow l_4$ , as the visiting orders and travel times are similar (though  $l_2$  and  $l_4$  is not consecutive in trajectory  $A$ ). When the occurrence of such a common sequence in a corpus, usually called support, exceeds a threshold, a *sequential trajectory pattern* is detected. Finding such kind of patterns can benefit travel recommendation [Zheng and Xie 2011b; Giannotti et al. 2007], life pattern understanding [Ye et al. 2009], next location prediction [Monreale et al. 2009], estimating user similarity [Xiao et al. 2014; Li et al. 2008], and trajectory compression [Song et al. 2014].

To detect the sequential patterns from trajectories, we first need to define a (common) location in a sequence. Ideally, in trajectory data, like user check-in sequences from a social networking service, each location is tagged with a unique identity (such as the name of a restaurant). If two locations share the same identity, they are common. In many GPS trajectories, however, each point is characterized by a pair of GPS coordinates, which do not repeat themselves exactly in every pattern instance. This makes the points from two different trajectories not directly comparable. In addition, a GPS trajectory may consist of thousands of points. Without handled properly, these points will result in a huge computational cost.

#### 6.3.1. Sequential Pattern Mining in a Free Space.

*Line-Simplification-Based Methods:* An early solution aiming to deal with the aforementioned issues was proposed in 2005 [Cao et al. 2005]. The solution first identifies key points shaping a trajectory, by using a line simplification algorithm like DP [Douglas and Peucker 1973]. It then groups the fragments of a trajectory that are close to each simplified line segment so as to count the support of each line segment. The travel time between two points in a trajectory is not considered.

*Clustering-Based Methods:* Recently, a more general way to solve the previously mentioned problems is to cluster points from different trajectories into regions of interest. A point from a trajectory is then represented by the cluster ID the point belongs to. As a consequence, a trajectory is re-formed as a sequence of cluster IDs, which are comparable among different trajectories. For example, as shown in Figure 17(a), the

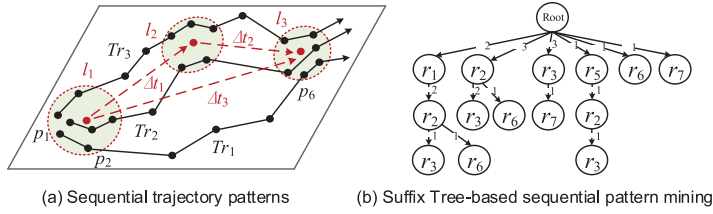


Fig. 17. Sequential pattern mining in trajectory data.

three trajectories can be represented as

$$Tr_1 : l_1 \xrightarrow{\Delta t_3} l_3, Tr_2 : l_1 \xrightarrow{\Delta t_1} l_2 \xrightarrow{\Delta t_2} l_3, Tr_3 : l_1 \xrightarrow{\Delta t'_1} l_2 \xrightarrow{\Delta t'_2} l_3,$$

where  $l_1$ ,  $l_2$ , and  $l_3$  are clusters of points. After the transformation, we can mine the sequential patterns from these sequences by using existing sequential pattern mining algorithms, such as PrefixSpan [Pei et al. 2011] and CloseSpan [Yan et al. 2003], with time constraints. In this example, setting the support threshold to 3, we can find  $l_1 \rightarrow l_3$  is a sequential pattern if

$$\frac{|\Delta t_3 - (\Delta t_1 + \Delta t_2)|}{\max(\Delta t_3, \Delta t_1 + \Delta t_2)} < \rho \text{ and } \frac{|\Delta t_3 - (\Delta t'_1 + \Delta t'_2)|}{\max(\Delta t_3, \Delta t'_1 + \Delta t'_2)} < \rho,$$

where  $\rho$  is a ratio threshold guaranteeing that two travel times are similar. Likewise, setting the threshold of support to 2,  $l_1 \rightarrow l_2 \rightarrow l_3$  is a sequential pattern, if  $\Delta t_1$  is similar to  $\Delta t'_1$  and  $\Delta t_2$  is similar to  $\Delta t'_2$ . Towards this direction, Giannotti et al. [2007] divide a city into uniform grids, grouping these grids into regions of interest based on the density of GPS points fallen in each grid. An a priori-like algorithm is then proposed to detect the sequential patterns of the region of interest.

With respect to the applications caring more about the semantic meaning of a location, we can first detect stay points from each trajectory, turning a trajectory into a sequence of stay points (see Section 3.2). Later, we can cluster these stay points to formulate regions of interest and use the cluster ID that a stay point belongs to represent a trajectory. Following this strategy, Ye et al. [2009] proposed to mine life patterns from an individual's GPS trajectories. Xiao et al. [2010, 2014] proposed a graph-based sequence matching algorithm to find the sequential pattern shared by two users' trajectories. These patterns are then used to estimate the similarity between two users.

**6.3.2. Sequential Pattern Mining in a Road Network.** When the sequential pattern mining problem is applied to a road network setting, we can first map each trajectory onto a road network by using map-matching algorithms. A trajectory is then represented by a sequence of road segment IDs, which can be regarded as strings. As a result, some sequential pattern mining algorithms, such as LCSS and Suffix Tree, designed for strings can be adapted to finding sequential trajectory patterns. Figure 17(b) presents a suffix tree that represents the four trajectories depicted in Figure 11. Here, a node is a road segment, and the path from the root to a node corresponds to a suffix of the string representing a trajectory. For example,  $Tr_1$  is represented by a string  $r_1 \rightarrow r_2 \rightarrow r_6$ , where  $r_2 \rightarrow r_6$  and  $r_6$  are suffixes of the string. The number associated with each link denotes the number of trajectories traversing the path, that is, the support of the string pattern. For instance, there are two trajectories ( $Tr_1$  and  $Tr_2$ ) traversing  $r_1 \rightarrow r_2$  and one trajectory traversing  $r_1 \rightarrow r_2 \rightarrow r_6$ . After building such a suffix tree, we can find the frequent patterns (i.e., the paths on the tree) with a support greater than a given threshold, with a complexity of  $O(n)$ . Note that the size of a suffix tree can be much bigger than the original trajectories. So, when the size of a trajectory dataset is

very large, we need to set a constraint on the depth of its suffix tree. Additionally, the sequential patterns derived from the suffix tree have to be consecutive. Though the temporal constraint is not explicitly considered, two objects' travel times on the same path should be similar, given the speed constraint of a path.

Towards this direction, Song et al. [2014] use Suffix Tree to detect frequent trajectory patterns, which are then leveraged to compress trajectories in conjunction with Huffman Encoding. Wang et al. [2014] employ Suffix Tree to find frequent trajectory patterns, which are used to reduce the candidates of a combination of subtrajectories when estimating the travel time of a query path.

#### 6.4. Mining Periodical Patterns from Trajectories

Moving objects usually have periodic activity patterns. For example, people go shopping every month and animals migrate yearly from one place to another. Such periodic behaviors provide an insightful and concise explanation over a long moving history, helping compress trajectory data and predict the future movement of a moving object.

Periodic pattern mining has been studied extensively for time series data. For example, Yang et al. tried to discover asynchronous patterns [Yang et al. 2003], surprising periodic patterns [Yang et al. 2001], and patterns with gap penalties [Yang et al. 2002], from (categorical) time series. Due to the fuzziness of spatial locations, existing methods designed for time series data are not directly applicable to trajectories. To this end, Cao et al. [2007] proposed an efficient algorithm for retrieving maximal periodic patterns from trajectories. This algorithm follows a paradigm that is similar to *frequent* pattern mining, where a (global) minimum support threshold is needed. In the real world, however, periodic behaviors could be more complicated, involving multiple interleaving periods, partial time span, and spatiotemporal noises and outliers.

To deal with these issues, Li et al. [2010c] proposed a two-stage detection method for trajectory data. In the first stage, the method detects a few reference spots, where a moving object has visited frequently, by using a density-based clustering algorithm, such as KDE. The trajectory of a moving object is then transformed into several binary time series, each of which indicates the “in” (1) and “out” (0) status of the moving object at a reference spot. Through applying Fourier transform and autocorrelation methods to each time series, the values of periods at each reference spot can be calculated. The second stage summarizes the periodic behaviors from partial movement sequences by using a hierarchical clustering algorithm. In 2012, Li et al. [2012] further extend the research [Li et al. 2010c] to mining periodic patterns from incomplete and sparse data sources.

### 7. TRAJECTORY CLASSIFICATION

Trajectory classification aims to differentiate between trajectories (or its segments) of different status, such as motions, transportation modes, and human activities. Tagging a raw trajectory (or its segment) with a semantic label raises the value of trajectories to the next level, which can facilitate many applications, such as trip recommendation, life experiences sharing, and context-aware computing.

In general, trajectory classification is comprised of three major steps: (1) Divide a trajectory into segments using segmentation methods. Sometimes, each single point is regarded as a minimum inference unit. (2) Extract features from each segment (or point). (3) Build a model to classify each segment (or point). As a trajectory is essentially a sequence, we can leverage existing sequence inference models, such as Dynamic Bayesian Network (DBN), HMM, and Conditional Random Field (CRF), which incorporate the information from local points (or segments) and the sequential patterns between adjacent points (or segments).

Using a sequence of 802.11 radio signals, LOCADIO [Krumm and Horvitz 2004] employs a hidden Markov model to classify the motion of a device into two statuses: Still

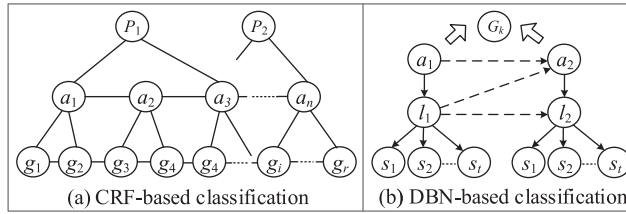


Fig. 18. Trajectory classification for activity recognition.

and Moving. Based on a trajectory of GSM signals, Timothy et al. [2006] attempted to classify the mobility of a user into three statuses, consisting of stationary, walking, and driving. Zhu et al. [2011] aim to infer the status of a taxi, consisting of *Occupied*, *Nonoccupied*, and *Parked*, according to its GPS trajectories. They first seek the possible *Parked* places in a trajectory, using a stay point-based detection method. A taxi trajectory is then partitioned into segments by these *Parked* places (refer to Figure 6(d) for an example). For each segment, they extract a set of features incorporating the knowledge of a single trajectory, historical trajectories of multiple taxis, and geographic data like road networks and Points of Interest (POIs). After that, a two-phase inference method is proposed to classify the status of a segment into either *Occupied* or *Nonoccupied*. The method first uses the identified features to train a local probabilistic classifier and then globally considers travel patterns via a hidden semi-Markov model.

Zheng et al. [2008a, 2008b] classify a user's trajectory by transportation modes, which is comprised of *Driving*, *Biking*, *Bus*, and *Walking*. As people usually change transportation modes in a single trip, a trajectory is first partitioned into segments based on the *Walk*-based segmentation method (refer to Figure 7 for details). A set of features, such as the heading change rate, stop rate, and velocity change rate, are extracted, being fed into a Decision Tree Classifier. Based on the inference results, a graph-based postprocessing step is conducted to fix the possibly wrong inference, considering the transition probability between different transportation modes at different places.

Liao et al. [2004] and Patterson et al. [2003] proposed a hierarchical inference model for location-based activity recognition and significant place discovery, as shown in Figure 18(a). A GPS trajectory is first divided into 10-m segments, each of which is then projected onto corresponding street patches by using a CRF-based map-matching algorithm. Based on the features extracted from these street patches, the model classifies a sequence of GPS points into a sequence of activities like  $a_1, a_2, \dots, a_n$  (such as *Walk*, *Driving*, and *Sleep*) and identifies a person's significant places like  $P_1$  and  $P_2$  (e.g., home, work, and bus stops), simultaneously. Yin et al. [2004] proposed a DBN-based inference model to infer a user's activities as well as high-level goals, according to a sequence of WiFi signals. Figure 18(b) presents the structure of the DBN, where the bottom layer contains the input of raw WiFi signals; the second layer is a list of locations where these signals are received; the top level corresponds to user activities. Finally, the high-level goal is inferred based on the sequence of inferred activities.

## 8. ANOMALIES DETECTION FROM TRAJECTORIES

Trajectory outliers (a.k.a. anomalies) can be items (e.g., a trajectory or a segment of trajectory) that are significantly different from other items in terms of some similarity metric. It can also be events or observations (represented by a collection of trajectories) that do not conform to an expected pattern (e.g., a traffic congestion caused by a car accident). A survey on general anomaly detection methods can be found in Chandola et al. [2009].



### 8.1. Detecting Outlier Trajectories

An outlier trajectory is a trajectory or a part of a trajectory that is significantly different from others in a corpus in terms of a distance metric, such as shape and travel time. The outlier trajectories could be a taxi driver's malicious detour [Liu et al. 2014; Zhang et al. 2011] or unexpected road changes (due to traffic accidents or construction). It can also remind people when traveling on a wrong path.

A general idea is to leverage existing trajectory clustering or frequent pattern mining methods. If a trajectory (or a segment) cannot be accommodated in any (density-based) clusters, or not frequent, it may be an outlier. Lee et al. [2008] proposed a partition-and-detection framework to find anomalous segments of trajectories from a trajectory dataset. This method can be an extension of the trajectory clustering proposed in Lee et al. [2007].

### 8.2. Identifying Anomalous Events by Trajectories

Another direction is to detect traffic anomalies (rather than trajectory itself) by using many trajectories. The traffic anomalies could be caused by accidents, controls, protests, sports, celebrations, disasters, and other events.

Liu et al. [2011] partition a city into disjointed regions with major roads and glean the anomalous links between two regions according to the trajectories of vehicles traveling between the two regions. They divide a day into time bins and identify for each link three features: the number of vehicles traveling a link in a time bin, the proportion of these vehicles among all vehicles entering the destination region, and that departing from the origin region. The three features of a time bin were respectively compared with those in the equivalent time bins of previous days to calculate the minimum distortion of each feature. Then, the link of the time bin can be represented in a 3D space, with each dimension denoting the minimum distort of a feature. Later, the Mahalanobis distance is used to measure the extreme points (in the 3D space), which are regarded as outliers. Following the aforementioned research, Chawla et al. [2012] proposed a two-step mining and optimization framework to detect traffic anomalies between two regions and explain an anomaly with the traffic flows passing the two regions (see Section 10 for details).

Pan et al. [2013] identify traffic anomalies according to drivers' routing behavior on an urban road network. Here, a detected anomaly is represented by a subgraph of a road network where drivers' routing behaviors significantly differ from their original patterns. They then tried to describe the detected anomaly by mining representative terms from the social media that people have posted when the anomaly was happening.

Pang et al. [2011, 2013] adapt likelihood ratio tests, which have previously been used in epidemiological studies, to describe traffic patterns. They partitioned a city into uniform grids and counted the number of vehicles arriving in a grid over a time period. The objective is to identify contiguous set of cells and time intervals that have the largest statistically significant departure from expected behavior (i.e., the number of vehicles). The regions whose log-likelihood ratio statistic value drops in the tail of  $\chi^2$  distribution are likely to be anomalous [Chandola et al. 2009].

## 9. TRANSFER TRAJECTORY TO OTHER REPRESENTATIONS

### 9.1. From Trajectory to Graph

Trajectories can be transformed into other data structures, besides being processed in its original form. This enriches the methodologies that can be used to discover knowledge from trajectories. Turning trajectories into graphs is one of the representative types of transformation. When conducting such a transformation, the main effort is to define what a node and an edge is in the transformed graph. The methods for

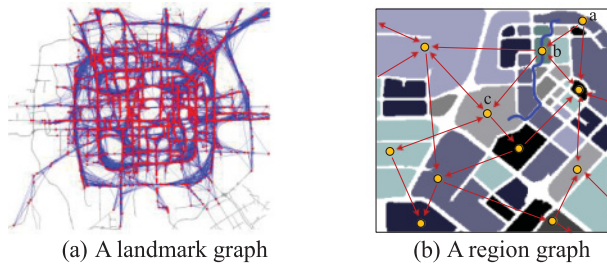


Fig. 19. Transforming trajectories into graphs.

transforming trajectories into a graph differentiate between one another, depending on whether a road network is involved in the transformation.

**9.1.1. In a Road Network Setting.** A road network is essentially a directed graph, where a node is an intersection and an edge denotes a road segment. Consequently, the most intuitive approach to turning trajectories into a graph is to project trajectories onto a road network. We can then calculate some weights, such as speed and traffic volume, for the edges based on the projected trajectories. Later, given the weighted graph, we can find the most likely route (traveled by people) between a few query points [Zheng et al. 2012a], identify the most popular route between a source and a destination [Luo et al. 2013], detect traffic anomalies [Pan et al. 2013], and update maps automatically.

The second approach is to build a *landmark graph*. For example, Yuan et al. [2013a, 2011a] proposed an intelligent driving direction system, entitled T-Drive, based on the GPS trajectories generated by a large number of taxicabs. After the map-matching process, T-Drive regards the top- $k$  road segments frequently traversed by taxicabs as *landmark* nodes (i.e., the red points shown in Figure 19(a)). The trajectories traversing two *landmarks* consecutively are aggregated into a *landmark edge* (denoted by a blue line), being used to estimate the travel time between two *landmarks*. A two-stage routing algorithm is proposed to find the fastest driving path. The algorithm first searches the *landmark graph* for a rough route (represented by a sequence of *landmarks*), and then finds a detailed route connecting consecutive *landmarks* on the original road network.

The third approach is to build a *region graph*, where a node denotes a region and an edge stands for the aggregation of commutes between the two regions. For instance, as illustrated in Figure 19(b), using an image segmentation-based algorithm, Yuan et al. [2012] and Zheng et al. [2011a] partition a city into regions by major roads so as to detect the underlying problems in a city's road network. A region bounded by major roads is then represented by a node, and two regions are connected with an edge if there are a certain number of commutes between them. After the transformation, they glean the region pairs (i.e., edges) that are not well connected, that is, with a huge traffic volume, a slow travel speed, and a long detour between them, using a skyline algorithm. The region graphs are also employed to detect traffic anomalies [Liu et al. 2011; Chawla et al. 2012] and urban functional regions [Yuan et al. 2012; Yuan et al. 2015].

**9.1.2. In Free Spaces.** Another branch of research transfers trajectories into a graph without using a road network, according to two major steps: (1) Identify key locations as vertexes from raw trajectories by using clustering methods. (2) Connect the vertexes to formulate a routable graph based on trajectories passing two locations.

**Travel Recommendation:** Zheng et al. [2009b] and Zheng and Xie [2011b] proposed to find the interesting locations and travel sequences from trajectories generated by many people. In the method, they first detect stay points from each trajectory and then cluster the stay points from different people into locations, as shown in Figure 20(a).

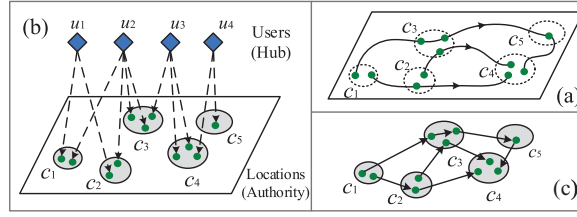


Fig. 20. Mining interesting locations and travel sequences.

Based on these locations and raw trajectories, they build a user-location bipartite graph as illustrated in Figure 20(b), as well as a routable graph between locations, as depicted in Figure 20(c).

In the bipartite graph, a user and a location are regarded as two different types of nodes. An edge is built between a user node and a location node if the user has visited the location. A HITS (Hypertext Induced Topic Search)-based model is then employed to infer the interest level of a location (i.e., the authority score) and the travel knowledge of a user (i.e., the hub score). According to the inferred scores, we can identify the top- $k$  most interesting locations and travel experts in a city. Bao et al. [2012] apply the similar idea in a CF framework to conduct the travel recommendation, concerned with a user's preferences, social environment, and current location.

In the location graph, as shown in Figure 20(c), an edge denotes the aggregation of raw trajectories traveling through two locations. To calculate the importance (or the representativeness) of an edge in this graph, three factors are considered: (1) the authority score of the source location (of the edge) weighted by the probability of people's moving out by this edge, (2) the authority score of a destination location (in the edge) weighted by the probability of people's moving in by this edge, and (3) the hub scores of the users who have traveled this edge. The score of a path is calculated by summing up the score of the edges the path contains.

Inspired by Zheng and Xie [2011b], a series of research was conducted to identify the popular routes from massive trajectories since 2010. Specifically, Yoon et al. [2012, 2011] suggest the best travel route, consisting of a sequence of locations with a typical stay time interval at each location, to a user, given the user's source and destination as well as the time period the user has. Chen et al. [2011] identify turning points from each raw trajectory, clustering these turning points into groups. These clusters are then used as vertexes to build a transfer network. Afterwards, the probability that people would travel from one vertex to another is calculated based on the counts of trajectories passing the two vertexes. Finally, given a source and a destination, the path with the maximum production of probabilities is found in the transfer network as the most popular route. However, the proposed method is not applicable to low-sampling-rate trajectories. To this end, Wei et al. [2012] divide a geographical space into uniform grids and then construct a routable graph based on the grids and raw trajectories. Refer to Section 5.1.2 for details.

Another branch of research is to detect the community of places based on the graph that is learned from trajectories, using some community discovery methods. A community of places is a cluster of locations with denser connections between locations in the cluster than between clusters. For example, Rinzivillo et al. [2012] aim to find the borders of human mobility at the lower spatial resolution of municipalities or counties. They mapped vehicle GPS tracks onto regions to formulate a complex network in Pisa. A community discovery algorithm, namely, Infomap, was then used

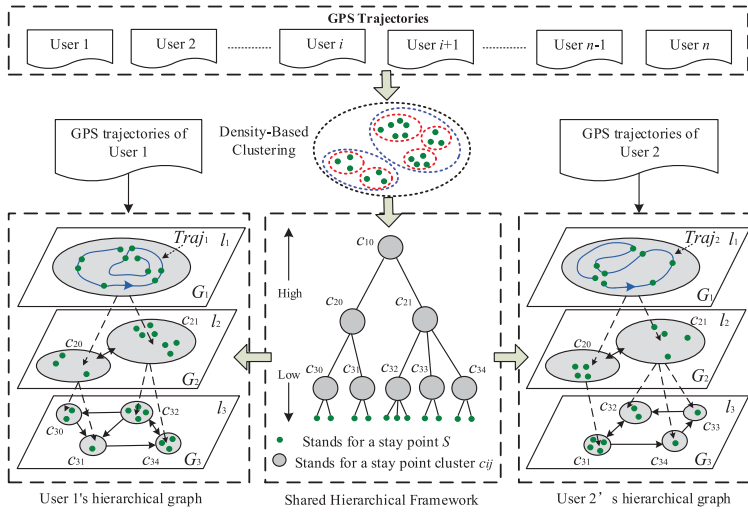


Fig. 21. Hierarchical graph-based user similarity estimation.

to partition the network into nonoverlapped subgraphs. More semantic meanings of a trajectory, such as a user's travel speed and experiences, have been considered in Liu et al. [2013] and Zheng et al. [2009c] to estimate the strength of interaction between two locations.

*Estimating User Similarity:* Another series of research transfers users' trajectories into hierarchical graphs so as to compute the similarity between different users. This is a foundation of many social applications, such as friend recommendation and community discovery.

As illustrated in Figure 21, Zheng et al. [2011c] deposit together the stay points detected from different users' trajectories, clustering them divisively by using a density-based clustering algorithm iteratively. As a result, a tree-based hierarchy is built, where a node on a higher level is a coarse-grained cluster (of stay points) and the nodes on a lower level are fine-grained clusters. The hierarchy is shared by different users as it is derived from all users' stay points. By projecting a user's trajectories onto this shared hierarchy, an individual hierarchical graph can be constructed for a user. As demonstrated in the bottom-left and bottom-right parts of Figure 21, two users' location histories are transformed from a collection of trajectories (which are not comparable between one another) to two individual graphs with common nodes. By matching the two graphs, common sequences of clusters are found on each level of the graphs. For example,  $c_{32} \rightarrow c_{31} \rightarrow c_{34}$  is a common sequence shared by the two users on the third level. Considering the popularity of a cluster in a common sequence, the length and the level (on the hierarchy) of the common sequences, a similarity score is calculated for a pair of users.

Xiao et al. [2014] extend the similarity computing from physical locations to a semantic space, aiming to facilitate the similarity estimation between users living in different cities or countries. A stay point detected from a trajectory is represented by the distribution of POIs (across different categories) within the scope of the stay point. The stay points from different users are then clustered into a hierarchy according to their distributions on different POI categories, in a similar way to that of Figure 21.

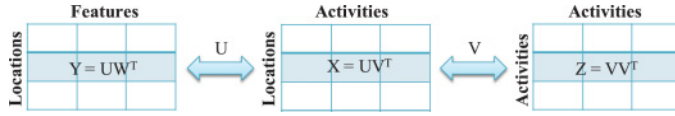


Fig. 22. Matrix factorization for recommendation.

## 9.2. From Trajectory to Matrix

Another form that we can transform trajectories into is a matrix. Using existing techniques, such as CF and MF, a matrix can help complement missing observations. A matrix can also be used as an input to identify anomalies. The key of the transformation lies in three aspects: (1) what does a row mean, (2) what is a column, and (3) what does an entry denote?

*Travel Recommendation.* Zheng et al. [2009b] and Zheng and Xie [2011b] transform users' GPS trajectory into a user-location matrix, where a row stands for a user and a column denotes a location (such as a cluster shown in Figure 21). The value of an entry means the number of visits of a user to a location. The matrix is very sparse, as a user can visit a very few locations. A CF model is then applied to the matrix to predict a user's interests in an unvisited location.

Zheng et al. [2010b] proposed a coupled-MF method to enable location-activity recommendation, using activity-tagged GPS trajectories. As illustrated in Figure 22, a location-activity matrix  $X$  is built, where a row stands for a venue (e.g., a cluster of GPS points) and a column represents a user-labeled activity (like shopping and dining). An entry in matrix  $X$  denotes the frequency of an activity that has been observed in users' labels in a particular location. Intuitively, this is a sparse matrix. A simple method to fill the missing entries is to decompose a matrix into the production of two low-rank matrices ( $U$  and  $V$ ) based on nonzero entries. After that, the missing entries can be filled by  $X = UV^T$ . Once this location-activity matrix is completely filled, given an activity, the top- $k$  locations, with a relatively high frequency from the column that corresponds to that activity, can be recommended. So does the activity recommendation for a location. To make a better recommendation, two context matrices, consisting of a location-feature matrix  $Y$  and an activity-activity matrix  $Z$ , are built based on additional data sources. The main idea is to propagate the information among  $X$ ,  $Y$ , and  $Z$  by requiring them to share low-rank matrices  $U$  and  $V$  in a collective MF model.

*Traffic Condition Estimation.* Shang et al. [2014] proposed a coupled-MF method to instantly estimate the travel speed on each road segment throughout an entire city, based on the GPS trajectory of a sample of vehicles (such as taxicabs). As shown in Figure 23(a), after map matching the GPS trajectories onto a road network, they formulate a matrix  $M'_r$  with a row denoting a time slot (e.g., 2pm–2:10pm) and a column standing for a road segment. Each entry in  $M'_r$  contains the travel speed on a particular road segment and in a particular time slot, calculated based on the recently received GPS trajectories. The goal is to fill the missing values in row  $t_j$ , which corresponds to the current time slot. Though we can achieve the goal by solely applying MF to  $M'_r$ , the accuracy of the inference is not very high as the majority of road segments are not covered by trajectories.

To address this issue, four context matrices ( $M_r$ ,  $M_G$ ,  $M'_G$ , and  $Z$ ) are built. Specifically,  $M_r$  stands for the historical traffic patterns on road segments. While the rows and columns of  $M_r$  have the same meaning as  $M'_r$ , an entry of  $M_r$  denotes the average travel speed derived from the historical data over a long period. The difference between the two corresponding entries from  $M'_r$  and  $M_r$  indicates the deviation of current traffic situation (on a road segment) from its average patterns. As depicted in Figure 23(b),  $Z$  contains the physical features of a road segment, such as the shape of a road,



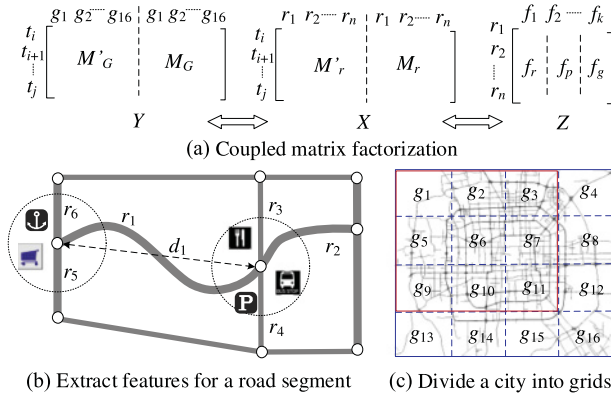


Fig. 23. Estimate traffic conditions based on trajectories.

number of lanes, speed constraint, and the distribution of surrounding POIs. The general assumption is that two road segments with similar geographical properties could have similar traffic conditions at the same time of day. To capture the high-level traffic conditions, as demonstrated in Figure 23(c), a city is divided into uniform grids. By projecting the recently received GPS trajectories into these grids, a matrix  $M'_G$  is built, with a column standing for a grid and a row denoting a time slot; an entry of  $M'_G$  means the number of vehicles traveling in a particular grid and at a particular time slot. Likewise, by projecting the historical trajectories over a long period into the grids, a similar  $M_G$  is built, where each entry means the average number of vehicles traveling in a particular grid and at a particular time slot. So,  $M'_G$  denotes the real-time high-level traffic conditions in a city and  $M_G$  indicates the historical high-level traffic patterns. The difference between the same entries of the two matrices suggests the deviation of current high-level traffic conditions from their historical averages. By combining these matrices, that is,  $X = M'_r || M_r$  and  $Y = M'_G || M_G$ , a coupled-MF is applied to  $X$ ,  $Y$ , and  $Z$ , with the objective function as follows:

$$L(T, R, G, F) = \frac{1}{2} \|Y - T(G; G)^T\|^2 + \frac{\lambda_1}{2} \|X - T(R; R)^T\|^2 + \frac{\lambda_2}{2} \|Z - RF^T\|^2 + \frac{\lambda_3}{2} (\|T\|^2 + \|R\|^2 + \|G\|^2 + \|F\|^2),$$

where  $\|\cdot\|$  denotes the Frobenius norm. The first three terms in the objective function control the loss in MF, and the last term is a regularization of penalty to prevent overfitting.

**Diagnosing Traffic Anomalies.** Chawla et al. [2012] aim to identify the traffic flows that cause an anomaly between two regions. In the methodology, they first partition a city into regions by major roads, building a region graph based on trajectories of taxicabs, as illustrated in Figure 24(a). A trajectory is then represented by a path on the graph, that is, a sequence of links between regions, as shown in Figure 24(b). Two matrices are built based on the trajectories and graph. One is a link-traffic matrix  $L$ , as shown in Figure 24(c), where a row is a link and a column corresponds to a time interval. An entry of  $L$  denotes the number of vehicles traversing a particular link at a specific time interval. The other is a link-path matrix  $A$ , with a row standing for a link and column denoting a path. An entry of  $A$  is set to 1 if a particular link is contained in a particular path. Given matrix  $L$ , they first use a Principal Component Analysis (PCA) algorithm to detect some anomalous links, which were represented by a column vector  $b$  with 1 denoting an anomaly detected on the link. Then, the

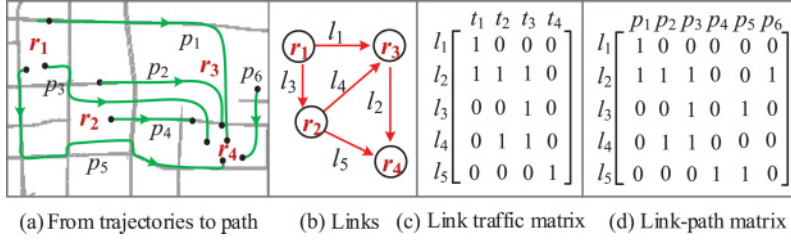


Fig. 24. From trajectories to matrices for detecting anomalies.

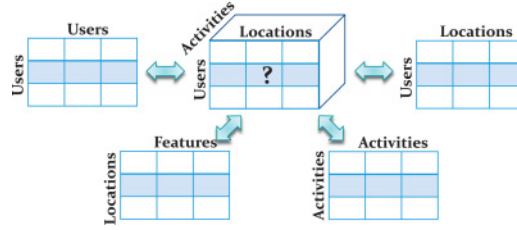


Fig. 25. Recommendation based on trajectories and tensors.

relationship between anomalous links and paths was captured by solving the equation  $Ax = b$ , where  $x$  is a column vector denoting which paths contribute to the emergency of these anomalies shown in  $b$ . Using  $L_1$  optimization techniques,  $x$  can be inferred.

### 9.3. From Trajectory to Tensor

A nature extension of the matrix-based transformation is turning trajectories into a (3D) tensor, where the third dimension is added to a matrix so as to accommodate additional information. The goal of the transformation is usually to fill the missing entries (in a tensor) or find the correlation between two objects, like two road segments or gas stations. A common approach to solving this problem is to decompose a tensor into the multiplication of a few (low-rank) matrices and a core tensor (or just a few vectors), based on the tensor's nonzero entries. When a tensor is very sparse, in order to achieve a better performance, the tensor is usually decomposed with other (context) matrices in a framework of CF.

Zheng et al. [2012b, 2010a] extend the generic location-activity research [Zheng et al. 2010b] into a personalized one, by adding a user dimension into the original location-activity matrix. As shown in Figure 25, a user-location-activity tensor  $A$  is built, with an entry denoting the times that a particular user has performed a particular activity in a particular location. If we can infer the value of every entry, personalized recommendation can be enabled. However, tensor  $A$  is very sparse as a user usually visits a few places. Thus, a simple tensor completion method cannot fill its missing entries very well. To address this issue, four context matrices are built based on additional data sources, such as road network and POI datasets, which are not sparse. In addition, these matrices share some dimension with tensor  $A$ . For instance, tensor  $A$  shares the *user* dimension with matrix  $B$  and the *location* dimension with matrix  $E$ . Consequently, the knowledge from these matrices can be transferred into the tensor to help in completing tensor  $A$ .

Wang et al. [2014] proposed a coupled tensor-decomposition-based method to instantly estimate the travel time of a path, based on a sample of vehicles' GPS trajectories. To model the traffic conditions of the current time slot, they construct a tensor  $\mathcal{A}_t \in \mathbb{R}^{N \times M \times L}$ , with the three dimensions standing for road segments, drivers, and time

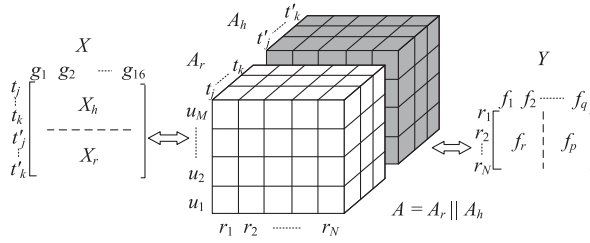


Fig. 26. Travel time estimation using tensor decomposition.

slots, respectively, based on the GPS trajectories received in the most recent  $L$  time slots and the road network data. As shown in Figure 26, an entry  $A_r(i, j, k) = c$  denotes the  $i$ th road segment is traveled by the  $j$ th driver with a time cost  $c$  in time slot  $k$  (e.g., 2–2:30pm). The last time slot denotes the present time slot, combined with the  $L-1$  time slots right before it formulates the tensor. Clearly, the tensor is very sparse as a driver can only travel a few road segments in a short time period. If the missing entries can be inferred based on the values of nonzero entries, we can obtain the travel time of any driver on any road segment in the present time slot.

To this end, another tensor  $A_h$  is built based on the historical trajectories over a long period of time (e.g., 1 month).  $A_h$  has the same structure as  $A_r$ , while an entry  $A_h(i, j, k) = c'$  denotes the  $j$ th driver's average travel time on the  $i$ th road segment in time slot  $k$  in the history. Intrinsically,  $A_h$  is much denser than  $A_r$ , denoting the historical traffic patterns and drivers' behavior on an entire road network. Besides, two context matrices ( $X$  and  $Y$ ) are built to help supplement the missing entries of  $A_r$ . Matrix  $X$  (consisting of  $X_r$  and  $X_h$ ) represents the correlation between different time slots in terms of the coarse-grained traffic conditions. This is similar to its correspondence shown in Figure 23(c). An entry of  $X_r$  denotes the number of vehicles traversing a particular grid in a particular time slot. A row of  $X_r$  represents coarse-grained traffic conditions in a city at a particular time slot. Consequently, the similarity of two different rows indicates the correlation of traffic flows between two time slots.  $X_h$  has the same structure as  $X_r$ , storing the historical average number of vehicles traversing a grid from  $t_i$  to  $t_j$ . Matrix  $Y$  stores each road segment's geographical features, which are similar to that of Matrix  $Z$  shown in Figure 23(a). Later, they decompose  $A = A_r || A_h$  with matrices  $X$  and  $Y$  collaboratively, by optimizing the following objective function:

$$\begin{aligned} \mathcal{L}(S, R, U, T, F, G) = & \frac{1}{2}A - S \times_R R \times_U U \times_T T^2 + \frac{\lambda_1}{2}X - TG^2 \\ & + \frac{\lambda_2}{2}Y - RF^2 + \frac{\lambda_3}{2}(S^2 + R^2 + U^2 + T^2 + F^2 + G^2). \end{aligned}$$

A similar idea was employed by Zhang et al. [2013, 2015] to estimate the queuing time in each gas station throughout a city. The queuing time is further used to estimate the number of vehicles that are being refueled. Specifically, the refueling events are first detected from a taxicab's GPS trajectories based on a stay point-based inference method. Then, as shown in Figure 27, a three dimension tensor  $F$  is built, with the first dimension denoting gas stations, the second one standing for time of day, and the third for the day of the week. An entry means the average waiting time (detected from taxi trajectories) at a station in a particular day of the week and at a particular time interval. This tensor is intrinsically sparse as we cannot guarantee to have a taxicab being refueled in each station anytime. A context matrix is built, incorporating the geographical features of a station. Intuitively, two gas stations with the similar

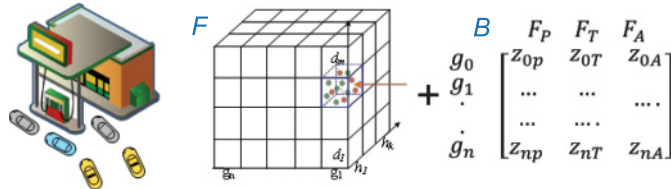


Fig. 27. Estimate the refueling behavior in a gas station.

surrounding environment (including road networks and POIs) and traffic flow could have the similar refueling pattern. The coupled TD method mentioned in previous examples is then applied to the tensor and matrix, filling the missing value in  $F$ .

## 10. MISCELLANEOUS

### 10.1. Public Trajectory Datasets

Collecting data is always the first priority of trajectory data mining. Thanks to researchers in this field, there are quite a few real trajectory datasets that are publicly available:

- GeoLife Trajectory Dataset* [GeoLife Data]: a GPS trajectory dataset from Microsoft Research GeoLife project [Zheng et al. 2010d], collected by 182 users from April 2007 to August 2012. The dataset has been used to estimate the similarity between users [Li et al. 2008], which enables friend and location recommendations [Zheng and Xie 2011b; Zheng et al. 2009c]. It was also used by Chen et al. [2010] for studying the problem of finding the nearest trajectory to a sequence of query points.
- T-Drive Taxi Trajectories* [T-Drive Data]: A sample of trajectories from Microsoft Research T-Drive project [Yuan et al. 2010a], generated by over 10,000 taxicabs in a week of 2008 in Beijing. The full dataset was used to suggest the practically fastest driving directions to normal drivers [Yuan et al. 2011a, 2013a], recommend passenger-pickup location for taxi drivers [Yuan et al. 2011a; Yuan et al. 2015], enable dynamic taxi ride-sharing [Ma et al. 2013; Ma et al. 2015], glean the problematic design in a city's transportation network [Zheng et al. 2011a], and identify urban functional regions [Yuan et al. 2012].
- GPS Trajectory with Transportation Labels* [Trajectory with transportation modes]: Each trajectory has a set of transportation mode labels, such as driving, taking a bus, riding a bike, and walking. The dataset can be used to evaluate trajectory classification and activity recognition [Zheng et al. 2008a, 2008b].
- Check-in Data from Location-based Social Networks* [User check-in data]: The dataset consists of the check-in data generated by over 49,000 users in New York City and 31,000 users in Los Angeles as well as the social structure of the users. Each check-in includes a venue ID, the category of the venue, a time stamp, and a user ID. As the check-in data of a user can be regarded as a low-sampling-rate trajectory, this dataset has been used to study the uncertainty of trajectories [Wei et al. 2012] and evaluate location recommendation [Bao et al. 2012].
- Hurricane Trajectories* [Hurricane trajectory (HURDAT)]: This dataset is provided by the National Hurricane Service (NHS), containing 1,740 trajectories of Atlantic Hurricanes (formally defined as tropical cyclone) from 1851 to 2012. NHS also provides annotations of typical hurricane tracks for each month throughout the annual hurricane season that spans from June to November. The dataset can be used to test trajectory clustering and uncertainty.

- The Greek Truck Trajectories* [The Greek Trucks Dataset]: This dataset contains 1,100 trajectories from 50 different trucks delivering concrete around Athens, Greece. It was used to evaluate trajectory pattern mining task in Giannotti et al. [2007].
- Movebank Animal Tracking Data* [Movebank data]: Movebank is a free, online database of animal tracking data, helping animal tracking researchers to manage, share, protect, analyze, and archive their data.

When needing massive trajectories to test the efficiency of a method, we can generate synthetic trajectories based on traffic generators, for example, BerlinMod [Duntgen et al. 2009] and Thomas-Brinkhoff [Brinkhoff and Str 2002]. There is also a web-based interface, called MNTG [Mokbel et al. 2014], which supports the two traffic generators to work on any arbitrary road networks. Ma et al. [2015] build a taxi ride request simulator based on the pickup and drop-off points of the real taxi trajectories generated in Beijing. The simulator is used to test the efficiency of a taxi ride-sharing service.

## 10.2. Conferences and Journals Concerning Trajectories

Research about trajectory data mining has a wide presence at the following venues:

- General data mining conferences: KDD, ICDM, SDM, PAKDD, and ICML-PKDD.
- General database conferences: ICDE, VLDB, SIGMOD, EDBT, and DASFAA.
- General artificial intelligence conferences: IJCAI and AAAI.
- Spatial-data-focusing conferences: ACM SIGSPATIAL GIS, SSTD, and MDM.
- Application-driven conferences and workshops: International Conference on Ubiquitous Computing, and the International Workshop on Urban Computing [Zheng et al. 2013].
- Journals and Transactions: IEEE TKDE, ACM TKDD, ACM TIST, VLDB, Data Mining and Knowledge Discovery, KAIS, DKE, and Journal on Personal and Ubiquitous Computing. Besides the journals in the computer science area, there are many journals in other disciplinaries, such as Transportation Research Part C, IEEE Transaction on Intelligent Transportation Systems, and Transportation Record.

## 10.3. Potential Future Direction

In the big data era, a data mining task needs to harness a diversity of data. This is calling for new technology that can unlock the power of knowledge from multiple data sources. Under such a circumstance, how to mine trajectory data together with other data sources is a new challenge. There two approaches towards this goal.

One is to combine trajectories with other data sources to fulfill a data mining task. For example, Zheng et al. [2014c, 2013] infer the fine-grained air quality, using trajectories of vehicles, POIs, and meteorological data. Fu et al. [2014a, 2014b] combine human mobility data represented by trajectories with social media and geographical data to rank the potential value of real estates. Yuan et al. [2012, 2015] aim to identify the functional regions in a city based on taxi trajectories, road network data, and POIs. Zheng et al. [2014c] diagnoses the urban noise, using check-in data, traffic, and 311 complaints. The other approach is to use other sources to enrich a trajectory. For instance, leveraging POIs and road network data, Wang et al. [2014] better estimates the travel time of a path based on sparse trajectories.

The new challenge calls for (1) data management techniques that can organize multi-modal data for an efficient retrieval and mining, (2) the cross-domain machine learning methods that can unlock the power of knowledge that cannot be discovered from a single data source, and (3) advanced visualization techniques that can suggest the insights across different sources.



## 11. CONCLUSION

The wide availability of trajectory data has fostered a diversity of applications, calling for algorithms that can discover knowledge from the data effectively and efficiently. This article surveys the techniques concerned with different stages of trajectory data mining, recapping them by categories and exploring the differences between one another. This article also suggests the approaches of transforming raw trajectories into other data structures, to which more existing data mining techniques can be applied. This article provides an overview on how to unlock the power of knowledge from trajectories, for researchers and professionals from not only computer sciences but also a broader range of communities dealing with trajectories. At the end of this article, a list of public trajectory datasets has been given and a few future directions have been suggested.

## ACKNOWLEDGMENTS

A small portion of the content of this article is derived from a few chapters of a book, entitled *Computing with Spatial Trajectories* [Zheng and Zhou 2011], which was co-edited by Xiaofang Zhou and myself. I appreciate the chapter authors of this book: C.-Y. Chow, K. Deng, C. S. Jensen, H. Jeung, J. Krumm, W.-C. Lee, M. F. Mokbel, K. Xie, K. Zheng, G. Trajcevski, Q. Yang, M. L. Yiu, V. W. Zheng, and Y. Zhu.

## REFERENCES

- O. Abul, F. Bonchi, and M. Nanni. 2008. Never walk alone: Uncertainty for anonymity in moving objects databases. In *Proceedings of the 24th IEEE International Conference on Data Engineering*. IEEE, 376–385.
- C. C. Aggarwal, J. Han, J. Wang, and P. S. Yu. 2003. A framework for clustering evolving data streams. In *Proceedings of the 29th International Conference on Very Large Data Bases*. VLDB Endowment 29, 81–92.
- R. Agrawal, C. Faloutsos, and A. Swami. 1993. *Efficient similarity search in sequence databases*. Springer, 69–84.
- H. Alt, A. Efrat, G. Rote, and C. Wenk. 2003. Matching planar maps. *Journal of Algorithms* 49, 2 (2003), 262–283.
- J. Bao, Y. Zheng, and M. F. Mokbel. 2012. Location-based and preference-aware recommendation using sparse geo-social networking data. In *Proceedings of the 20th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems*. ACM, 199–208.
- J. Bao, Y. Zheng, D. Wilkie, and M. F. Mokbel. 2015. A survey on recommendations in location-based social networks. *GeoInformatica*, 19, 3, 525–565.
- R. Bellman. 1961. On the approximation of curves by line segments using dynamic programming. *Communications of the ACM* 4, 6 (1961), 284.
- A. R. Beresford and F. Stajano. 2003. Location privacy in pervasive computing. *IEEE Pervasive Computing* 2, 1 (2003), 46–55.
- S. Brakatsouls, D. Pfoser, R. Salas, and C. Wenk. 2005. On map-matching vehicle tracking data. In *Proceedings of the 31st International Conference on Very Large Data Bases*. VLDB Endowment, 853–864.
- T. Brinkhoff and O. Str. 2002. A framework for generating network-based moving objects. *Geoinformatica*, 6, 2 (2002), 153–180.
- H. Cao, N. Mamoulis, and D. W. Cheung. 2005. Mining frequent spatio-temporal sequential patterns. In *Proceedings of the 5th IEEE International Conference on Data Mining*. IEEE, 82–89.
- H. Cao, N. Mamoulis, and D. W. Cheung. 2007. Discovery of periodic patterns in spatiotemporal sequences. *IEEE Transactions on Knowledge and Data Engineering* 19, 4 (2007), 453–467.
- I. V. Cadez, S. Gaffney, and P. Smyth. 2000. A general probabilistic framework for clustering individuals and objects. In *Proceedings of the 6th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. ACM, 140–149.
- V. Chandola, A. Banerjee, and V. Kumar. 2009. Anomaly detection: A survey. *ACM Computing Surveys* 41, 3 (2009), 1–58.
- S. Chawla, Y. Zheng, and J. Hu. 2012. Inferring the root cause in road traffic anomalies. In *Proceedings of the 12th IEEE International Conference on Data Mining*. IEEE, 141–150.

- S. S. Chawathe. 2007. Segment-based map matching. *2007 IEEE Intelligent Vehicles Symposium*. IEEE, 1190–1197.
- Y. Chen, K. Jiang, Y. Zheng, C. Li, and N. Yu. 2009. Trajectory simplification method for location-based social networking services. In *Proceedings of the ACM SIGSPATIAL Workshop on Location-Based Social Networking Services*. ACM, 33–40.
- L. Chen and R. Ng. 2004. On the marriage of lp-norms and edit distance. In *Proceedings of the 30th International Conference on Very Large Data Bases*. VLDB Endowment, 792–803.
- L. Chen, M. T. Ozsu, and V. Oria. 2005. Robust and fast similarity search for moving object trajectories. In *Proceedings of the 24th ACM SIGMOD International Conference on Management of Data*. ACM, 491–502.
- Z. Chen, H. T. Shen, and X. Zhou. 2011. Discovering popular routes from trajectories. In *Proceedings of the 27th IEEE International Conference on Data Engineering*. IEEE, 900–911.
- Z. Chen, H. T. Shen, X. Zhou, Y. Zheng, and X. Xie. 2010. Searching trajectories by locations—An efficient study. In *Proceedings of the 29th ACM SIGMOD International Conference on Management of Data*. ACM, 255–266.
- W. Chen, M. Yu, Z. Li, and Y. Chen. 2003. Integrated vehicle navigation system for urban applications. In *Proceedings of the International Conference Global Navigation Satellite System*. CGNS, 15–22.
- R. Cheng, J. Chen, M. F. Mokbel, and C. Y. Chow. 2008. Probabilistic verifiers: Evaluating constrained nearest-neighbor queries over uncertain data. In *Proceedings of the IEEE 24th Conference on Data Engineering*. IEEE, 973–982.
- R. Cheng, D. V. Kalashnikov, and S. Prabhakar. 2004. Querying imprecise data in moving objects environments. *IEEE Transactions on Knowledge and Data Engineering* 16, 9 (2004).
- C. Y. Chow and M. F. Mokbel. 2011. Privacy of spatial trajectories. *Computing with Spatial Trajectories*, Y. Zheng and X. Zhou (Eds.). Springer, 109–141.
- A. Civilis, C. S. Jensen, J. Nenortaitė, and S. Pakalnis. 2005. Techniques for efficient road-network-based tracking of moving objects. *IEEE Transactions on Knowledge and Data Engineering* 17, 5 (2005), 698–711.
- K. Deng, K. Xie, K. Zheng, and X. Zhou. 2011. Trajectory indexing and retrieval. *Computing with Spatial Trajectories*. Y. Zheng and X. Zhou (Eds.). Springer, 35–60.
- D. Douglas and T. Peucker. 1973. Algorithms for the reduction of the number of points required to represent a line or its caricature. *Cartographica: The International Journal for Geographic Information and Geovisualization* 10, 2 (1973), 112–122.
- C. Duntgen, T. Behr, and R. H. Guting. 2009. BerlinMOD: A benchmark for moving object databases. *The VLDB Journal* 18, 6 (2009), 1335–1368.
- T. Emrich, H. P. Kriegel, N. Mamoulis, M. Renz, and A. Züfle. 2012. Querying uncertain spatio-temporal data. In *Proceedings of the 28th IEEE International Conference on Data Engineering*. IEEE, 354–365.
- Y. Fu, Y. Ge, Y. Zheng, Z. Yao, Y. Liu, H. Xiong, and N. J. Yuan. 2014a. Sparse real estate ranking with online user reviews and offline moving behaviors. In *Proceedings of the 14th IEEE International Conference on Data Mining*. IEEE, 120–129.
- Y. Fu, H. Xiong, Y. Ge, Z. Yao, and Y. Zheng. 2014b. Exploiting geographic dependencies for real estate appraisal: A mutual perspective of ranking and clustering. In *Proceedings of the 20th SIGKDD Conference on Knowledge Discovery and Data Mining*. ACM, 1047–1056.
- S. Gaffney and P. Smyth. 1999. Trajectory clustering with mixtures of regression models. In *Proceedings of the 5th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 63–67.
- F. Giannotti, M. Nanni, D. Pedreschi, and F. Pinelli. 2007. Trajectory pattern mining. In *Proceedings of the 13th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 330–339.
- G. Gid'ofalvi, X. Huang, and T. B. Pedersen. 2007. Privacy-preserving data mining on moving object trajectories. In *Proceedings of the 8th IEEE International Conference on Mobile Data Management*. IEEE, 60–68.
- J. S. Greenfeld. 2002. Matching GPS observations to locations on a digital map. In *Proceedings of the 81st Annual Meeting of the Transportation Research Board*. 576–582.
- J. Gudmundsson and M. V. Kreveld. 2006. Computing longest duration flocks in trajectory data. In *Proceedings of the 14th Annual ACM International Symposium on Advances in Geographic Information Systems*. ACM, 35–42.
- J. Gudmundsson, M. V. Kreveld, and B. Speckmann. 2004. Efficient detection of motion patterns in spatio-temporal data sets. In *Proceedings of the 12th Annual ACM International Symposium on Advances in Geographic Information Systems*. ACM, 250–257.
- J. Hersherberger and J. Snoeyink. 1992. Speeding up the Douglas-Peucker line simplification algorithm. In *Proceedings of the International Symposium on Spatial Data Handling*. 134–143.

- B. Hoh, M. Gruteser, H. Xiong, and A. Alrabady. 2010. Achieving guaranteed anonymity in GPS traces via uncertainty-aware path cloaking. *IEEE Transactions on Mobile Computing* 9, 8 (2010), 1089–1107.
- C. S. Jensen, D. Lin, and B. C. Ooi. 2007. Continuous clustering of moving objects. *IEEE Transaction on Knowledge and Data Engineering* 19, 9 (2007), 1161–1174.
- H. Jeung, H. Shen, and X. Zhou. 2008a. Convoy queries in spatio-temporal databases. In *Proceedings of the 24th IEEE International Conference on Data Engineering*. IEEE, 1457–1459.
- H. Jeung, M. L. Yiu, and C. S. Jensen. 2011. Trajectory pattern mining. *Computing with Spatial Trajectories*. Y. Zheng and X. Zhou (Eds.). Springer, 143–177.
- H. Jeung, M. Yiu, X. Zhou, C. Jensen, and H. Shen. 2008b. Discovery of convoys in trajectory databases. *Proceedings of the VLDB Endowment* 1, 1 (2008), 1068–1080.
- G. Kellaris, N. Pelekis, and Y. Theodoridis. 2009. Trajectory compression under network constraints. In *Proceedings of the International Symposium on Advances in Spatial and Temporal Databases*. 392–398.
- E. J. Keogh, S. Chu, D. Hart, and M. J. Pazzani. 2001. An on-line algorithm for segmenting time series. In *Proceedings of the IEEE International Conference on Data Engineering*. IEEE, 289–296.
- A. Kharrat, I. S. Popa, K. Zeitouni, and S. Faiz. 2008. Clustering algorithm for network constraint trajectories. *Headway in Spatial Data Handling*. 631–647.
- H. Kido, Y. Yanagisawa, and T. Satoh. 2005. An anonymous communication technique using dummies for location-based services. In *Proceedings of the 3rd International Conference on Pervasive Services*. IEEE, 88–97.
- J. Krumm. 2011. Trajectory analysis for driving. *Computing with Spatial Trajectories*, Y. Zheng and X. Zhou (Eds.). Springer, 213–241.
- J. Krumm and E. Horvitz. 2004. LOCADIO: Inferring motion and location from Wi-Fi signal strengths. In *Proceedings of the International Conference on Mobile and Ubiquitous Systems*. IEEE, 4–13.
- J. G. Lee, J. Han, and K. Y. Whang. 2007. Trajectory clustering: A partition-and-group framework. In *Proceedings of the ACM SIGMOD Conference on Management of Data*. ACM, 593–604.
- J. Lee, J. Han, and X. Li. 2008. Trajectory outlier detection: A partition-and-detect framework. In *Proceedings of the 24th IEEE International Conference on Data Engineering*. IEEE, 140–149.
- W.-C. Lee and J. Krumm. 2011. Trajectory preprocessing. *Computing with Spatial Trajectories*, Y. Zheng and X. Zhou (Eds.). Springer, 1–31.
- Q. Li, Y. Zheng, X. Xie, Y. Chen, W. Liu, and M. Ma. 2008. Mining user similarity based on location history. In *Proceedings of the 16th Annual ACM International Conference on Advances in Geographic Information Systems*. ACM, 34.
- Z. Li, B. Ding, J. Han, and R. Kays. 2010a. Swarm: Mining relaxed temporal moving object clusters. *Proceedings of the VLDB Endowment* 3, 1–2 (2010), 723–734.
- Z. Li, J. Lee, X. Li, and J. Han. 2010b. Incremental clustering for trajectories. *Database Systems for Advanced Applications*. 32–46.
- Z. Li, B. Ding, J. Han, R. Kays, and P. Nye. 2010c. Mining periodic behaviors for moving objects. In *Proceedings of the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 1099–1108.
- Z. Li, J. Wang, and J. Han. 2012. Mining event periodicity from incomplete observations. In *Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 444–452.
- L. Liao, D. Fox, and H. Kautz. 2004. Learning and inferring transportation routines. In *Proceedings of the National Conference on Artificial Intelligence*. 348–353.
- S. Liu, K. Jayarajah, A. Misra, and R. Krishnan. 2013. TODMIS: Mining communities from trajectories. In *Proceedings of the 22nd ACM CIKM International Conference on Information and Knowledge Management*. ACM, 2109–2118.
- S. Liu, L. Ni, and R. Krishnan. 2014. Fraud detection from Taxis' driving behaviors. *IEEE Transactions on Vehicular Technology* 63, 1 (2014), 464–472.
- W. Liu, Y. Zheng, S. Chawla, J. Yuan, and X. Xie. 2011. Discovering spatio-temporal causal interactions in traffic data streams. In *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 1010–1018.
- Y. Lou, C. Zhang, Y. Zheng, X. Xie, W. Wang, and Y. Huang. 2009. Map-matching for low-sampling-rate GPS trajectories. In *Proceedings of the 17th ACM SIGSPATIAL International Conference on Geographical Information Systems*. ACM, 352–361.
- W. Luo, H. Tan, L. Chen, and M. N. Lionel. 2013. Finding time period-based most frequent path in big trajectory data. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*. ACM, 713–724.

- S. Ma, Y. Zheng, and O. Wolfson. 2013. T-Share: A large-scale dynamic taxi ridesharing service. In *Proceedings of the 29th IEEE International Conference on Data Engineering*. IEEE, 410–421.
- S. Ma, Y. Zheng, and O. Wolfson. 2015. Real-time city-scale taxi ridesharing. *IEEE Transactions on Knowledge and Data Engineering* 99. DOI: <http://doi.ieeecomputersociety.org/10.1109/TKDE.2014.2334313>
- N. Maratnia and R. A. de By. 2004. Spatio-temporal compression techniques for moving point objects. In *Proceedings of the 9th International Conference on Extending Database Technology*. 765–782.
- R. B. McMaster. 1986. A statistical analysis of mathematical measures of linear simplification. *The American Cartographer* 13, 2 (1986), 103–116.
- M. F. Mokbel, C. Y. Chow, and W. G. Aref. 2007. The new Casper: Query processing for location services without compromising privacy. In *Proceedings of the 23rd IEEE International Conference on Data Engineering*. IEEE, 1499–1500.
- M. Mokbel, L. Alarabi, J. Bao, A. Eldawy, A. Magdy, M. Sarwat, E. Waytas, and S. Yackel. 2014. A demonstration of MNTG —A Web-based road network traffic generator. In *Proceedings of the 30th IEEE International Conference on Data Engineering*, IEEE, 1246–1249.
- A. Monreale, F. Pinelli, R. Trasarti, and F. Giannotti. 2009. WhereNext: A location predictor on trajectory pattern mining. In *Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 637–646.
- M. E. Nergiz, M. Atzori, Y. Saygin, and B. Guc. 2009. Towards trajectory anonymization: A generalization-based approach. *Transactions on Data Privacy* 2, 1 (2009), 47–75.
- P. Newson, J. Krumm. 2009. Hidden Markov map matching through noise and sparseness. In *Proceedings of the 17th ACM SIGSPATIAL International Conference on Geographical Information Systems*. ACM, 336–343.
- J. Niedermayer, A. Zuffe, T. Emrich, M. Renz, N. Mamouliso, L. Chen, and H. Kriegel. 2014. Probabilistic nearest neighbor queries on uncertain moving object trajectories. *Proceedings of the VLDB Endowment* 7, 3 (2014), 205–216.
- W. Y. Ochieng, M. A. Quddus, and R. B. Noland. 2004. Map-matching in complex urban road networks. *Brazilian Journal of Cartography* 55, 2 (2004), 1–18.
- B. Pan, Y. Zheng, D. Wilkie, and C. Shahabi. 2013. Crowd sensing of traffic anomalies based on human mobility and social media. In *Proceedings of the 21st Annual ACM International Conference on Advances in Geographic Information Systems*. ACM, 334–343.
- L. X. Pang, S. Chawla, W. Liu, and Y. Zheng. 2011. On mining anomalous patterns in road traffic streams. In *Proceedings of the International Conference on Advanced Data Mining and Applications*. 237–251.
- L. X. Pang, S. Chawla, W. Liu, and Y. Zheng. 2013. On detection of emerging anomalous traffic patterns using GPS data. *Data & Knowledge Engineering*, 87 (2013), 357–373.
- D. J. Patterson, L. Liao, D. Fox, and H. Kaut. 2003. Inferring high-level behavior from low-level sensors. In *Proceedings of the 5th International Conference on Ubiquitous Computing*. ACM, 73–89.
- J. Pei, J. Han, B. Mortazavi-Asl, and H. Pinto. 2011. PrefixSpan: Mining sequential patterns efficiently by prefix-projected pattern growth. In *Proceedings of the 29th IEEE International Conference on Data Engineering*. IEEE, 215.
- D. Pfoser and C. S. Jensen. 1999. Capturing the uncertainty of moving objects representation. In *Proceedings of the International Symposium on Advances in Spatial Databases*. 111–131.
- D. Pfoser, C. S. Jensen, and Y. Theodoridis. 2000. Novel approaches to the indexing of moving object trajectories. In *Proceedings of the 26th International Conference on Very Large Data Bases*. VLDB Endowment, 395–406.
- O. Pink and B. Hummel. 2008. A statistical approach to map matching using road network geometry, topology and vehicular motion constraints. In *Proceedings of the 11th International IEEE Conference on Intelligent Transportation Systems*. IEEE, 862–867.
- M. Potamias, K. Patroumpas, and T. Sellis. 2006. Sampling trajectory streams with spatio-temporal criteria. In *Proceedings of the 18th International Conference on Scientific and Statistical Database Management*. IEEE, 275–284.
- S. Qiao, C. Tang, H. Jin, T. Long, S. Dai, Y. Ku, and M. Chau. 2010. Putmode: Prediction of uncertain trajectories in moving objects databases. *Applied Intelligence* 33, 3 (2010), 370–386.
- M. A. Quddus, W. Y. Ochieng, and R. B. Noland. 2006. A high accuracy fuzzy logic-based map-matching algorithm for road transport. *Journal of Intelligent Transportation Systems: Technology, Planning, and Operations* 10, 3 (2006), 103–115.
- K. Richter, F. Schmid, and P. Laube. 2012. Semantic trajectory compression: Representing urban movement in a nutshell. *Journal of Spatial Information Science*, 4 (2012), 3–30.



- S. Rinzivillo, S. Mainardi, F. Pezzoni, M. Coscia, D. Pedreschi, and F. Giannotti. 2012. Discovering the geographical borders of human mobility. *Künstl. Intell.* 26, 3 (2012), 253–260.
- J. Shang, Y. Zheng, W. Tong, E. Chang, and Y. Yu. 2014. Inferring gas consumption and pollution emission of vehicles throughout a city. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 1027–1036.
- R. Song, W. Sun, B. Zheng, and Y. Zheng. 2014. PRESS: A novel framework of trajectory compression in road networks. *Proceedings of the VLDB Endowment* 7, 9 (2014), 661–672.
- H. Su, K. Zheng, H. Wang, J. Huang, and X. Zhou. 2013. Calibrating trajectory data for similarity-based analysis. In *Proceedings of the 39th International Conference on Very Large Data Bases*. VLDB Endowment, 833–844.
- Y. Tao and D. Papadias. 2001a. Efficient historical R-trees. In *Proceedings of the 13th International Conference on Scientific and Statistical Database Management*, 223–232.
- Y. Tao and D. Papadias. 2001b. Mv3r-tree: A spatio-temporal access method for timestamp and interval queries. In *Proceedings of the 27th International Conference on Very Large Data Bases*. VLDB Endowment, 431–440.
- Y. Tao, D. Papadias, and Q. Shen. 2002. Continuous nearest neighbour search. In *Proceedings of the 28th International Conference on Very Large Data Bases*. VLDB Endowment, 287–298.
- L. A. Tang, Y. Zheng, X. Xie, J. Yuan, X. Yu, and J. Han. 2011. Retrieving k-nearest neighboring trajectories by a set of point locations. In *Proceedings of the 12th Symposium on Spatial and Temporal Databases*. Springer, 223–241.
- L. A. Tang, Y. Zheng, J. Yuan, J. Han, A. Leung, C. Hung, and W. Peng. 2012a. Discovery of traveling companions from streaming trajectories. In *Proceedings of the 28th IEEE International Conference on Data Engineering*. IEEE, 186–197.
- L. A. Tang, Y. Zheng, J. Yuan, J. Han, A. Leung, W. Peng, and T. L. Porta. 2012b. A framework of traveling companion discovery on trajectory data streams. *ACM Transactions on Intelligent Systems and Technology* 5, 1 (2012).
- M. Terrovitis and N. Mamoulis. 2008. Privacy preservation in the publication of trajectories. In *Proceedings of the 9th IEEE International Conference on Mobile Data Management*. IEEE, 65–72.
- G. Trajcevski, A. N. Choudhary, O. Wolfson, L. Ye, and G. Li. 2010. Uncertain range queries for necklaces. In *Proceedings of the 11th IEEE International Conference on Mobile Data Management*. IEEE, 199–208.
- G. Trajcevski, R. Tamassia, H. Ding, P. Scheuermann, and I. F. Cruz. 2009. Continuous probabilistic nearest-neighbor queries for uncertain trajectories. In *Proceedings of the 12th International Conference on Extending Database Technology: Advances in Database Technology*. ACM, 874–885.
- G. Trajcevski, O. Wolfson, K. Hinrichs, and S. Chamberlain. 2004. Managing uncertainty in moving objects databases. *ACM Transactions on Database Systems* 29, 3(04), 463–507.
- S. Timothy, A. Varshavsky, A. Lamarca, M. Y. Chen, and T. Chounhury. 2006. Mobility detection using everyday GSM traces. In *Proceedings of the 8th International Conference on Ubiquitous Computing*. ACM, 212–224.
- L. Wang, Y. Zheng, X. Xie, and W. Ma. 2008. A flexible spatio-temporal indexing scheme for large-scale GPS track retrieval. In *Proceedings of the 8th IEEE International Conference on Mobile Data Management*. IEEE, 1–8.
- Y. Wang, Y. Zheng, and Y. Xue. 2014. Travel time estimation of a path using sparse trajectories. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 25–34.
- L. Wei, Y. Zheng, and W. Peng. 2012. Constructing popular routes from uncertain trajectories. In *Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 195–203.
- X. Xiao, Y. Zheng, Q. Luo, and X. Xie. 2010. Finding similar users using category-based location history. In *Proceedings of the 18th Annual ACM International Conference on Advances in Geographic Information Systems*. ACM, 442–445.
- X. Xiao, Y. Zheng, Q. Luo, and X. Xie. 2014. Inferring social ties between users with human location history. *Journal of Ambient Intelligence and Humanized Computing* 5, 1 (2014), 3–19.
- H. Xie, L. Kulik, and E. Tanin. 2010. Privacy-aware traffic monitoring. *IEEE Transactions on Intelligent Transportation Systems* 11, 1 (2010), 61–70.
- C. Xu, Y. Gu, L. Chen, J. Qiao, and G. Yu. 2013. Interval reverse nearest neighbor queries on uncertain data with Markov correlations. In *Proceedings of the 29th IEEE International Conference on Data Mining*. IEEE, 170–181.



- X. Xu, J. Han, and W. Lu. 1990. RT-tree: An improved R-Tree indexing structure for temporal spatial databases. In *Proceedings of International Symposium on Spatial Data Handling*. 1040–1049.
- A. Y. Xue, R. Zhang, Y. Zheng, X. Xie, J. Huang, and Z. Xu. 2013. Destination prediction by sub-trajectory synthesis and privacy protection against such prediction. In *Proceedings of the 29th IEEE International Conference on Data Engineering*. IEEE, 254–265.
- X. Yan, J. Han, and R. Afshar. 2003. CloSpan: Mining closed sequential patterns in large datasets. In *Proceedings of the 3rd SIAM International Conference on Data Mining*. IEEE, 166–177.
- J. Yang, W. Wang, and P. S. Yu. 2003. Mining asynchronous periodic patterns in time series data. *IEEE Transactions on Knowledge and Data Engineering* 15, 3 (2003), 613–628.
- J. Yang, W. Wang, and S. Y. Philip. 2001. Infominer: Mining surprising periodic patterns. In *Proceedings of the 7th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 395–400.
- J. Yang, W. Wang, and P. S. Yu. 2002. Infominer+: Mining partial periodic patterns with gap penalties. In *Proceedings of the IEEE International Conference on Data Mining*. IEEE, 725–728.
- Y. Ye, Y. Zheng, Y. Chen, J. Feng, and X. Xie. 2009. Mining individual life pattern based on location history. In *Proceedings of the 10th IEEE International Conference on Mobile Data Management*. IEEE, 1–10.
- B. K. Yi, H. Jagadish, and C. Faloutsos. 1998. Efficient retrieval of similar time sequences under time warping. In *Proceedings of the 14th IEEE International Conference on Data Engineering*. IEEE, 201–208.
- H. B. Yin and O. Wolfson. 2004. A weight-based map matching method in moving objects databases1. In *Proceedings of the 16th International Conference on Scientific and Statistical Database Management*. IEEE, 437–410.
- J. Yin, X. Chai, and Q. Yang. 2004. High-level goal recognition in a wireless Lan. In *Proceedings of the National Conference on Artificial Intelligence*. AAAI, 578–584.
- H. Yoon, Y. Zheng, X. Xie, and W. Woo. 2012. Social itinerary recommendation from user-generated digital trails. *Journal on Personal and Ubiquitous Computing* 16, 5 (2012), 469–484.
- H. Yoon, Y. Zheng, X. Xie, and W. Woo. 2011. Smart itinerary recommendation based on user-generated GPS trajectories. In *Proceedings of the 8th IEEE International Conference on Ubiquitous Intelligence and Computing*. IEEE, 19–34.
- J. Yuan, Y. Zheng, and X. Xie. 2012. Discovering regions of different functions in a city using human mobility and POIs. In *Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 186–194.
- J. Yuan, Y. Zheng, C. Zhang, W. Xie, X. Xie, G. Sun, and Y. Huang. 2010a. T-Drive: Driving directions based on taxi trajectories. In *Proceedings of the 18th Annual ACM International Conference on Advances in Geographic Information Systems*. ACM, 99–108.
- J. Yuan, Y. Zheng, X. Xie, and G. Sun. 2011a. Driving with knowledge from the physical world. In *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 316–324.
- J. Yuan, Y. Zheng, X. Xie, and G. Sun. 2013a. T-Drive: Enhancing driving directions with taxi drivers' intelligence. *IEEE Transaction on Knowledge and Data Engineering* 25, 1 (2013), 220–232.
- J. Yuan, Y. Zheng, C. Zhang, X. Xie and G. Sun. 2010b. An interactive-voting based map matching algorithm. In *Proceedings of the 11th IEEE International Conference on Mobile Data Management*. IEEE, 43–52.
- J. Yuan, Y. Zheng, L. Zhang, X. Xie, and G. Sun. 2011b. Where to find my next passenger? In *Proceedings of the 13th International Conference on Ubiquitous Computing*. ACM, 109–118.
- N. J. Yuan, Y. Zheng, and X. Xie. 2012. *Segmentation of Urban Areas using Road Networks*. Technical Report MSR-TR-2012-65.
- N. J. Yuan, Y. Zheng, L. Zhang, and X. Xie. 2013b. T-Finder: A recommender system for finding passengers and vacant taxis. *IEEE Transaction on Knowledge and Data Engineering* 25, 10 (2013), 2390–2403.
- N. J. Yuan, Y. Zheng, X. Xie, Y. Wang, K. Zheng, and H. Xiong. 2015. Discovering urban functional zones using latent activity trajectories. *IEEE Transactions on Knowledge and Data Engineering* 27, 3 (2015), 1041–4347.
- D. Zhang, N. Li, Z. Zhou, C. Chen, L. Sun, and S. Li. 2011. iBAT: Detecting anomalous taxi trajectories from GPS traces. In *Proceedings of the 13th International Conference on Ubiquitous Computing*. ACM, 99–108.
- F. Zhang, D. Wilkie, Y. Zheng, and X. Xie. 2013. Sensing the pulse of urban refueling behavior. In *Proceedings of the 15th International Conference on Ubiquitous Computing*. ACM, 13–22.
- F. Zhang, N. J. Yuan, D. Wilkie, Y. Zheng, and X. Xie. 2015. Sensing the pulse of urban refueling behavior: A perspective from taxi mobility. *ACM Transactions on Intelligent Systems and Technology* 6 (2015), 3.
- K. Zheng, Y. Zheng, X. Xie, and X. Zhou. 2012a. Reducing uncertainty of low-sampling-rate trajectories. In *Proceedings of the 28th IEEE International Conference on Data Engineering*. IEEE, 1144–1155.

- K. Zheng, Y. Zheng, N. J. Yuan, and S. Shang. 2013a. On discovery of gathering patterns from trajectories. In *Proceedings of the 29th IEEE International Conference on Data Engineering*. IEEE, 242–253.
- K. Zheng, Y. Zheng, N. J. Yuan, S. Shang, and X. Zhou. 2014a. Online discovery of gathering patterns over trajectories. *IEEE Transaction on Knowledge and Data Engineering* 26, 8 (2014), 1974–1988.
- V. W. Zheng, B. Cao, Y. Zheng, X. Xie, and Q. Yang. 2010a. Collaborative filtering meets mobile recommendation: A user-centered approach. In *Proceedings of the 24th AAAI Conference on Artificial Intelligence*. AAAI, 236–241.
- V. W. Zheng, Y. Zheng, X. Xie, and Q. Yang. 2010b. Collaborative location and activity recommendations with gps history data. In *Proceedings of the 19th International Conference on World Wide Web*. ACM, 1029–1038.
- V. W. Zheng, Y. Zheng, X. Xie, and Q. Yang. 2012b. Towards mobile intelligence: Learning from GPS history data for collaborative recommendation. *Artificial Intelligence* 184–185 (2012), 17–37.
- Y. Zheng. 2011. Location-based social networks: users. *Computing with Spatial Trajectories*, Y. Zheng and X. Zhou (Eds.). Springer, 243–276.
- Y. Zheng. 2012. Tutorial on location-based social networks. In *Proceedings of the 21st International Conference on World Wide Web*. ACM.
- Y. Zheng, L. Capra, O. Wolfson, and H. Yang. 2014b. Urban computing: Concepts, methodologies, and applications. *ACM Transactions on Intelligent Systems and Technology* 5, 3 (2014), 38–55.
- Y. Zheng, X. Chen, Q. Jin, Y. Chen, X. Qu, X. Liu, E. Chang, W. Ma, Y. Rui, and W. Sun. 2014c. A Cloud-based knowledge discovery system for monitoring fine-grained air quality. MSR-TR-2014-40.
- Y. Zheng, Y. Chen, Q. Li, X. Xie, and W.-Y. Ma. 2010c. Understanding transportation modes based on GPS data for Web applications. *ACM Transactions on the Web* 4, 1 (2010), 1–36.
- Y. Zheng, Y. Chen, X. Xie, and W.-Y. Ma. 2009a. GeoLife2.0: A location-based social networking service. In *Proceedings of the 10th IEEE International Conference on Mobile Data Management*. IEEE, 357–358.
- Y. Zheng, S. E. Koonin, and O. E. Wolfson. 2013. *Proceedings of the 2nd ACM SIGKDD International Workshop on Urban Computing*. ACM.
- Y. Zheng, Q. Li, Y. Chen, and X. Xie. 2008a. Understanding mobility based on GPS data. In *Proceedings of the 11th International Conference on Ubiquitous Computing*. ACM, 312–321.
- Y. Zheng, L. Liu, L. Wang, and X. Xie. 2008b. Learning transportation mode from raw GPS data for geographic application on the Web. In *Proceedings of the 17th International Conference on World Wide Web*. ACM, 247–256.
- Y. Zheng, F. Liu, and H. P. Hsieh. 2013b. U-Air: When urban air quality inference meets big data. In *Proceedings of the 19th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. ACM, 1436–1444.
- Y. Zheng, T. Liu, Y. Wang, Y. Liu, Y. Zhu, and E. Chang. 2014c. Diagnosing New York City’s noises with ubiquitous data. In *Proceedings of the 16th ACM International Joint Conference on Pervasive and Ubiquitous Computing*. ACM, 715–725.
- Y. Zheng, Y. Liu, J. Yuan, and X. Xie. 2011a. Urban computing with taxicabs. In *Proceedings of the 13th International Conference on Ubiquitous Computing*. ACM, 89–98.
- Y. Zheng and X. Xie. 2011b. Learning travel recommendations from user-generated GPS traces. *ACM Transactions on Intelligent Systems and Technology* 2, 1 (2011), 2–19.
- Y. Zheng, X. Xie, and W.-Y. Ma. 2010d. GeoLife: A collaborative social networking service among user, location and trajectory. *IEEE Data Engineering Bulletin* 33, 2 (2010), 32–40.
- Y. Zheng, L. Zhang, Z. Ma, X. Xie, and W.-Y. Ma. 2011c. Recommending friends and locations based on individual location history. *ACM Transaction on the Web* 5, 1 (2011), 5–44.
- Y. Zheng, L. Zhang, X. Xie, and W.-Y. Ma. 2009b. Mining interesting locations and travel sequences from GPS trajectories. In *Proceedings of the 18th International Conference on World Wide Web*. ACM, 791–800.
- Y. Zheng, L. Zhang, X. Xie, and W.-Y. Ma. 2009c. Mining correlation between locations using human location history. In *Proceedings of the 17th Annual ACM International Conference on Advances in Geographic Information Systems*. ACM, 352–361.
- Y. Zheng and X. Zhou. 2011. *Computing with Spatial Trajectories*. Springer.
- Y. Zhu, Y. Zheng, L. Zhang, D. Santani, X. Xie, and Q. Yang. 2011. *Inferring Taxi Status using GPS Trajectories*. Technical Report MSR-TR-2011-144.
- GeoLife Data: <http://research.microsoft.com/en-us/downloads/b16d359d-d164-469e-9fd4-daa38f2b2e13/default.aspx>.
- T-Drive Data: <http://research.microsoft.com/apps/pubs/?id=152883>.

Trajectory with transportation modes: <http://research.microsoft.com/apps/pubs/?id=141896>.

User check-in data: <https://www.dropbox.com/s/4nwb7zpsj25ibyh/check-in%20data.zip>.

Hurricane trajectory (HURDAT): <http://www.nhc.noaa.gov/data/hurdat>.

The Greek Trucks Dataset,” <http://www.chorochronos.org>.

Movebank data: <https://www.movebank.org/>.

Received October 2013; revised May 2014; accepted November 2014