



Trajectory splicing

Qiang Lu¹ · Rencai Wang^{1,3} · Bin Yang² · Zhiguang Wang¹

Received: 23 September 2018 / Revised: 25 June 2019 / Accepted: 30 June 2019 / Published online: 18 July 2019
© Springer-Verlag London Ltd., part of Springer Nature 2019

Abstract

With continued development of location-based services, large amount of trajectory become available, which records moving object location at a time. If the trajectory collected by different location-based services come from the same moving object, they are *spliceable trajectories*, which contribute to interesting holistic behavior of the moving object. In this paper, we consider how to efficiently identify spliceable trajectory. More specifically, we first formalize a splicing model to characterize spliceable trajectory, where their time are disjoint, and the distance between them are close. Next, to efficiently implement the model, we design three components: a disjoint time index, a directed acyclic graph of trajectory location connection, and trajectory splicing algorithm. The disjoint time index is a disjoint time set of each trajectory for finding disjoint trajectory efficiently. The directed acyclic graph contributes to discovering groups of spliceable trajectory. Based on the identified groups, the splicing algorithm *findmaxCTR* find maximal groups containing all spliceable trajectory. Although the splicing algorithm is efficient in some practical application, its running time is exponential. Therefore, an approximate algorithm *findApproxMaxCTR* is proposed to find trajectory which can be spliced with each other with a certain probability in polynomial time. Finally, experiment on trajectory data set demonstrate that the model and its components are effective and efficient.

Keywords Trajectory compression · Trajectory fusion · Trajectory reconstruction · Trajectory linking

✉ Qiang Lu
luqiang@cis.umd.edu.cn

Rencai Wang
rcwang3@infotech.com

Bin Yang
byang@cbaa.dk

Zhiguang Wang
zgwang@cis.umd.edu.cn

¹ Beijing Key Lab of Petroleum Data Mining, China University of Petroleum-Beijing, Beijing, China

² Department of Computer Science, Aalborg University, Aalborg, Denmark

³ IFLYTEK CO.,LTD, Hefei, China

1 Introduction

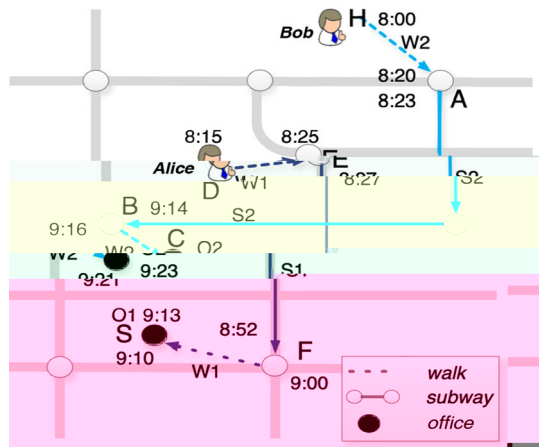
Information technology is almost everywhere in our daily life, which collects a lot of information from different digital devices [4,10]. Especially, the location-based system based on mobile devices, such as GPS, mobile phone, and near-field communication (NFC) terminal, generate large amount of trajectory of moving object. Usually, each individual item in the network is identified by a unique ID code to identify each trajectory. For example, a mobile phone network can identify a trajectory by its telephone number, while an NFC terminal identifies it by its device ID. Since multiple items may take the same moving object at different time and place, each system gathers the object's partial trajectory. Recording a **complete trajectory** of a moving object from the partial trajectory collected in a system, named **trajectory splicing**, is essential for many applications, such as anomaly behavior detection [21,22], data fusion, and trajectory data mining [46]. The following case history in Fig. 1 elaborates trajectory splicing.

Every weekday, Alice and Bob go to work by walking and taking the subway, a history in Fig. 1. Their movement generates six partial trajectories: W1, S1, O1, W2, S2, and O2, where the mobile context of Alice are W1 and W2; the subway check-in system context are S1 and S2; the office check-in system context are O1 and O2. Their complete trajectory can be reconstructed based on the temporal location of the partial trajectory. For example, S2 is more likely to be with W2 than W1, because the end point of S2 is close to the end of W2, and the time interval of S2 [8:23,9:14] can be embedded into the time gap of W2 (8:20, 9:16). Similarly, O2 can be with W2. So, connecting W2, S2, and O2 can be for Bob's whole trajectory.

According to the above case, finding a group of traceable trajectory must satisfy the following three requirements. The first is the **disjoint time constraint** that requires that the time interval of traceable trajectory in the group should not overlap with each other. The second is the **spatial constraint** that requires that the distance between their end point should be near by with each other. The third is the **maximal group constraint** that requires that the group of traceable trajectory should be maximal and should not be contained by other group. That means connecting a maximal traceable trajectory is a possible to reconstruct a complete trajectory.

However, it is non-trivial to find traceable trajectory to satisfy the above constraints owing to the following three challenges. The first challenge is that the process of finding trajectory that satisfy the disjoint time constraint is a time-consuming. The process includes the following: selecting trajectories in all time gaps of a trajectory and counting the number of trajectories that belong to the same trajectory. For example, in Fig. 1, W2 has three time gaps: $(-\infty$

Fig. 1 The case of trajectory linking



which connect trajectories without linking other traceable trajectories. The other is the **indirect splice**, which connect trajectories by linking other traceable trajectories. For example, in Fig. 1, W2 and S2 are connected directly, while W2 and O2 are connected by S2. The indirect splice make the process of linking trajectories complicated because it need to find other trajectories to determine whether the trajectories can be connected or not. To the best of our knowledge, knowledge graph pattern mining [8,9,19,24,27,36,45] or trajectory clustering [24,25] cannot find groups of traceable trajectories, because the discrete groups of trajectories according to the similarity between them rather than the relation of direct (indirect) splice. Although trajectory linking [38] is close to the challenge, it can only find trajectory splice trajectories and is not suitable for mining multiple trajectories that are the direct or indirect splice.

The third challenge is that it must find a man traceable trajectory of a moving object available. In general, if a method want to acquire a group of traceable trajectories which are not contained by other groups, it need to traverse all possible combination of traceable trajectories for a moving object. For example, in the above case, to recognize Bob trajectory, the groups, such as {W2, S2}, {W2, O2}, {S2, O2}, and {W2, S2, O2}, must be traversed. Namely, it need to find a group of traceable trajectories which must completely fill a specific spatiotemporal range. So, it is a bin-packing problem and is NP-hard [23]. The design of an approximation scheme heuristic method is the key to deal with the problem.

In order to deal with the above challenge, a sliced model is defined to formalize the above requirement of traceable trajectories. Based on the sliced model, trajectories are segmented into sub-trajectories according to a fixed threshold. A B^+ -tree [7] is used to store the sub-trajectories. For speeding up the process of finding disjoint time set, the index of disjoint time set called **DT-index** is constructed to keep intermediate result of exchanging the disjoint time set in each time slice. Moreover, the DT-index is a multi-resolution structure like a pyramid and can be intermediate result of time slice with different length, representing slice with different time interval. For example, assuming that the DT-index consists of intermediate result of one, two, and four days, if a query time interval is 4.5 days, the DT-index can find disjoint time set within the four days, and the B^+ -tree can find disjoint time set within the 0.5 days. Based on the above two indexes, an algorithm **DTTR** is proposed to obtain all disjoint time set within a specific time interval.

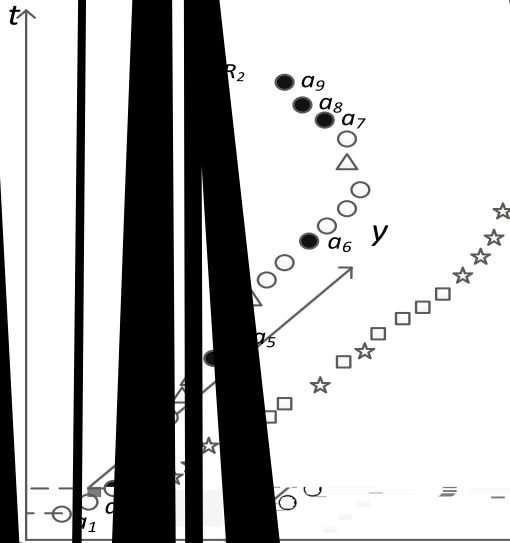
In order to find feasible trajectories, a directed acyclic graph of feasible location connections called *STLC-DAG* is created to connect feasible trajectories by their time and location. Once the algorithm *createSTLC-DAG* has created the graph, it can obtain the feasible set of trajectories that can coincide with a specific trajectory. For example, in the above case, the algorithm can find $S2$ feasible set $\{W2\}$, $W2$ feasible set $\{S2, O2\}$, and $O2$ feasible set $\{W2\}$. Moreover, the feasible set forms a **splice graph**, where each node is a trajectory, and the edge between two nodes represents that the two trajectories are feasible. For instance, the node $S2$ has one edge which connects the node $W2$, and $W2$ has edges which connect $S2$ and $O2$. Thus, in the splice graph, a cycle is a group of feasible trajectories. For addressing the third challenge, an algorithm *dMCTC* is proposed to find all maximal

plication

ation

Definition

- A *TR* database
- A sample point (i, a_i)
- The i th *TR* in Ω
- A TR_i disjoint time interval
- A set of *TR* that can be
- The j th trajectory
- The time interval of
- The number of *TR* in
- The number of all *STR*
- A complete trajectory
- Return the first element
- Return the last element
- The Euclidean distance
- The Euclidean distance between
- The time interval of a
- The time interval of a
- The gap between two
- The gap of $ti(TR_i)$



moving object : $CTR_1 = \{TR_A, TR_B, TR_C\}$, which include the trajectory with identifier A, B , and C , and $CTR_2 = \{TR_D, TR_E\}$, which include the trajectory with identifier D and E .

In a trajectory, two adjacent, p_i and p_{i+1} , are **connectable** if $speed(p_i, p_{i+1}) \geq e$, where e is a speed threshold and

$$speed(p_i, p_{i+1}) = \frac{d(p_i, p_{i+1})}{|p_{i+1}.t - p_i.t|} \tag{1}$$

where $d(p_i, p_{i+1})$ is the Euclidean distance between adjacent p_i and p_{i+1} . Given a sequence of adjacent points in a trajectory TR_i , if an adjacent pair of points in the sequence are connectable, the sequence is **connectable** in that it has one continuous movement. Moreover, if other connectable sequence do not contain a connectable sequence, the connectable sequence is called **sub-trajectory** (denoted as STR). In particular, we use STR_j^i to denote the j th sub-trajectory in trajectory TR_i . For example, trajectory TR_A in Fig. 2 has 4 sub-trajectories: $STR_A^1 = \langle a_1, a_2, a_3 \rangle$, $STR_A^2 = \langle a_4, a_5 \rangle$, $STR_A^3 = \langle a_6 \rangle$, and $STR_A^4 = \langle a_7, a_8, a_9 \rangle$. A sub-trajectory is the **atomic** computational unit in this paper.

The **time interval** of the sub-trajectory, denoted as $ti(STR)$, is $[first(STR).t, last(STR).t]$, where the function $first(\cdot)$ and $last(\cdot)$ denote the first and last adjacent points in the sub-trajectory STR , respectively. The **time interval** of the trajectory is the set of time interval of all its sub-trajectories, denoted as $ti(TR_i) = \bigcup_{STR_j^i \in TR_i} ti(STR_j^i)$.

The **gap** between two adjacent sub-trajectories STR_j^i and STR_m^n , denoted as $gap(STR_j^i, STR_m^n)$, is defined by Eq. 2.

$$gap(STR_j^i, STR_m^n) = (last(STR_j^i).t, first(STR_m^n).t) \tag{2}$$

Moreover, the **gap** of trajectory TR_i in the time interval T , denoted as $ga(TR_i)$, is defined by Eq. 3.

$$gap(TR_i) = T - ti(TR_i) = T - \bigcup_{STR_j^i \in TR_i} ti(STR_j^i) \tag{3}$$

For example, the time interval of trajectory TR_A , denoted as $ti(TR_A)$, is $\{[t_1, t_2], [t_3, t_4], [t_5, t_5], [t_6, t_7]\}$. Given $T = [t_0, t_8]$, we have $gap(TR_A) = \{(t_0, t_1), (t_2, t_3), (t_4, t_5), (t_5, t_6), (t_7, t_8)\}$.

2.2 Spliceable trajectories

If two trajectories TR_i and TR_j can be divided into a complete trajectory, the must meet the **disjoint time constraint** that is, they are that their interval time should not overlap each other, namely $ti(TR_i) \subset gap(TR_j)$. Given a trajectory TR_i , all the trajectories that meet the disjoint time constraint with TR_i constitute the **disjoint time set** of TR_i , denoted as DT_i . In Fig. 2, since $ti(TR_B) \subset gap(TR_A)$ and $ti(TR_C) \subset gap(TR_A)$, we have $DT_A = \{TR_B, TR_C\}$.

In addition to the aforementioned temporal constraint, if TR_i and TR_j are spliceable, the must also meet the **spatial constraint**, meaning that the sub-trajectories from TR_i and TR_j must be close enough to each other. To formally define the spatial constraint, we introduce two concepts: **iceable** and **iceable**.

Definition 1 Given two sub-trajectories STR_j^i and STR_m^n from trajectories, respectively, and a distance threshold γ , if they do not overlap each other on the time dimension and

the distance between them is less than γ^1 , the \mathbb{V} -observable trajectory STR_i^j and STR_m^n form a *iceable pair*, denoted as (STR_i^j, STR_m^n) .

Definition 2 Given observable \mathbb{V} , if the \mathbb{V} -observable in the given trajectory can contain a \mathbb{V} -observable sequence $(STR_i^j, \dots, STR_m^n)$ such that an \mathbb{V} -observable neighborhood \mathbb{V} -trajectory is a *spliceable pair*, the trajectory is called *iceable trajectory*.

Based on the above definition, we first introduce the concept of *iceable trajectory* to form the maximal good connection, which ensures that the good of spliceable trajectory should not be contained by other good. Then, we define the *ice degree* to analyze the complete trajectory.

Definition 3 If other good of spliceable trajectory do not contain a good of spliceable trajectory, the good form a **complete trajectory**, denoted as *CTR*.

Definition 4 The *ice degree*, which consists of \mathbb{V} -observable factors: the ratio of the sum of the distance between different trajectory to the distance of *CTR* and the ratio of the sum of time gap to the time interval of *CTR*, is denoted as *ice degree* of connection between trajectory in a *CTR*, defined by Eq. 4.

$$\begin{aligned}
 dg(CTR) &= \frac{\sum_{(STR_i^j, STR_m^n) \in CTR} d(STR_i^j, STR_m^n)}{distance(CTR)} \\
 &\times \frac{\sum_{(STR_i^j, STR_m^n) \in CTR} gap(STR_i^j, STR_m^n)}{time(CTR)} \tag{4}
 \end{aligned}$$

Here (STR_i^j, STR_m^n) is a *spliceable pair* in the *CTR*; $d(STR_i^j, STR_m^n)$ is the distance between \mathbb{V} -observable trajectory STR_i^j and STR_m^n ; $distance(CTR)$ is the sum of distance between consecutive elements in *CTR*, namely $distance(CTR) = \sum_{p_i \in CTR} d(p_i, p_{i+1})$, in which p_i and p_{i+1} are \mathbb{V} -observable consecutive elements in the *CTR*; $time(CTR) = last(CTR).t - first(CTR).t$.

Based on the definition, $dg(CTR) \in (0, 1)$ and the smaller the *ice degree* $dg(CTR)$, the closer trajectory in the complete trajectory *CTR*. For example, in Fig. 2, assuming that the distance factor in Alice and Bob are the same as 0.02, $dg(Alice) = 0.02 \times (((8 : 27 - 8 : 25) + (9 : 00 - 8 : 52) + (9 : 13 - 9 : 10)) / (9 : 13 - 8 : 15)) \approx 0.0448$, and $dg(Bob) = 0.02 \times ((8 : 23 - 8 : 20) + (9 : 16 - 9 : 14) + (9 : 23 - 9 : 21)) / (9 : 23 - 8 : 00) \approx 0.0017$. So, due to $dg(Bob) < dg(Alice)$, the complete trajectory of Bob is better than that of Alice.

2.3 Problem definition

According to the above definition, we form the problem of trajectory splicing by the trajectory splicing problem.

Definition 5 From a dataset of trajectory, according to a given time interval, the *ice degree* of each complete trajectory sequence $CTRS = \langle CTR_1, \dots, CTR_n \rangle_{\mathbb{V}}$ here each complete trajectory *CTR* is ranked by its *splice degree*.

¹ Namely $(ti(STR_m^n) \subset gap(STR_i^j, STR_i^{j+1})) \cap (ti(STR_i^j) \subset gap(STR_m^{n-1}, STR_m^n)) \cap (d(last(STR_i^j), first(STR_m^n)) \leq \gamma)$.

1



(5)

each in the
 $\{B, C, D, E\}$
 $DT_A^{2,d} = \{D, E\}$

g. 4.

i_i of each trajectory TR_i
of in A (endix A).

$\dots \cup DF_i^{n,d}]$ (6) Fig.



where $|T| = n \times d$, d is the length of the time slice, n is the number of time slices, and P_i is a set which contains all trajectories that appear in T except the trajectory TR_i . For example, in Fig. 4, if $T = [0, 3d]$, $DT_D(T) = P_D^{0,3d} - [(P_D^{0,3d} - DT_D^{1,d}) \cup DF_D^{2,d} \cup DF_D^{3,d}] = \{A, B, C, E\} - [\{A, B, C, E\} - \{E\}] \cup \{A\} \cup \emptyset = \{E\}$.

If T is too long, there are many time slices in T , and E. 6 contains many non-operations of DF so that the computation of E. 6 is time-consuming. To alleviate the situation, we partition the time dimension into multiple levels of time slices. For instance, one level of time slice is a day, and another level is a week or a month. So, if $|T|$ is one month, E. 6 can be computed by only one DF on the monthly level of time slices rather than about 30 DF on the daily level.

(2) The recursive definition

Based on the above analysis, we design the disjoint time index (called **DT-index**), which includes a *DT-tree* and a *DF-tree* that store the disjoint time set DT of each trajectory and its decomposition DF on different levels of time slices, respectively, as shown in Fig. 4. The two trees have the same structure. The *DT-tree* (*DF-tree*) consists of a single root node, leaf nodes, and non-root, non-leaf nodes. The detailed data structure of the nodes are as follows.

A **node**, which may have multiple children, are their *ID*. A *ID* is both a time interval and a filename, hence denoting a time interval T , its children and their files are located in T itself.

A **leaf node** is a pair of $\langle i, DT_i \rangle$ or $\langle i, DF_i \rangle$ in a specific time slice. For example, in Fig. 4, $DT^{3,d}$ records pairs $\langle A, \{B, C\} \rangle$, $\langle B, \{A, C\} \rangle$ and $\langle C, \{A, B\} \rangle$.

A **non-leaf node**, a **node** only has two children. It stores its children *ID* and pairs of $\langle i, DT_i \rangle$ or $\langle i, DF_i \rangle$, where DT_i or DF_i can be computed by E. 6 or 5, respectively.

3.2 Processing query

With the B^+ -tree and the DT -index, we implement an algorithm $QueryDTsTR$, which is able to find the disjoint time set DT of each trajectory and all sub-trajectory (denoted as $STRSet$) in a time interval T , as shown in Algorithm 1.

Algorithm 1: $queryDTsTR$

Input: B^+ -tree, DT -Index, T
Output: $DT(T)$, $STRSet$

```

1  $STRSet, DT(T_1), R(T_1), R(T_2), P = readsTR(B^+tree, T)$ ;
2  $DT(T_2) = Equation\ 7$ ;
3  $DT = (DT(T_1) \cup R(T_1)) \cap (DT(T_2) \cup R(T_2))$ ;
4 return  $DT, STRSet$ ;

```

The execution time interval T consists of two parts: One is a set of two time intervals, which is an interval in the DT -index, denoted as $T_1 = \{t_1, t_2\}$; the other is the time interval that contains n intervals in the DT -index, denoted as T_2 . For example, given $T = [8 : 35 : 11 : 25]$ and the minimal time intervals are $T_1 = \{[8 : 35 : 9 : 00], [11 : 00 : 11 : 25]\}$, and $T_2 = [9 : 00 : 11 : 00]$. With the B^+ -tree, it is easy to find all trajectories P and their sub-trajectory $STRSet$ in T . Meanwhile, exchanging the sub-trajectory can obtain a trajectory set $R(T_1)$, where each trajectory is a pair in T_1 but not in T_2 , a trajectory set $R(T_2)$, where each trajectory is a pair in T_2 but not in T_1 , and a disjoint time set $DT(T_1)$ in the set T_1 . The function $readSTR$ at Line 1 implements the above process. Then, with the DT -index, the code at Line 2 computes the disjoint time set $DT(T_2)$ by Eq. 7. At last, the code at Line 3 gets the disjoint time set DT in T .

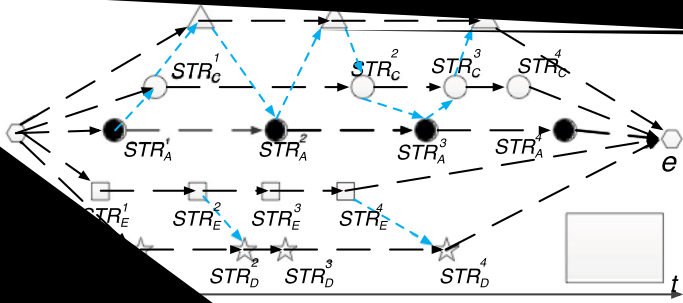
The algorithm can be refined based on the following observations. One is that, in general, compared with the set T_2 , the set T_1 is easier to check that there are effective sub-trajectory (STR) in T_1 . Hence, finding the disjoint time set $DT(T_1)$ is faster. The other is that, since the disjoint time set DT of each trajectory has been established on multiple time scales in the DT -index, only a small amount of nodes need to be searched from the index in order to compute the disjoint time set $DT(T_2)$ by Eq. 7. So, finding $DT(T_2)$ is also faster.

3.3 Splicing trajectory

3.3.1 Finding spliceable trajectories

We design an algorithm $createSTL-DAG$ to discover spliceable trajectory by constructing a directed acyclic graph of sub-trajectory location connection ($STLC-DAG$), which is defined as $STLC-DAG = (V, E)$, where

- the vertex set V consists of all sub-trajectory ($STRSet$), a start vertex s , and an end vertex e , namely $V = \{STRSet\} \cup \{s, e\}$;
- the edge set E consists of two categories of directed edges. One, denoted as E_s , is the directed edge that connects two sub-trajectories in the same trajectory. The other, denoted as E_d , is the directed edge that connects a spliceable pair $\langle STR_i^j, STR_m^n \rangle$, as shown in Fig. 5.



Since there are the multiple existence in the graph, all first edges from the existence contribute a **candidate vertex set (CVS)**, which is defined by Eq. 8.

$$CVS(STR_i^j) = \{STR_m^n | STR_m^n = first(\{ti(STR_m^k) \subset gap(STR_i^{j+1}, STR_i^j)\}), m \in DT_i\} \tag{8}$$

For example, in Fig. 5, $CVS(STR_A^1) = \{STR_B^1, STR_C^1, STR_D^1, STR_E^2\}$.

Lemma 3 shows that when a trajectory cannot coincide with another trajectory, the edge between the two trajectories can be deleted. Moreover, the deletion does not cause the result of the lizable trajectory to change.

The procedure of constructing the graph *STLC-DAG* is shown in Algorithm 2. The input argument is the set of trajectories *STRSet* and the disjoint time set *DT*, respectively. In the algorithm *queryDTsTR*, and γ is a distance threshold. The algorithm will return a set $SP = \{SP_1, \dots, SP_n\}$, where each SP_i is a group of lizable trajectories.

Algorithm 2: *createSTLC-DAG*

```

Input: STRSet,  $\gamma$ , SP = DT
Output: SP
1 sortByStartTime(STRSet);
2 DAG.V = STRSet  $\cup$  {s, e};
3 DAG.E.Es = createEsEdge(STRSet, s, e);
4 C =  $\phi$ ;
5 for k = 0; k < len(STRSet); k++ do
6   STR_i^j = STRSet[k];
7   for each STR_k^v  $\in$  sortByDes(C.get(STR_i^j)) do
8     sg = 0;
9     repeat
10      if !existPath(STR_k^v, STR_i^j, SP_k, DAG) then
11        DAG.E.Ed.delEdges(TR_k, TR_i);
12        SP_i = SP_i - k;
13        SP_k = SP_k - i;
14        C.del((TR_i, TR_m));
15      sg = |C|;
16    else
17      sg = sg - 1;
18    (STR_k^v, STR_i^j)  $\leftarrow$  C.next(STR_k^v, STR_i^j);
19  until (STR_k^v, STR_i^j)  $\neq$   $\phi$  && sg > 0;
20 canTRSet = CVS(STR_i^j);
21 for each STR_m^n  $\in$  canTRSet do
22   if d(STR_i^j, STR_m^n)  $\leq$   $\gamma$  then
23     DAG.E.Ed.addEdge(STR_i^j, STR_m^n);
24   else
25     C.add((STR_m^n, STR_i^j));
26 return SP;
    
```

Initially, the algorithm extract all trajectories in *STRSet* by their start time, create all edges, and connect the edges that belong to the same trajectory (Line 1-3). *C* is a container that a pair of trajectories which are likely to be indirectly lizable by

other β -trajectory (Line 4). For each β -trajectory STR_i^j in $STRSet$, its candidate set $CVS(STR_i^j)$ is finally obtained by E. 8. Then, the algorithm creates a directed edge between the two β -trajectories STR_i^j and STR_m^n

After Algorithm 2 finishes its running, if there exists an edge between two trajectories in the graph *STLC-DAG*, the trajectories can be spliced according to Theorem 1. At the same time, the algorithm can find groups of spliceable trajectories SP_i , where each SP_i is a set of trajectories that can be directly or indirectly spliced with the trajectory TR_i based on Theorem 2.

Theorem 1 *If there exists a directed edge between two trajectories in the graph *STLC-DAG*, the two trajectories can be spliced.*

Theorem 2 *For each $SP_i \in SP$, where SP is one of the output parameters of algorithm 2, SP_i is a set of trajectories that can splice with the trajectory TR_i .*

The above two proofs are provided in Appendix B.

3.3.2 Finding maximum trajectories

Algorithm 5: *findApproxMaxCTR*

```

Input:  $SP, SUBG = V, CAND = V, d, k, c = 0, fCTR = \phi$ 
Output:  $fCTRSet: a fCTR \in$ 
1 if  $SUBG \neq \phi$  then
2   if  $c = k$  then
3     if  $|CAND| \leq (d - k)$  then
4        $fCTR \leftarrow CAND;$ 
5     else
6        $fCTR \leftarrow takeFirst(CAND, d - k);$ 
7      $fCTRSet \leftarrow fCTR;$ 
8     return ;
9    $i = subscript(max|SUBG \cap SP_i|), i \in SUBG;$ 
10   $branch = CAND - SP_i;$ 
11  while  $branch \neq null$  do
12     $b = takeFirst(branch);$ 
13     $fCTR \leftarrow b;$ 
14     $SUBG_b = SUBG \cap SP_b;$ 
15     $CAND_b = CAND \cap SP_b;$ 
16     $fCTRSet = findApproxMaxCTR(SP, SUBG_b, CAND_b, d, k, c + 1, fCTR);$ 
17     $CAND = CAND - \{b\};$ 
18 else
19    $fCTRSet \leftarrow fCTR;$ 
20 et ; et ;  $n fCTRSet;$ 

```

Based on the above analysis, we design an algorithm *findApproxMaxCTR* to find a ϵ -approximate maximal sliced path. The detailed pseudocode of *findApproxMaxCTR* is listed in Algorithm 5. The algorithm is similar to Algorithm 4 except the code on Line 2–8. The additional parameters are a follow: d, k , and c , here d is used to limit the number of sliceable trajectory in one complete trajectory; k , which is used to limit the time of interaction between two SP , is a recursive depth of the algorithm; and c records the current time of completing interaction in a sliced path $fCTR$. The code on Line 2–8 helps help to deal with trajectory in $CAND$ when $c = k$. If the size of $CAND$ is less than $d - k$, all trajectory in $CAND$ are added into $fCTR$ (Line 3–4). If the size is more than $d - k$, the first $(d - k)$ trajectory are added into $fCTR$ (Line 6).

4 Time complexity analysis

In this section, we analyze the running time of the above algorithm and ignore algorithm in the recursive step, such as the construction of B^+ -tree and DT-index, because they can be done offline. Let $T(function)$ be the running time of the function, M be the number of sub-trajectory, and N be the number of trajectory.

Lemma 7 *For the algorithm queryDTsTR, if the query time interval T consists of time slices from the DT-index, namely $T_1 = 0$ and $T_2 \neq 0$, the running time of queryDTsTR is $O(N^2)$; if the query time interval T does not contain the time slice for the DT-index, namely $T_2 = 0$ and $T_1 \neq 0$, the running time of queryDTsTR is $O(M^2)$.*

Proof Since all sub-trajectory are indexed by B^+ -tree, the time of finding m sub-trajectory is $O(\log_b^{|\Omega|} + M)$. $|\Omega|$ and b are constant. And, $\log_b^{|\Omega|} \ll M$. So, the running

time of reading all sub-objects in T is $O(M)$. At the same time, $R(T_1)$ and $R(T_2)$ can be obtained. If $T_1 = 0$, $DT(T_1)$ does not need to be computed. Therefore, $T(readSTR) = O(M)$. If $T_1 \neq 0$, the running time of computing $DT(T_1)$ is $O(M^2)$. And, $T(readSTR) = O(M^2)$. If $T_2 = 0$, Eq. 7 does not need to be computed. So, $T(queryDTsTR) = O(M^2)$.

If $T_2 \neq 0$, given that T_2 consists of k time slices, which are in different levels in DT -index, k nodes in the DT -index need to be read. Each node contains no more than N items in which there are at most N TR . According to Eq. 7, $T(Eq. 7) = O(kN^2)$. The running time of interaction between $DT(T_1)$ and $DT(T_2)$ is $O(N^2)$. So, $T(queryDTsTR)$ is $O(N^2)$. \square

Lemma 8 *The running time of the algorithm createSTLC-DAG is $O(M^2N^2)$.*

Proof Let $P = \sum_{i=1}^N |DT_i|$, where $DT_i \in DT$. So, $N \leq P \leq N^2$. The running time of creating vertices (Line 3) and edge (Line 4) both are $O(M)$. In each loop (Line 5), $T(getCandSet) = O(m_k)$, where $m_k = |CVS(i, j)|$. And, the number of loops between Line 21 and 25 at $O(m_k)$. $T(addEdge)$ and $T(add)$ both are $O(1)$. The number of creating all edge in E_d (Line 20–25) is $\sum_{k=1}^M m_k$ since $len(STRSet) = M$. According to $CVS(STR_i^j)$ (Eq. 8), $m_k \leq DT_i$.

Since more sub-objects in TR exist in level $|DT_i|$, the number of all edge is $\sum_{k=1}^M m_k$ and $\sum_{k=1}^M m_k \leq \frac{kM}{N} \times P$, where $k \ll N$. Moreover, the running time of *pseudocode* on Line 20–25 is $O(\frac{M}{N} \times P)$. If all edge are added into DAG (Line 23), C is empty. If all edge are added into C (Line 25), the longest time that *existPath* costs is $\frac{M}{N} \times P$ because *delEdges* (Line 11) can delete some edge. $T(existPath)$ depends on the number of vertices and edge between the two sub-objects STR_k^v and STR_m^u . So, $T(existPath) = O(M + \frac{M}{N} \times P)$. The running time of operation on Line 11–17 all is $O(1)$. The running time of *pseudocode* on Line 5–19 is $O(\frac{M}{N} \times P \times (M + \frac{M}{N} \times P)) = O(\frac{M^2}{N} \times P + \frac{M^2}{N^2} \times P^2)$.

Thus, $T(createSTLC-DAG) = O(M + \frac{M}{N} \times P + \frac{M^2}{N} \times P + \frac{M^2}{N^2} \times P^2) = O(\frac{M^2}{N} \times P + \frac{M^2}{N^2} \times P^2) = O(\frac{M^2}{N} \times (P + \frac{P^2}{N}))$. Owing to $P \leq N^2$, $T(createSTLC-DAG) = O(M^2N^2)$ \square

Lemma 9 *The running time of the algorithm findMaxCTR is $O(3^{N/3})$.*

Proof See Theorem 3 of [34]. \square

Lemma 10 *Let D be a maximal degree of vertices in the SP-set graph. The running time of the algorithm findApproxMaxCTR is $O(N(N - D)C_{k-1}^{D-1})$. Moreover, if k in Eq. 11 is a small numerical value, the running time of the algorithm findApproxMaxCTR is $O(CN^2)$, where C is a constant.*

Proof When the algorithm executes the code on Line 11 for the first time, $|branch| = N - D$. The algorithm will go to the branch SP_b , where b is the maximal degree of vertices $b \leq D$. Therefore, $|SUBG_b| \leq D$. When it executes the code on Line 11 for the second time, $|branch| \leq D - 1$. When it executes the code on Line 11 for the third time, $|branch| \leq D - 2$.

Each branch contains the above code until the death of iteration: each k . At the death of iteration, $|branch|$ decreases. Moreover, in the $k - 1$, $|branch| \leq D - k + 1$. According to Theorem 1 of [34], the algorithm generates all maximal cliques without duplication. So, each branch in the death is looked at a combination C_{k-1}^{D-1} . The running time of $SUBG \cap SP_i$ on Line 9 is $O(N)$. Thus, $T(findApproxMaxCTR) = O(N(N - D)C_{k-1}^{D-1})$. When k is small, C_{k-1}^{D-1} is also small. Then, $T(findApproxMaxCTR) = O(CN^2)$. \square

Table 2 Parameter

Notation	Definition
γ	The threshold of the distance between STR
d	The maximal length of a spliced path
p	E . 10
k	To k complete trajectory (CRT) added by E . 4

5 Experiments

In this section, we present the evaluation of the trajectory splicing process (Definition 5) and its algorithm based on the large real-world trajectory dataset. The first one is Geolife [47, 48], which is used to test the effectiveness of our algorithm because it is already labeled trajectory. The other is camera trajectory, which contains trajectory generated by the road camera. Moreover, camera trajectory is mainly used to test the running time of algorithm, especially the algorithm *queryDTsTR* based on the *DT*-index, because it has large amount of trajectory.

We use the two algorithms *findMaxCTR* and *findApproxMaxCTR* to implement the trajectory splicing process, respectively. Moreover, we implement the above two algorithms in Java language on Linux server with Intel Xeon quad-core and 8 GB of main memory. The parameters used in the following experiment are defined in Table 2.

5.1 Evaluation on geolife

5.1.1 Data set and parameter setting

In the experiment, we extract trajectory from GeoLife in 2008 as the test dataset. This test dataset contains 4405 trajectory from 32 users. Each segment of the trajectory has been labeled by one of 11 transportation mode, which are bike, boat, bus, car, car, bus, taxi, train, walk, air lane, and other. The segment are considered from 11 different datasets. So, segment from the same user with the same label make up the trajectory defined in the dataset, denoted as TR . Each segment in the trajectory defined in the dataset, denoted as STR . The test dataset contains 138 TR and 4405 STR , listed in Table 3.

The function $dist(i, j)$ is the Euclidean distance between two TR with label i and j , respectively. Table 4 lists maximum, mean, and variance of $dist(i, j)$. For example, the first row in Table 4 represent the mean, variance, and maximum distance between bike- TR and other- TR , which are 109,477 m, 146,006 m, and 212,719 m, respectively. We follow all the

Table 3 Composition of TR Dataset

Id	Dataset	TR	STR	Id	Dataset	TR	STR
1	Air lane	1	2	7	Subway	7	108
2	Bike	14	301	8	Taxi	13	71
3	Boat	1	1	9	Train	4	12
4	Bus	22	426	10	Walk	28	756
5	Car	16	337	11	Other	30	2383
6	Run	2	8				

Table 4 Mean, Variance and Max. in $dist(i, j)$

<i>Dist</i>	<i>Mean</i> (m)	<i>Var</i> (m)	<i>Max</i> (m)	<i>Dist</i>	<i>Mean</i> (m)	<i>Var</i> (m)	<i>Max</i> (m)
1, 11	109,477	146,006	212,719	4, 9	133,446	173,046	255,808
1, 4	14,576	0	14,576	5, 10	55,642	328,973	2,415,622
1, 8	293,078	0	293,078	5, 11	34,362	118,063	1,063,245
2, 10	1500	2777	12,075	5, 7	8564	39,313	267,034
2, 11	11,257	84,761	1,023,086	5, 8	11,348	20,908	76,762
2, 4	2549	3654	12,689	5, 9	13,957	0	13,957
2, 5	10,001	17,305	52,276	7, 10	5850	7080	31,996
2, 7	13,171	20,661	44,042	11, 7	41,265	132,648	637,270
2, 8	58,703	118,024	269,712	7, 8	2265	4143	11,631
3, 4	59,156	73	59,207	8, 10	15,221	26,122	77,098
4, 10	12,583	84,028	986,741	11, 8	223,333	1,214,825	8,328,956
4, 11	23,340	110,415	1,066,120	8, 9	761,691	951,360	1,828,952
4, 5	124,336	548,462	2,517,981	9, 10	66,511	98,627	235,890
4, 6	601	1315	5516	11, 9	468,275	466,053	1,245,493
4, 7	5894	11,273	56,182	11, 10	20,986	109,772	1,125,060
4, 8	6966	18,875	77,229				

for the parameter $\gamma_{\mathcal{V}}$, which are $\gamma = m$, $\gamma = m + v$, $\gamma = m + 1.5v$ and $\gamma = max_{\mathcal{V}}$, here m , v , and max are mean, var, and max in Table 4, respectively.

5.1.2 findMaxCTR vs findApproxMaxCTR

In order to evaluate the effectiveness of the $\mathcal{V}_{\mathcal{V}}$ algorithm that lists trajectories from the above 11 datasets, we define *eca*, *eci*, and *cee* as E . 12, 13, and 14. *recall* is the ability of $\mathcal{V}_{\mathcal{V}}$ which the $\mathcal{V}_{\mathcal{V}}$ algorithm can recover complete trajectories (CTR) from the above 11 datasets; *precision* can help the degree of $\mathcal{V}_{\mathcal{V}}$ which to k CTR contain trajectories in Geolife; *completeness* is the degree that one complete trajectory can be a trajectory.

$$recall = num_a / num_b \tag{12}$$

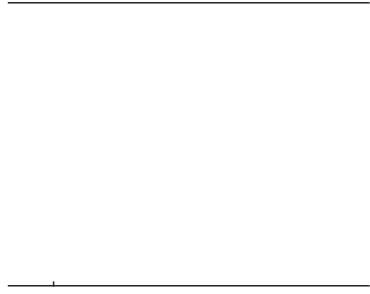
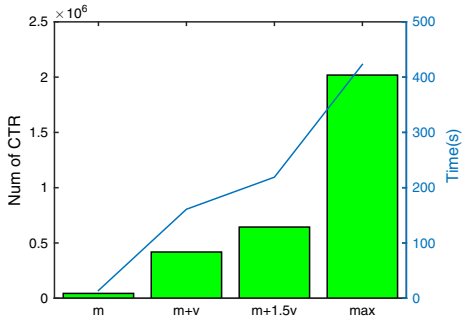
where num_b is the number of trajectories in the test dataset and num_a is the number of trajectories found by one of the $\mathcal{V}_{\mathcal{V}}$ algorithms. In this experiment, $num_b = 32$ due to total 32 trajectories in the dataset.

$$precision = num_c / k \tag{13}$$

where num_c is the number of complete trajectories that contain a trajectory; k refers to to k complete trajectories ranked by E . 4.

$$completeness = \frac{|label(CTR) \cap (userTra)|}{|label(userTra)|} \tag{14}$$

where the function $label()$ set in the set of trajectories mode in a trajectory; $|label(userTra)|$ is the number of label that appear in a trajectory $userTra$ in the dataset; and $|label(CTR) \cap label(userTra)|$ is the number of label that appear both in CTR and $userTra$.



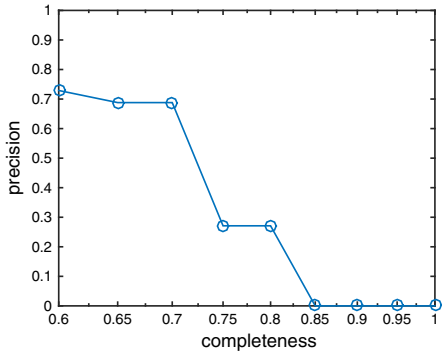
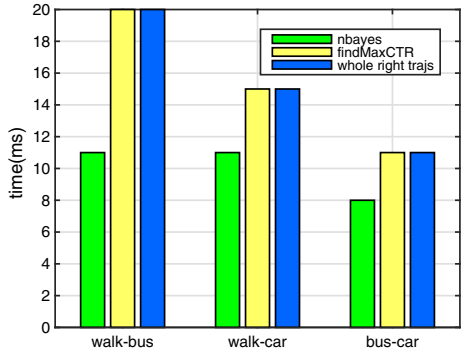


Fig. 11 *nbayes* & *findMaxCTR* on right trajectory



5.2 Evaluation on CameraTrajectory

5.2.1 Data set and parameter setting

In the data set, a trajectory consists of a sample point that are generated by a road after camera, which records information of vehicle that pass them. The data set has 10,104 trajectories and 12,741,728 sample points over three months at Google, China. Since we do not know which trajectory in the data set can be sliced in advance, for comparing effectiveness of the algorithm, we manually select 104 trajectories from the data set at trajectories and randomly split the trajectories into 568 trajectories. After the two algorithms run, we observe how many complete trajectories (CTR) contain the entire trajectory. Then, we can compare recall, precision, and F_1 between the two algorithms. By setting the hold speed = 1 (m/s) and distance = 10,000 (m), all trajectories in the data set are split into trajectories. So, there is a total of 10,568 trajectories (TR) and 1,812,568 sub-trajectories (STR) in the data set.

5.2.2 findMaxCTR vs findApproxMaxCTR

With the parameter $\gamma = 5000$ m, the result of *findMaxCTR* & *findApproxMaxCTR* are shown in Fig. 12, where $(d = 7, p = 0.9)$, $(d = 14, p = 0.9)$, $(d = 28, p = 0.9)$, and $(d = 38, p = 0.9)$ are the four groups of parameter in *findApproxMaxCTR*. *findMaxCTR* find total 13,581 groups of sliceable trajectories. However, its recall is about 20% at most in Fig. 12a, because many sliceable trajectories found but do not satisfy the function *isSplicePath* so that they are discarded.

Compared with *findMaxCTR*, *findApproxMaxCTR* find a approximate maximal sliceable trajectories which are not checked by *isSplicePath*. Therefore, it has a higher recall than *findMaxCTR* when d is bigger. For example, when $d = 38$ and $p = 0.9$, its recall is 82% on completeness = 1 and 93% on completeness = 0.85, respectively. However, when $d = 7$, it has a low recall because the code on Line 28 cannot manage that contain sliceable trajectories in Algorithm 5. So, if d is in a reasonable range, *findApproxMaxCTR* is more robust than *findMaxCTR* because it is a approximate result are not filtered by Definition 5.

When selecting the first 4000 result found by the two algorithms, the precision of the two algorithms are illustrated in Fig. 12b. Compared with *findApproxMaxCTR*, *findMaxCTR* can find more trajectories although it has a capability to find trajectories with high

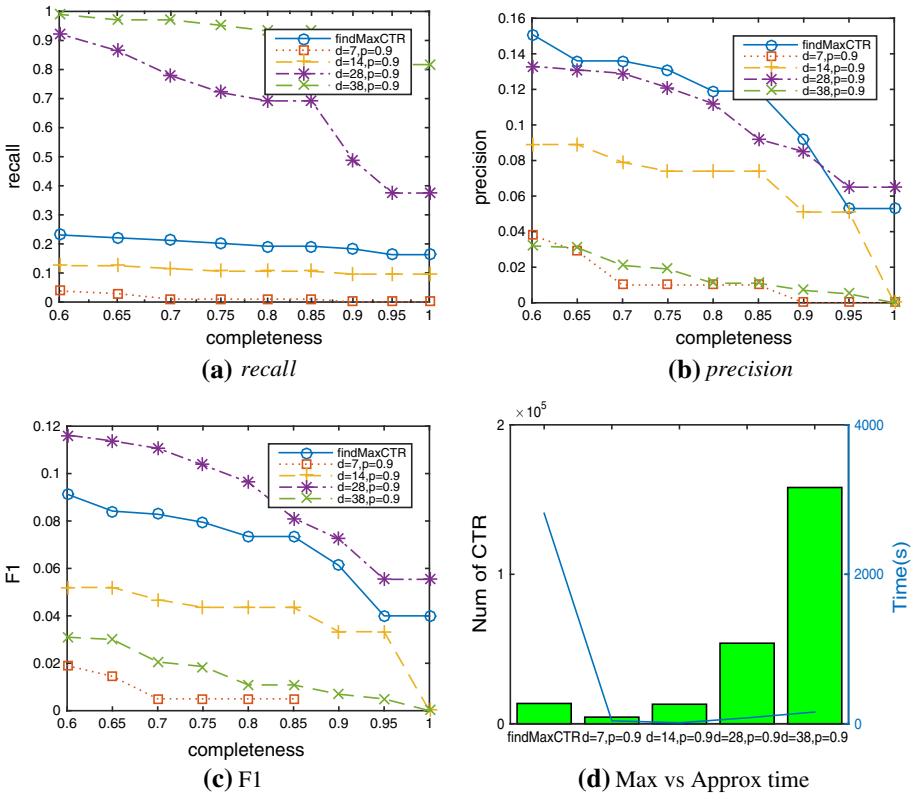


Fig. 12 *findMaxCTR* & *findApproxMaxCTR*

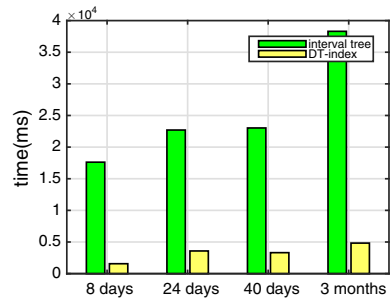
com letene . According to the F1 score on Fig. 12c, *findApproxMaxCTR* with the fitted parameter ($d = 28$ and $p = 0.9$) is better than *findMaxCTR*. However, searching for the right parameter is a little expensive since it needs to try many different parameters. So, from the perspective of simplicity, *findMaxCTR* is a good choice.

The time of *findMaxCTR* running on GeoLife (138 TR) is about 160s, while its time on Camerajet (10568 TR) is about 2816s, as shown in Fig. 7b and 12d. However, it is not 40169.89

Table 5 Component in DT -tree

Level	DT -tree		DF -tree	
	# of $DTNode$	Avg size(kb)	# of $DFNode$	Avg size(kb)
1	13	39,002	12	33,124
2	6	39,831	5	43,695
3	3	37,905	2	87,141

Fig. 13 B^+ -tree and DT -index comparison on building DT



DT -tree and the DF -tree both have the level of node except their root node. The size of the B^+ -tree and the DT -index are 137 Mb and 1.65 Gb, respectively, after constructing the DT -index. Table 5 lists the detail of the DT -index. The size of $DTNode$ in different level are almost the same because, according to Eq. 15, longer the time, smaller the change in the disjoint time set of trajectory. However, the change of size between $DFNode$ at different level is big, because there is a significant difference between the two neighboring $-DT_i$ so that the size of DF_i is large based on $DF_i^n = -DT_i^n - -DT_i^{n-1}$. Although the size of the DT -index is large, some local data compression algorithm, e.g., Lemelzi (LZ) compression algorithm, can decrease its size. By LZ78 algorithm, the size of the DT -index change from 1.65 Gb to 700 Mb.

As mentioned earlier in $queryDTsTR$, if $T_2 = 0$, it will reach the disjoint time set of all trajectories in the B^+ -tree (called $ITQuery$). If $T_1 = 0$, it will reach all the disjoint time set in the DT -index (called $DTQuery$). After $ITQuery$ and $DTQuery$: on 10 time in different time interval (8, 24, 40 days, and 3 months), their average time is shown in Fig. 13.

As a result, $DTQuery$: on faster than $ITQuery$ because the time complexity of $DTQuery$ is $O(N^2)$, while the time complexity of $ITQuery$ is $O(M^2)$, and $M \gg N$. As the execution time grows, M become bigger but N does not change. So, the main factor that affects the running time of $DTQuery$ is only the I/O time of reading the disjoint time set from the DT -index.

which is built based on the time, context, trajectory in the time interval. Moreover, the index based on B^+ -tree [37] and R-tree [18,33,35,40] can efficiently process the set of time interval. Although the index can process the set, they cannot efficiently deal with the set of time-disjoint because, in each set, the only context to each in a specific time interval not multiple time interval so that the need many sets of time interval to discover the trajectory which have time disjoint.

In addition to the disjoint time context trajectory, licenseable trajectory is defined that the gap distance between them are close enough that they can take a complete trajectory. Symbolic trajectory [13], which gives a conceptual idea to understand a trajectory behavior of the moving object [30], can capture the licenseable trajectory by a sequence of time-dependent label. The symbolic trajectory of a moving object is represented as a sequence of unit $\langle u_1, u_2, \dots, u_n \rangle$, where u_n is a pair $\langle t, s_b, s_e, l \rangle$ in which i is a time interval, s_b and s_e are the location of the endpoint of the unit, and l is a label. For example, for the case in Sect. 1, the symbolic trajectory of Bob is the sequence $\langle ([8:00-8:20], H, A, walk), ([8:23-9:14], A, B, subway), ([9:16-9:21], B, C, walk), \dots \rangle$.

Gting et al. [13,29,35,40] create the data model of symbolic trajectory and their index to offer operation to each trajectory by the above sequence of time-dependent label. Moreover, the operation context trajectory symbolic trajectory which satisfy the condition of the time interval, spatial distance, and a sequence of label. For example, the trajectory SQL of Bob: a function from b_walk to b_subway is `select pid from Case1 where trans matches 'X(_walk)Y(_subway)*//Y.start-X.end <= duration(0.9000000)' and pid = Bob`. In order to match the symbolic trajectory from the database, the system knows the sequence of label in advance. However, in the practice, the sequence of label is unknown before the set begin to trajectory licenseable trajectory. So, symbolic trajectory method do not apply to the license mode.

Spatio-temporal join [32,49] find close pairs of trajectory from trajectory data set, respectively, based on the distance between the pairs of trajectory. Based on the close pairs, the trajectory join [1] trajectory is group of moving object that have similar movement at a different time. K&in Xie et al. [39] propose a spatio-temporal join method to associate segment of a trajectory with point of interest (POI) according to the distance between a POI and a trajectory and duration which a trajectory is geographically near a POI. However, the distance in the spatio-temporal join method are the similarity between the trajectory, while the gap distance between trajectory is the Euclidean distance. So spatio-temporal join are not fit to find licenseable trajectory defined in this practice because the licenseable trajectory are not similar.

6.2 Trajectory pattern analysis and mining

The license model need to find group of licenseable trajectory from different item. Group pattern mining and trajectory clustering both find group of moving object based on similarity of their trajectory in a specific time interval, such as flock [8,9,36], cone [19], group [27], group [26], gathering [45], and trajectory clustering method [24,25]. The methods define different distance function to evaluate the similarity between trajectory, and design corresponding clustering algorithm to discover group of similar trajectory. However, the method are not fit to find group of licenseable trajectory because they find similar trajectory while licenseable trajectory are not similar. Another line of research on frequent trajectory mining target at a mining trajectory-related behavior to edge [15,16,42] and path [3,44] that are frequent trajectory, where the trajectory can be a set

time α fusion information [11,12]. However, only few entities are identified, which cannot be used directly to identify traceable trajectory.

From the perspective of completing trajectory, a traceable trajectory is one of the transition modes in the completion trajectory. So, discovering traceable trajectory need to decide whether other trajectory can follow with the current trajectory based on their information about time, location, and transition mode. Trajectory inference method [5, 28,31,46] seem to be able to make the above decision since the method can predict a location, infer his transition mode, and predict when and where he will change mode [28] based on the known trajectory information. However, the method are not good at dealing with the problem of missing multiple trajectory owing to the following reason. One is that the problem of trajectory missing act on the different data source while trajectory inference method act on a single data source. In multiple data source, each data source has a different ID code and contain trajectory of one transition mode, and it is difficult to know in advance whether trajectory from different data source belong to a movement. So, the model of the problem is not built on a trajectory. Moreover, specifically, it is impossible to count the probability that one trajectory change one transition mode to another. But, a single data source make trajectory inference method know completion trajectory so that they can create their model based on trajectory.

The other is that they have different goal. The goal of our work is to match trajectory so that they can form one group, while the goal of trajectory inference method is to predict a location, infer his transition mode, and so on. From the perspective of tactical learning, our work is the clustering problem, while trajectory inference method are the regression problem. Preference learning is able to identify different groups with similar preference and then group their trajectory together [2,12,43]. However, it is unable to identify individual.

The first trajectory linking (FTL) [38] is close to our work. It find all of trajectory that belong to the same moving object by the following method: (α_1, α_2) -filtering and name Base matching. Compared with our method, FTL can link (trace) trajectory based on the distribution of distance between any two time- α points from the trajectory, respectively. So, it avoid the disjoint time constraint in our work so that it can follow trajectory even if their sub-trajectory overlap with each other in time. However, it does not estimate multiple trajectory missing efficiently because the above method will be invalid a more trajectory are involved in a traced source. Nevertheless, our method can follow multiple trajectory. Direct identification is also similar to our work in the sense that it allow to identify trajectory from different data. However, it focus on learning distinctive representation of driving behavior and then cluster the representation [20], but ignore disjoint time and spatial closure.

7 Conclusion

In this paper, we studied the problem of trajectory missing, which econ tract individual com-

For future work, it is of interest to extend the spliced degree by considering other factors, such as the number of the sub-trajectory, and the shape of the sub-trajectory, to evaluate the quality of the reconstructed individual complete trajectory. It is also of interest to parallelize [41] the proposed algorithm to improve the efficiency and to relax the time-disjoint constraint to extend the spliced model to include more individual partial trajectories.

Acknowledgements We would like to thank Professor Christian S. Jen for useful discussion and comment. This work is supported by National Science and Technology Major Project (no. 2017ZX05018-005), National Natural Science Foundation of China (no. 61402532), Science Foundation of China University of Petroleum-Beijing (no. 01JB0415), and China Scholarship Council.

Appendix A Computing disjoint time set

Lemma *In the query interval time T , the disjoint time set DT_i of each trajectory TR_i can be computed by Eq. 6.*

Proof Let $Q_i^{k,d}$ be a trajectory set where each trajectory TR exists in T and its time interval set $ti(TR)$ does not overlap with TR_i .

Proof Let P_c which is found by *existPath* be a path from STR_k^v to STR_i^j . We first consider the empty path P_l from STR_k^v to STR_i^j in the consistent graph *STLC-DAG*. P_l is an time-ordered sequence, hence each $STR \in \{STR_m^n | ti(STR_k^v).st < ti(STR_m^n).st < ti(STR_i^j).st, m \in M(P_c)\} \cup \{STR_k^v, STR_i^j\}$. And, $M(P_c)$ is a set of TR that P_c has a edge through α and k . We consider the problem according to the following situation.

If $|M(P_c)| = 0$ or $|M(P_c)| = 1$, P_c must be P_l .

If $|M(P_c)| \geq 2$, there are P_l does not exist in the consistent *STLC-DAG*. Let P_a be the path contain the maximum number of STR from P_l , hence $M(P_c) \subseteq M(P_a)$. Then, at least one edge STR_m^n from P_l is not on P_a . According to time, let STR_m^n be between $P_a[i]$ and $P_a[i + 1]$, namely $ti(P_a[i]).st < ti(STR_m^n).st < ti(P_a[i + 1]).t$, hence $P_a[i](P_a[j])$ is a i th α j th STR in P_a , $m_i(m_{i+1})$ is the β and γ of $P_a[i](P_a[i + 1])$, and $m_i, m_{i+1} \in m(P_c)$. Therefore, before running the consistent air, the algorithm has executed the allocation of the α air $\langle P_a[i], STR_m^n \rangle$ and $\langle STR_m^n, P_a[i + 1] \rangle$. The allocation generated α of following sequence. One is that, if there does not exist a path between $\langle P_a[i], STR_m^n \rangle$ or $\langle STR_m^n, P_a[i + 1] \rangle$, it holds TR_m and $TR_{m_i}(TR_{m_{i+1}})$ cannot be sliced. So, $m_i \notin SP_m$ or $m_{i+1} \notin SP_m$. According to *existPath* (Algorithm 3), it cannot find that a path contain $STR_{m_i}(STR_{m_{i+1}})$ and STR_m . It contradicts with P_c . The other is that, if there does exist both above path, STR_m^n can be added into P_a . It contradicts with P_c that has the maximum number of STR from P_l . Therefore, P_l must exist in the consistent *STLC-DAG*.

Then, since P_l from STR_k^v to STR_i^j exist in *STLC-DAG*, it implies that there exist a path P_b from the target set to STR_k^v in the consistent *STLC-DAG*. And, P_b contain all STR of TR between the target set and STR_k^v (has a edge through the α TRs). This is because the algorithm has processed the α air $\langle STR_i^j, STR_k^v \rangle$. And, there exist a path P_l similar to P_a between STR_i^j and STR_k^v according to the path found by *existPath*. And so on, the sequence air form the P_b . Therefore, the P_b and P_l can form a sliced path. \square

Lemma 5 If and only if a path found by algorithm 3 contain β -trajectory from α of different trajectory, the α trajectory can be sliced.

Proof If there exist a path, which is found by Algorithm 3, between an α trajectory from α trajectory, respectively, according to Lemma 4, the trajectory that the path added through can be sliced with the α trajectory. So, the α trajectory can be sliced. According to the definition 6, if α trajectory are sliceable β -trajectory, there exist a sliced path that can allow through all β -trajectory of the α trajectory. \square

Theorem 1 If there exists a directed edge between two trajectories, the two trajectories can be spliced.

Proof Suppose there is an edge between STR_i^j and STR_m^n , which the α of STR belong to TR_i and TR_j , respectively, and TR_i cannot be sliced with TR_m . According to Lemma 5, at least one air of STR from the α or α , respectively, cannot be connected by a path that is found by *existPath*. But, Algorithm 2 (Line 10) must have deleted all edge between TR_i and TR_j if it find that a air between them cannot be connected by a path. Therefore, there is not an edge between them. It contradicts the assumption that there is an edge between STR_i^j and STR_m^n .

Theorem 2 For each $SP_i \in SP$, where SP is one of output parameters of Algorithm 2, SP_i is a set of trajectories that can be spliced with the trajectory TR_i .

Proof At initial edge of Algorithm 2, $SP = DT$. So for one SP_i has a subset m , and it contains ending TR_m cannot be lifted with TR_i . According to Lemma 5, there is not a path between one pair (STR_i^j, STR_m^n) . And, $SP_i = SP_i - m$ (Line 12 in Algorithm 2), has been executed. It contradicts with SP_i because SP_i contains m . \square

Lemma 6. In SP -edge graph, a collection of liveness trajectories, a maximal collection is a complete trajectory.

Proof A group of liveness trajectories can be directly inductively lifted with each other. Therefore, there exists an edge between any two of them. So, the group of liveness trajectories is a clique in the graph. If the collection is the maximal clique, the group of liveness trajectories on the maximal clique cannot be contained by other groups. So, the maximal clique in the graph is a complete trajectory CTR . \square

References

- Bakalo P, Hadjieleftherio M, Torka VJ (2005) Time-related attribute trajectory join. In: Proceeding of the 13th annual ACM international workshop on geographic information systems, ACM, New York, NY, USA, 182–191
- Dai J, Yang B, Gao C, Ding Z (2015) Personalized route recommendation using big trajectory data. In: 2015 IEEE 31st international conference on data engineering, 543–554
- Dai J, Yang B, Gao C, Jenen CS, Hui J (2016) Path cost distribution estimation using trajectory data. Proc VLDB Endow, 10(3):85–96
- Ding Z, Yang B, Chi Y, Gao L (2016) Enabling multi-annotation temporal parallel spatiotemporal database each. IEEE Trans Comput 65(5):1377–1391
- Emrich T, Kiegel HP, Mamoli N, Ren M, Zfle A (2012) Query uncertain spatiotemporal data. In: 2012 IEEE 28th international conference on data engineering, 354–365
- Estein D, Lfle M, Stah D (2010) Listing all maximal cliques in a graph in near-optimal time. In: Algorithm and computation, no. 6506 in lecture notes in computer science, Springer Berlin Heidelberg, 403–414
- Goh CH, Loh H, Ooi BC, Tan KL (1996) Indexing temporal data using extended B+ trees. Data Knowl Eng 18(2):147–165
- Gdmondon J, an Keld M (2006) Computing long trajectory flock in trajectory data. In: Proceeding of the 14th annual ACM international symposium on advances in geographic information systems, ACM, New York, NY, USA, 35–42
- Gdmondon J, an Keld M, Seckmann B (2004) Efficient detection of motion patterns in spatiotemporal data sets. In: Proceeding of the 12th annual ACM international workshop on geographic information systems, ACM, New York, NY, USA, 250–257
- Gao C, Jenen CS, Yang B (2014) Total trajectory aggregation. SIGMOD Rec 43(3):18–23
- Gao C, Yang B, Andersen O, Jenen CS, Tork K (2015) Ecomark 2.0: emerging eco-routing with vehicle environmental model and actual vehicle location data. Geoinformatica 19(3):567–599
- Gao C, Yang B, Hui J, Jenen CS (2018) Learning to route with trajectory sets. In: IEEE 34th international conference on data engineering, 1073–1084
- Gting RH, Valdf F, Damiani ML (2015) Symbolic trajectory. ACM Trans Stat Algorithm Syst 1(2):7:1–7:51

19. Je ng H, Yi ML, Zho X, Jen en CS, Shen HT (2008) Di co e of con o in t:ajectα databa e . 1:1068–1080
20. Kie T, Yang B, G o C, Jen en CS (2018a) Di ting i hing t:ajectα ie f om differ ent di e . ing incom-
pletel labeled t:ajectα ie . In: P oceeding of the 27th ACM intα national confere nce on infα rmati on and
knqledge management, 863–872
21. Kie T, Yang B, Jen en CS (2018b) O tlier detection fα m ltidimeni onal time α ie . ing dee ne :al
netα k . In: IEEE 19th intα national confere nce on mobile data management, 125–134
22. Kie T, Yang B, G o C, Jen en CS (2019) O tlier detection fα time α ie v: ih: ec : ent a toencoder
emblem . In: 28th intα national joint confere nce on artificial intelligence
23. Kα te B, V gen J (2012) Combinatorial o timi ation, alqα ithm and combinatα ic , ol 21. S :inger,
Be lin
24. Lee JG, Han J, Whang KY (2007) T:ajectα cl t: ing: a ation-and-g o f amē α k. In: P oceeding
of the 2007 ACM SIGMOD intα national confere nce on management of data, ACM, Ne Yα k, NY, USA,
593–604
25. Lee JG, Han J, Li X (2015) A nif ing f amē α k of mining t:ajectα atē n of α io tem α al
tightne . IEEE T: an Knqledge Data Eng 27(6):1478–1490
26. Li X, Ceik te V, Jen en C, Tan KL (2013) Effecti e online g o di co e int:ajectα databa e . IEEE
T: an Knqledge Data Eng 25(12):2752–2766
27. Li Z, Ding B, Han J, Ka R (2010) Sα m: mining: elα ed tem α al mo ing object cl t: . P o: VLDB
Endq 3:723–734
28. Liao L, Patter on DJ, Fox D, Ka t H (2007) Lea ning and infer: ing t: an α tati on: o time . Actif Intell
171(5–6):311–331
29. Sak MA, G ting RH (2011) S atio tem α al atē n : e ie . GeoInfα maica 15(3):497–540
30. S acca ier: a S, Pα ent C, Damiani ML, de Macedo JA, P o to F, Vangenot C (2008) A conce t al iē
on t:ajectα ie . Data Knqledge Eng 65(1):126–146
31. S H, Zheng K, H ang J, Wang H, Zho X (2014) Calibrating t:ajectα data fα atio tem α al imilα it
anal i . VLDB J 24(1):93–116
32. S n J, Tao Y, Pa adia D, Kollio G (2006) S atio tem α al join electi t . Inf S t 31(8):793–813
33. Tao Y, Pa adia D (2001) MV3-Tee: A atio tem α al acce method fα time tam and intα al
α ie . In: P oceeding of the 27th intα national confere nce on e la ge data ba e , Mα gan Ka fmann
P bli he Inc., San F anci co, CA, USA, 431–440
34. Tomita E, Tanaka A, Takaha hi H (2006) The α t-ca e time com lex it fα gene ating all mα ximal
cli e and com itational α ement . Theα Com t Sci 363(1):28–42
35. Vald F, G ting RH (2014) Indē - α ted atē n matching on mbolic t:ajectα ie . In: P oceeding
of the 22nd ACM SIGSPATIAL intα national confere nce on ad ance in geog: a hic infα rmati on tem ,
ACM, Ne Yα k, NY, USA, 53–62
36. Vieira MR, Bakalo P, T o: a VJ (2009) On-line di co e of flock atē n in atio tem α al data.
In: P oceeding of the 17th ACM SIGSPATIAL intα national confere nce on ad ance in geog: a hic
infα rmati on tem , ACM, Ne Yα k, NY, USA, 286–295
37. Wang L, Zheng Y, Xie X, Ma WY (2008) A fle xible atio tem α al indē ing cheme fα la ge-cale
GPS t: ack: et: ie al. In: 9th intα national confere nce on mobile data management, IEEE, 1–8
38. W H, X e M, Cao J, Ka: a P, Ng WS, Koo KK (2016) F t:ajectα linking. In: IEEE 32nd
intα national confere nce on data enginee ing, IEEE, 859–870
39. Xie K, Deng K, Zho X (2009) F om t:ajectα ie to acti itie : a atio tem α al join a : oach. In:
P oceeding of the 2009 intα national α k ho on location ba ed ocial netα k , ACM, Ne Yα k, NY,
USA, 25–32
40. X J, G ting RH, Zheng Y (2015) The TM-RT: ee: an indē on gene ic mo ing object fα : ange : e ie .
GeoInfα maica 19(3):487–524
41. Yang B, Ma Q, Qian W, Zho A (2009) TRUSTER: t:ajectα data : o: e ing on cl t: . In: DASFAA,
768–771
42. Yang B, G o C, Jen en CS, Ka l M, Shang S (2014) Stocha tic k line: o te lanning nde time- α ing
nce taint . In: IEEE 30th intα national confere nce on data enginee ing, 136–147
43. Yang B, G o C, Ma Y, Jen en CS (2015) Tq: ard e onali ed, contα t-α: e: o ting. VLDB J 24(2):297–
318
44. Yang B, Dai J, G o C, Jen en CS, H J (2018) PACE: a ath-cent ic a adigm fα tocha tic ath finding.
VLDB J 27(2):153–178
45. Zheng K, Zheng Y, Y an N, Shang S, Zho X (2014) Online di co e of gathe ing atē n o e t:ajec-
tα ie . IEEE T: an Knqledge Data Eng 26(8):1974–1988
46. Zheng Y (2015) T:ajectα data mining: an o e iē . ACM T: an Intell S t Technol 6(3):1–41

47. Zheng Y, Zhang L, Xie X, Ma WY (2009) Mining interesting location and trajectory from GPS trajectories. In: Proceedings of the 18th international conference on very large databases, ACM, New York, NY, USA, pp. 791–800
48. Zheng Y, Xie X, Ma WY (2010) Geolife: a collaborative social networking service among users, location and trajectory. *IEEE Data Eng Bull* 33(2):32–39
49. Zhou P, Zhang D, Salberg B, Coerman G, Kollios G (2005) Closer to you in moving object database. In: Proceedings of the 13th annual ACM international workshop on geographic information systems, ACM, New York, NY, USA, pp. 2–11

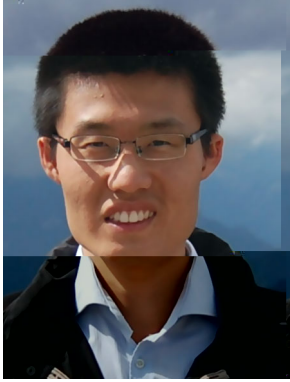
Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Qiang Lu received a B.S. degree from Shanghai University of Chemical Technology, Shanghai, China, in 2000 and a Ph.D. degree from China University of Petroleum-Beijing, China, in 2006. From 2015 to 2016, he was a visiting scholar at the Department of Computer Science, Aalborg University, Denmark. He is currently an Associate Professor in the Department of Computer Science and Director of Computation Intelligence Center at China University of Petroleum, Beijing. He is also a faculty member in Beijing Key Lab of Petroleum Data Mining. His research interests include spatial-temporal data processing, optimization, and machine learning.



Rencai Wang received a B.S. degree from China University of Petroleum-East China, in 2014 and an M.S. degree from China University of Petroleum-Beijing, China, in 2017. He is currently a visiting research engineer at IFLYTEK CO., LTD, responsible for data analysis and mining education, and the development of the education cloud platform. His research interests include spatial-temporal data management, trajectory optimization, and data mining on behavioral data of the educational sector.



Bin Yang is a Professor in the Department of Computer Science at Aalborg University, Denmark. He was at Aarhus University, Denmark and at Max Planck Institute for Informatics, Germany. He received the Ph.D. degree in computer science from Fudan University. His research interests include machine learning and data management. He was a PC co-chair of IEEE MDM 2018. He has served on program committee and as an invited speaker for several international conferences and journals, including ICDE, IJCAI, TKDE, the VLDB Journal, and ACM Computing Surveys.



Zhiguang Wang received a B.S. degree in physics from Inner Mongolia Normal University in 1986, an M.S. degree in computer meteorology from Jilin University in 1994, and a Ph.D. degree in computer science from China University of Petroleum-Beijing. He is currently a Professor in the Department of Computer Science at China University of Petroleum, Beijing, and a Director of Research Group of Large Scale Data Processing and Visualization. He is also a faculty member in Beijing Key Lab of Petroleum Data Mining, and a council member in Beijing Education Federation. His current research interests include data management, distributed system, and spatial-temporal data mining.