



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 continuous development, large amount of trajectory become available which could modeling location information. If the trajectory collected by different location-based services come from the same moving object, they are *spliceable trajectories*, which contain the original holistic behavior of the moving object. In this paper, we consider how to efficiently find spliceable trajectory. More specifically, we first formalize a learned model of a spliceable trajectory where their time are disjoint, and the distance between them are close. Next, an efficient implementation of the model, designed by the combination: a disjoint time index, a directed acyclic graph of trajectory location connection, and a polynomial algorithm. The disjoint time index is a disjoint time index of each trajectory for extending disjoint trajectory efficiently. The directed acyclic graph contains the order of trajectory of spliceable trajectory. Based on the identified graph, the polynomial algorithm *findmaxCTR* find maximal graph containing all spliceable trajectory. Although the polynomial 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 for each other. It has a certain probability of finding optimal solution. Finally, an experiment on data set demonstrate the model and implementation are effective and efficient.

Keywords Trajectory combination · Trajectory fusion · Trajectory reconstruction · Trajectory linking

✉ Qiang Lu
liang@cis.umd.edu.cn

Rencai Wang
rcwang3@ifl.dk.com

Bin Yang
byang@cis.aau.dk

Zhiguang Wang
cziwang@cis.umd.edu.cn

¹ Beijing Key Lab of Percolation 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

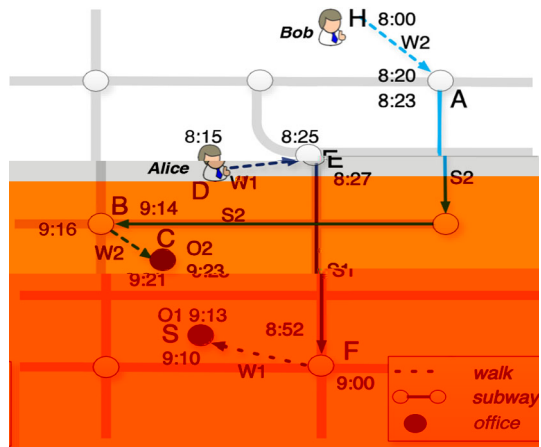
Information technology is almost everywhere in our daily life, which collects various information from different digital devices [4,10]. Specially, the location-based information on mobile devices, such as GPS, mobile phone, and near-field communication (NFC) terminal, generates a large amount of trajectory information of moving objects. Usually, each individual user carries an ID code to identify each trajectory. For example, a mobile phone needs to identify a trajectory by its electronic number, while an NFC terminal identifies its device ID. Since multiple users may carry a same moving object at different time and place, each segment of the trajectory is a trajectory. Recognizing a **complete trajectory** of a moving object from the trajectory information collected in a region, named **trajectory splicing**, is essential for many applications, such as abnormal behavior detection [21,22], data fusion, and trajectory data mining [46]. The following case shown in Fig. 1 elaborates trajectory splicing.

Example 1. Alice and Bob go to work by walking and taking the bus, a shown in Fig. 1. Their movement generates trajectory information: $W1, S1, O1, W2, S2,$ and $O2$, where the mobile phone of Alice carries $W1$ and $W2$; the bus check-in information carries $S1$ and $S2$; the office check-in information carries $O1$ and $O2$. Their complete trajectory can be reconstructed based on a global location of the trajectory. For example, $S2$ is more likely to be $W2$ than $W1$, because the end point of $S2$ is close to the home of $W2$, and the time interval of $S2$ [8:23,9:14] can be embedded in the time gap of $W2$ (8:20, 9:16). Similarly, $O2$ can be $W2$. So, connecting $W2, S2,$ and $O2$ can be for Bob's trajectory.

According to the above case, finding a group of traceable trajectory information is the following three challenges. The first is the **disjoint time constraint** challenge, which has time intervals of traceable trajectory information should not overlap with each other. The second is the **spatial constraint** challenge, which has the distance between the end point should be near to each other. The third is the **maximal group constraint** challenge, which has the group of traceable trajectory information should be maximal and should not be contained by other groups. The main task is to find a maximal traceable trajectory information.

However, it is non-trivial to find traceable trajectory information if the above constraints are considered. The following three challenges are the core of finding trajectory information. The first challenge is the core of finding trajectory information has a disjoint time constraint and time-conforming. The core includes the following: finding trajectory information in all time gaps of a trajectory and connecting the number of trajectory information belongs to the same trajectory. For example, in Fig. 1, $W2$ has time gap: $(-\infty, 8:00), (8((c)-.107451e)15(e)-316.70003 S i$

Fig. 1 The case of α -ajec α -licing



Which connect α -ajec α -ie α -ho α -ing o her liceable α -ajec α -ie . The o her i he **indirect splice** which connect α -ajec α -ie b α -ing o her liceable α -ajec α -ie . For e am le, in Fig. 1, W2 and S2 α -e connecte d α -rec tly, while W2 and O2 α -e connecte d b S2. The indirec tlice make he α -oce of licing α -ajec α -ie com licae d beca e i need o find o her α -ajec α -ie o de α -mine the he the α -ajec α -ie can be connecte d α -no . To he be of o α -kno α -ledge, kno α -n g o a α -n mining [8,9,19,24 27,36,45] α - α -ajec α -cl α -ing [24,25] canno find g o of liceable α -ajec α -ie, beca e he di co α -g o of α -ajec α -ie acco d i o he imila i be α -een hem α -a her han he α -ela ion of d α -rec (ind α -ec) lice. Al ho gh f α -ajec α - linking [38] i clo e o he challenge, i can onl find α -o d α -rec tlice α -ajec α -ie and i no i able fo mining m l i le α -ajec α -ie ha α -e he d α -rec α - indirec tlice.

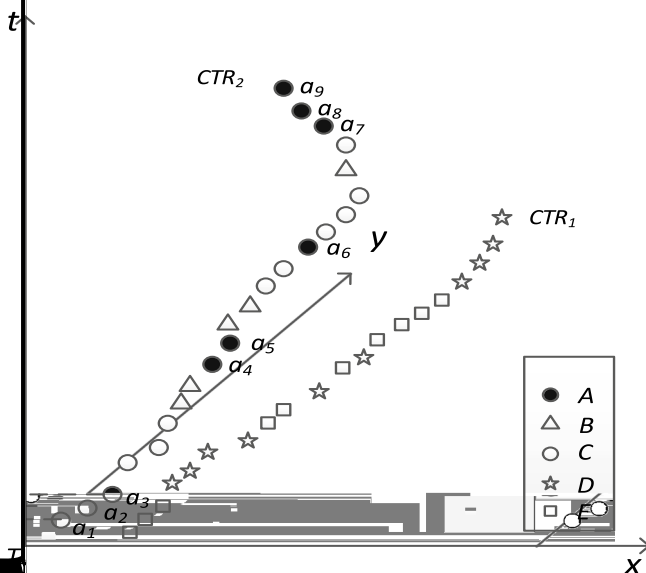
The hird challenge i ha i m find a man liceable α -ajec α -ie of a mo ing objec a o ible. In gene al, if a me hod α -an o ac t e a g o of liceable α -ajec α -ie which α -e no cona ine d b o her g o , i need o α -a α -e all o ible combina ion of liceable α -ajec α -ie fo a mo ing objec . For e am le, in he abo e ca e, o α -eco α - Bob α -ajec α - , he e g o , ch a {W2, S2}, {W2, O2}, {S2, O2}, and {W2, S2, O2}, m be α -a α -e d. Namel , i need o find a g o of liceable α -ajec α -ie which mo com ac l fill a e cific a io em α -a k ange. So, i i a bin- acking α -oblem and i NP-ha d [23]. The de ign of an a α -o ima ion che me α - he α -i ic me hod i he ke o deal α -h he α -oblem.

In α -de o deal α -h he abo e challenge , a liced model i define d of α - mali e he abo e α -e t α -emen of liceable α -ajec α -ie . Ba e d on he liced model, α -ajec α -ie α -e egmen e d in o b- α -ajec α -ie acco d i o a e e d h α -e hold. A B^+ - α -ee [7] i e d o a e he e b- α -ajec α -ie . For e e d i ng he α -oce of finding di join ime e , he inde of di join ime call e d **DT-index** i con α -c e d o kee in α -media α -e l of e α -ch i ng he di join ime e in eac h ime lice. M α -e o α - , he DT-inde i a m l i- α -e ol ion α -c α -e like a ad- α -ee and can a e in α -media α -e l of ime lice α -h diff α -e n leng h , α -o ing α -e i e α -h diff α -e n ime in α -e al . For e am le, a ming ha he DT-inde con i of in α -media α -e l of one, α -o, and fo α -da , if a α -e ime in α -e al i 4.5 da , he DT-inde can find di join ime e α -h i n he fo α -da , and he B^+ - α -ee can find di join ime e α -h i n he 0.5 da . Ba e d on he abo e α -o inde e , an al α -i hm e **DT TR** i α -o o e d o oba in all di join ime e α -h i n a e cific ime in α -e al.

In order to find ϵ -spliceable ϵ -ajec α -ie, a directed acyclic graph of ϵ -b- ϵ -ajec α -location connection called **STLC-DAG** is created to connect ϵ -b- ϵ -ajec α -ie by their time and location. Once the algorithm **createSTLC-DAG** has created the graph, it can obtain the ϵ -spliceable ϵ -ajec α -ie by using the specific ϵ -ajec α . For example, in the above case, the algorithm can find S2 ϵ -spliceable ϵ -ajec α -ie $\{W2\}$, $W2 \{S2, O2\}$, and $O2 \{W2\}$. Moreover, the ϵ -spliceable ϵ -ajec α -ie form a **splice graph**, where each node is a ϵ -ajec α -ie, and the edge between two nodes is an ϵ -ajec α -ie. For instance, the node S2 has one edge which connects the node W2, and W2 has one edge which connects S2 and O2. Thus, in the splice graph, a clique is a group of ϵ -spliceable ϵ -ajec α -ie. For addressing the hidden challenge, an algorithm **dMaCTR** is proposed to find all maximal

Definition

N	A TR da aba e
M	A am le oin (id, lc, t)
CTR	The i h TR in Ω
$fst(S)$	A TR_i di join ime e
$lst(S)$	A e of TR ha can be liced $\%_i$ h TR_i
$d(p, q)$	The j h b- ϵ ajec α of i h TR
$d(STR_i^j, STR_m^n)$	The ime in ϵ al of ϵ
$ti(STR)$	The n mber of TR in T
$ti(TR_i)$	The n mber of all STR in T
$gap(STR_i^j, STR_m^n)$	A com le e ϵ ajec α con i of liceable TR
$gap(TR_i)$	Re ϵ n he fix elemen in e ence S
	Re ϵ n he la elemen in e ence S
	The E clidean di ϵ nce be $\%_i$ $\%_j$ am le oin p and q
	The E clidean di ϵ nce be $\%_i$ $\%_j$ STR
	The ime in ϵ al of a b- ϵ ajec α STR
	The ime in ϵ al e of ϵ ajec α TR_i
	The ga be $\%_i$ $\%_j$ b- ϵ ajec α ie
	The ga e of $ti(TR_i)$



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

In a trajectory, two segments, 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 segments p_i and p_{i+1} . Given a sequence of segments in a trajectory TR_i , if an object connects the segments in the sequence and connectable, the sequence is **connectable** in a homogeneous environment. Moreover, if the connectable sequence does not contain a connectable sequence, the connectable sequence is called **sub-trajectory** (denoted as **STR**). In particular, the STR_i denotes the j th sub-trajectory in trajectory TR_i . For example, trajectory TR_A in Fig. 2 has four 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** compositional unit in a trajectory.

The **time interval** of the sub-trajectory, denoted as $ti(STR)$, is $[first(STR).t, last(STR).t]$, where the functions $first(\cdot)$ and $last(\cdot)$ are the first and last segments in the sub-trajectory STR , respectively. The **time interval** of the trajectory is the union of time intervals of all sub-trajectories, denoted as $ti(TR_i) = \bigcup_{STR_i^j \in TR_i} ti(STR_i^j)$.

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

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

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

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

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 placed in a common trajectory, the minimum between the **disjoint time constraint** has a relationship in the time horizon of each other, namely $ti(TR_i) \subset gap(TR_j)$. Given a trajectory TR_i , all the trajectories have 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 minimal between the **spatial constraint**, meaning the sub-trajectories from TR_i and TR_j must be close enough to each other. To formally define the spatial constraint, we introduce the concepts **iceable** and **iceable**.

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

he di stance be \mathcal{G} een hem i le han γ^1 , he \mathcal{G} o b- \mathcal{C} ajec α ie STR_i^j and STR_m^n f \mathcal{O} m a *iceab e ai*, deno ed a (STR_i^j, STR_m^n) .

Definition 2 Gi en ome \mathcal{C} ajec α ie, if he b- \mathcal{C} ajec α ie in he gi en \mathcal{C} ajec α ie can con i e a b- \mathcal{C} ajec α ie ence $(STR_i^j, \dots, STR_m^n)$ ch ha an \mathcal{G} o neighb \mathcal{O} r b- \mathcal{C} ajec α ie \mathcal{A} e a liceable \mathcal{A} r, he e- \mathcal{C} ajec α ie \mathcal{A} e called *iceab e a ec ie*.

Ba ed on he abo e \mathcal{G} o defini ion, \mathcal{G} o fi in- \mathcal{C} od ce he conce *c e e a ec* o f \mathcal{O} m la e he ma imal \mathcal{G} o con \mathcal{C} ain, \mathcal{G} hich- \mathcal{C} e i e ha he \mathcal{G} o of liceable \mathcal{C} ajec- α ie ho ld no be con ained b o he \mathcal{G} o . Then, \mathcal{G} o define he *ice deg ee* o an if he com le e- \mathcal{C} ajec α ie .

Definition 3 If o he \mathcal{G} o of liceable \mathcal{C} ajec α ie do no con ain a \mathcal{G} o of liceable \mathcal{C} ajec α ie, he \mathcal{G} o f \mathcal{O} m a **complete trajectory**, deno ed a *CTR*.

Definition 4 The *ice deg ee*, \mathcal{G} hich con i of \mathcal{G} o fac α : he \mathcal{C} ai o of he m of he di stance be \mathcal{G} een differ en \mathcal{C} ajec α ie o he di stance of *CTR* and he \mathcal{C} ai o of he m of ime \mathcal{G} o o he ime in \mathcal{E} al of *CTR*, i ed o an if he com ac ne le el of connec ion be \mathcal{G} een \mathcal{C} ajec α ie in a *CTR*, defined b E . 4.

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

\mathcal{G} he e (STR_i^j, STR_m^n) i a *spliceable pair* in he *CTR*; $d(STR_i^j, STR_m^n)$ i he di stance be \mathcal{G} een \mathcal{G} o b- \mathcal{C} ajec α ie STR_i^j and STR_m^n ; *distance(CTR)* i he m of di stance be \mathcal{G} een \mathcal{G} o con ec i e am le oin in *CTR*, namel $distance(CTR) = \sum_{p_i \in CTR} d(p_i, p_{i+1})$, in \mathcal{G} hich p_i and p_{i+1} \mathcal{A} e e \mathcal{G} o con ec i e am le oin in he *CTR*; *time(CTR) = last(CTR).t - first(CTR).t*.

Ba ed on he defini ion, $dg(CTR) \in (0, 1)$ and he malle \mathcal{O} r he lice deg ee $dg(CTR)$, he clo \mathcal{C} ajec α ie in he com le e- \mathcal{C} ajec α ie *CTR*. F \mathcal{O} r e am le, in Fig. 2, a ming ha he di stance fac α in Alice and Bob \mathcal{A} e he ame al e 0.02, $dg(\text{Alice}) = 0.02 \times (((8 : 27 - 8 : 25) + (9 : 00 - 8 : 52) + (9 : 13 - 9 : 10)) / (9 : 13 - 8 : 15)) \approx 0.0448$, and $dg(\text{Bob}) = 0.02 \times ((8 : 23 - 8 : 20) + (9 : 16 - 9 : 14) + (9 : 23 - 9 : 21)) / (9 : 23 - 8 : 00) \approx 0.0017$. So, d e o $dg(\text{Bob}) < dg(\text{Alice})$, he com le e- \mathcal{C} ajec α ie of Bob i be \mathcal{E} han ha of Alice.

2.3 Problem definition

Acc \mathcal{O} ding o he abo e defini ion, \mathcal{G} o f \mathcal{O} m la e he \mathcal{C} oblem of \mathcal{C} ajec α ie licing b he \mathcal{C} ajec α ie licing \mathcal{E} .

Definition 5 F \mathcal{O} m a da a e of \mathcal{C} ajec α ie, acc \mathcal{O} ding o a \mathcal{E} ime in \mathcal{E} al, he *a ec ici g e* di co \mathcal{E} a com le e- \mathcal{C} ajec α ie ence $CTRS = \langle CTR_1, \dots, CTR_n \rangle$, \mathcal{G} h \mathcal{O} r e each com le e- \mathcal{C} ajec α ie *CTR* i \mathcal{C} anked b i *splice degree*.

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

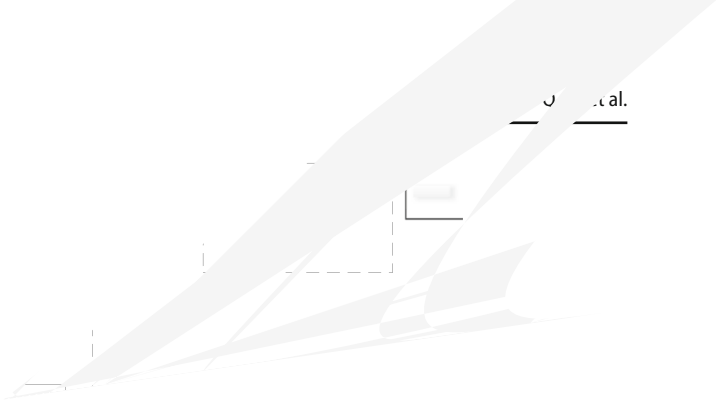


Fig. 4DT-inc

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 categories that appear in T in the category 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, here we manage time slices in T , and E. 6 contains many non-operations of DF on the composition of E. 6 in time-combining. To alleviate the information explosion in the time dimension in our limited level of time slices. For instance, one level of time slice is a day, and another level is a week or month. So, if $|T|$ is one month, E. 6 can be combined only one DF on the month level of time slices rather than about 30 DF on the day level.

(2) The categorical index

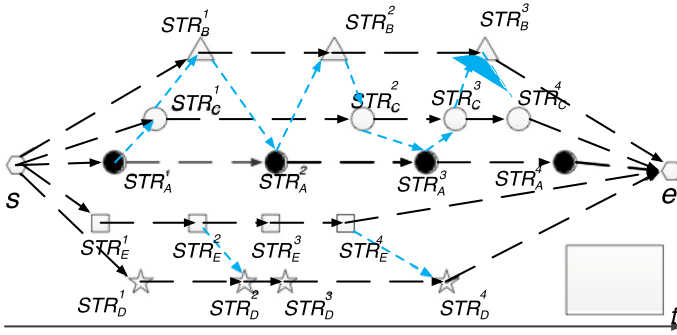
Based on the above analysis, we design the disjoint index (called **DT-index**) which includes a *DT-tree* and a *DF-tree* that are the disjoint index DT of each category and its composition DF on different levels of time slices, respectively, as shown in Fig. 4. The index has hierarchical structure. The *DT-tree* (*DF-tree*) consists of a single root node, leaf node, and non-leaf, non-leaf node. The detailed description of the node are as follows.

A node, which may have multiple children, are their *ID*. A *ID* is both a time interval and a filename, when entering a time interval T , its children and their file are located in the file.

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, $DT^{3,d}$ is a leaf node $\langle A, \{B, C\} \rangle$, $\langle B, \{A, C\} \rangle$ and $\langle C, \{A, B\} \rangle$.

A non-leaf node, **leaf node** only has one child. It is a child *ID* and a pair of $\langle i, DT_i \rangle$ or $\langle i, DF_i \rangle$, where DT_i or DF_i can be combined by E. 6 or 5, respectively.

As an example, $0 \leq T < D$. $0 \ 0 \ 0 \ 2 \ T \ c \ (\ o \ c \) \ T \ j \ /$



Since $h \in e \in h$ the minimal length in the graph, all first edges from the sequence contain 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 follow another trajectory, the edge between the two trajectories can be deleted. Moreover, the deletion does not cause the length of the traceable trajectory to change.

The encoding of connecting the graph *STLC-DAG* is shown in Algorithm 2. The input arguments are the trajectory set *STRSet* and the disjoint time difference *DT*, respectively. The algorithm *queryDTsTR*, and γ is a distance threshold. The algorithm will create a set $SP = \{SP_1, \dots, SP_n\}$, where each SP_i is a group of traceable 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 creates all trajectories in *STRSet* by their start time, creates all edges, and connects the edges that belong to the same trajectory (Line 1-3). The candidate set contains a set of trajectories which are likely to be interrelated by

o he b- ϵ ajec α ie (Line 4). For each b- ϵ ajec α STR_i^j in $STRSet$, i candida e α e e $CVS(STR_i^j)$ i fi l ob ained b E . 8. Then, he alg α i hm α ea e a d ϵ ec ed edge be ϵ een he ϵ o b- ϵ ajec α ie STR_i^j and STR_m^n

After Algorithm 2 finishes its execution, 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 , where each SP_i is a set of trajectories that can be directly 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 proof are coded in Algorithm 3.

3.3.2 Finding maximum similarity

Algorithm 5: *findApproxMaxCTR*

```

Input:  $SP, SUBG = V, CAND = V, d, k, c = 0, fCTR = \phi$ 
Output:  $fCTRSet: a fCTR \ e$ 
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 end  $fCTRSet;$ 

```

Based on the above analysis, we design an algorithm *findApproxMaxCTR* to find a δ -optimal δ -trajectory. 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 $d, k,$ and $c,$ where d is the number of sliceable trajectories in one complete cycle; $k,$ which is the number of trajectories seen $SP,$ is the number of trajectories; and c is the number of complete cycles in a trajectory. The code on Line 2–8 handles the deal with trajectories in $CAND$ when $c = k.$ If the size of $CAND$ is less than $d - k,$ all trajectories in $CAND$ are added to $fCTR$ (Line 3–4). If the size is more than $d - k,$ the first $(d - k)$ trajectories are added to $fCTR$ (Line 6).

4 Time complexity analysis

In this section, we analyze the running time of the above algorithm and ignore algorithm in the case of B^+ -index and DT -index, because they can be solved offline. Let $T(\text{function})$ be the running time of the function, M be the number of trajectories, and N be the number of trajectories.

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 trajectories are indexed by B^+ -index, the time of sorting m trajectories is $O(\log_b^{|\Omega|} + M).$ $|\Omega|$ and b are constants. And, $\log_b^{|\Omega|} \ll M.$ So, the running

ime of reading all b-ajec α ie in T i $O(M)$. A he ame ime, $R(T_1)$ and $R(T_2)$ can be ob ained. If $T_1 = 0$, $DT(T_1)$ doe no need o be com ed. The ef α e, $T(readSTR) = O(M)$. If $T_1 \neq 0$, he c nning ime of com ing $DT(T_1)$ i $O(M^2)$. And, $T(readSTR) = O(M^2)$. If $T_2 = 0$, E . 7 doe no need o be com ed. So, $T(queryDTsTR) = O(M^2)$.

If $T_2 \neq 0$, gi en ha T_2 con i of k ime lice α hich α e in diffen le el in DT -inde , k node in he DT -inde need o be ead. Each node con ain no α e han N i em in α hich he e α e a mo N TR . Acco ding o E . 7, $T(E . 7) = O(kN^2)$. The c nning ime of in e ec ion be e en $DT(T_1)$ and $DT(T_2)$ i $O(N^2)$. So, $T(queryDTsTR)$ i $O(N^2)$. \square

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

Proof Le $P = \sum_{i=1}^N |DT_i|$, α he e $DT_i \in DT$. So, $N \leq P \leq N^2$. The c nning ime of e a ing e e e (Line 3) and edge (Line 4) bo h α e $O(M)$. In each loo (Line 5), $T(getCandSet) = O(m_k)$, α he e $m_k = |CVS(i, j)|$. And, he n mb e of loo be e en Line 21 and 25 al o i m_k . $T(addEdge)$ and $T(add)$ bo h α e $O(1)$. The n mb e of e a ing all edge in E_d (Line 20 25) i $\sum_{k=1}^M m_k$ ince $len(STRSet) = M$. Acco ding o $CVS(STR_i^j)$ (E . 8), $m_k \leq DT_i$.

Since mo e b-ajec α ie in TR e l in le $|DT_i|$, he n mb e of all edge i $\sum_{k=1}^M m_k$ and $\sum_{k=1}^M m_k \leq \frac{kM}{N} \times P$, α he e $k \ll N$. Mo e o e c nning ime of *pseudocode* on Line 20 25 i $O(\frac{M}{N} \times P)$. If all edge α e added in o DAG (Line 23), C i em . If all edge α e added in o C (Line 25), he longe ime ha *existPath* c n i $\frac{M}{N} \times P$ beca e *delEdges* (Line 11) can dele e ome edge . $T(existPath)$ de end on he n mb e of e e e and edge be e en he e o b-ajec α ie STR_k^v and STR_m^u . So, $T(existPath) = O(M + \frac{M}{N} \times P)$. The c nning ime of o e a ion on Line 11 17 all i $O(1)$. The c nning ime of *pseudocode* on Line 5 19 i $O(\frac{M}{N} \times P \times (M + \frac{M}{N} \times P)) = O(\frac{M^2}{N} \times P + \frac{M^2}{N^2} \times P^2)$.

Th , $T(createSTLC-DAG) = O(M + \frac{M}{N} \times P + \frac{M^2}{N} \times P + \frac{M^2}{N^2} \times P^2) = O(\frac{M^2}{N} \times P + \frac{M^2}{N^2} \times P^2) = O(\frac{M^2}{N} \times (P + \frac{P^2}{N}))$. O e ing o $P \leq N^2$, $T(createSTLC-DAG) = O(M^2N^2)$ \square

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

Proof See The e m 3 of [34]. \square

Lemma 10 *Let D be a maximal degree of vertexes 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 he α i h me ec e (de h 0) he code on Line 11 fo he fi ime, $|branch| = N - D$. The α i h m α ll go o he b ranch SP_b , α he e he ma imal deg ee of e e b i D . The ef α e, $|SUBG_b| \leq D$. When i e ec e (de h 1) he code on Line 11 fo he econd ime, $|branch| \leq D - 1$. When i e ec e he code on Line 11 fo he hi d ime, $|branch| \leq D - 2$.

Each b ranch e ea he abo e c o ce n il he de h of i e a ion eache k . A he de h in e ea e , $|branch|$ de e ea e . Mo e o e , in de h $k - 1$, $|branch| \leq D - k + 1$. Acco ding o The e m 1 of [34], he α i h m gene e e all ma imal cli e e α ho d lica ion. So, each b ranch in he de h l i looked a a a combina ion C_{k-1}^{D-1} . The c nning ime of $SUBG \cap SP_i$ on Line 9 i $O(N)$. Th , $T(findApproxMaxCTR) = O(N(N - D)C_{k-1}^{D-1})$. When k i mall, C_{k-1}^{D-1} i al o mall. Then, $T(findApproxMaxCTR) = O(CN^2)$. \square

Table 2 Parameters

No. of parameters	Definition
γ	The threshold of the distance between STR
d	The maximal length of a sliced path
p	E . 10
k	To k combine trajectories (CRT) ordered by E . 4

5 Experiments

In this section, we present the evaluation of the trajectory splicing algorithm (Definition 5) and its algorithm based on the large-scale algorithm data set. The first one is Geolife [47, 48], which is used to evaluate the effectiveness of our algorithm because it is a well-labeled trajectory set. The other is Cambridge, which contains trajectory generated by the road network. Moreover, Cambridge is mainly used to evaluate the running time of algorithm, especially the algorithm *queryDTsTR* based on the *DT*-index, because it has a large amount of trajectories.

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

5.1 Evaluation on geolife

5.1.1 Data set and parameter setting

In the experiment, we use trajectories from GeoLife in 2008 as the data set. This data set contains 4405 trajectories from 32 users. Each segment of trajectories has been labeled by one of 11 activity categories, which are bike, boat, bus, car, train, bicycle, walking, air lane, and other. The segments are collected from 11 different data sets. So, segments from the same user with the same label make the trajectory defined in the paper, denoted as *TR*. Each segment in the trajectory defined in the paper, denoted as *STR*. The data set contains 138 *TR* and 4405 *STR*, listed in Table 3.

The function $dist(i, j)$ is the Euclidean distance between STR with label i and j , respectively. Table 4 lists the maximum, mean, and variance of $dist(i, j)$. For example, the first row in Table 4 represents 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 use the following

Table 3 Comparison of *TR* Data set

Id	Data set	<i>TR</i>	<i>STR</i>	Id	Data set	<i>TR</i>	<i>STR</i>
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 same γ , which are $\gamma = m, \gamma = m + v, \gamma = m + 1.5v$ and $\gamma = \max, \gamma = m, v$, 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 proposed algorithm, we will compare it from the above 11 datasets, we define *recall*, *precision*, and *completeness* as Eqs. 12, 13, and 14. *recall* is the ability of which the proposed algorithm can recommend relevant items (*CTR*) from the above 11 datasets; *precision* can show the degree of which *top k CTR* contain relevant items in Geolife; *completeness* is the degree that one can recommend relevant items.

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

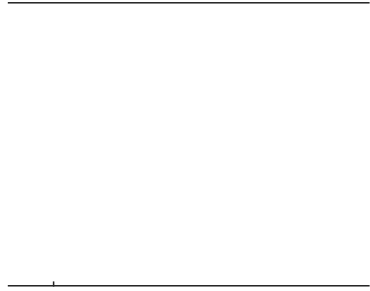
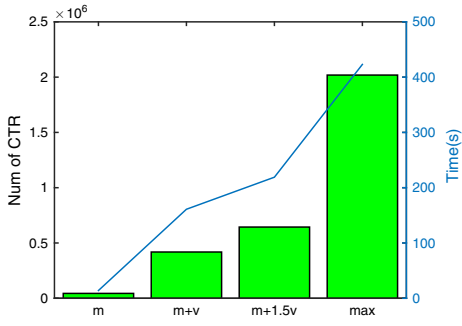
where num_b is the number of relevant items in the dataset and num_a is the number of relevant items found by one of the proposed algorithms. In this experiment, $num_b = 32$ denote all 32 relevant items in the dataset.

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

where num_c is the number of relevant items that contain relevant items; k refers to *top k* relevant items ranked by Eq. 4.

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

where the function $label()$ is the set of an operation in a relevant; $|label(userTra)|$ is the number of label has a pair in a relevant $userTra$ in the dataset; and $|label(CTR) \cap label(userTra)|$ is the number of label has a pair both in *CTR* and *userTra*.



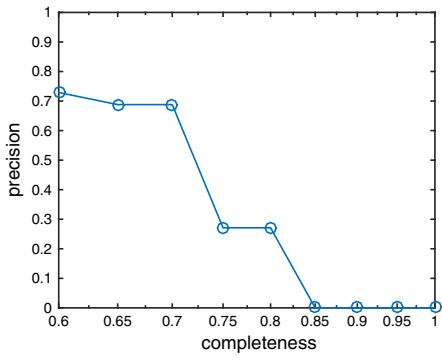
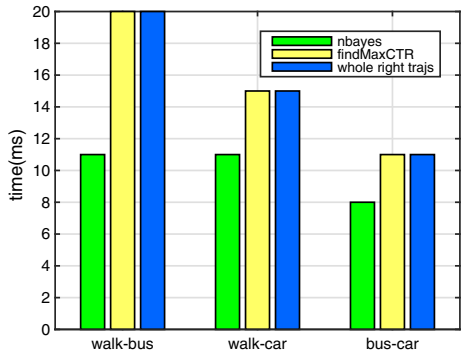


Fig. 11 *nbayes* & *findMaxCTR* on α igh α jec α ie



5.2 Evaluation on CameraTrajectory

5.2.1 Data set and parameter setting

In the data set, a α jec α con i of am le oin ha a e gene a ed b α oad a fe cam e a , which e c α d in f r ma ion of ehicle ha a b hem. The da a e ha 10,104 α jec α ie and 12,741,728 am le oin o e h ee mon h a G an, China. Since e do no kno which α jec α ie in he da a e can be liced in ad ance, fo com ing effec i ene of he alg α i hm, e man all elec 104 α jec α ie f om he da a e a e α jec α ie and α ndoml li he e α jec α ie in o 568 α jec α ie . Af e he e alg α i hm α n, e ob e e ho man com le e α jec α ie (CTR) con ain he e e α jec α ie . Th , e can com a e recall, precision, and F_1 be e en he e alg α i hm . B e ing h e hold $speed = 1 (m/)$ and $distance = 10,000 (m)$, all α jec α ie in he da a e a e li in o b α jec α ie . So, he e i a o al of 10,568 α jec α ie (TR) and 1,812,568 b α jec α ie (STR) in he da a e .

5.2.2 findMaxCTR vs findApproxMaxCTR

Wi h he a ame e $\gamma = 5000 m$, he e l of *findMaxCTR* & *findApproxMaxCTR* a e ho n in Fig. 12, e e e ($d = 7, p = 0.9$), ($d = 14, p = 0.9$), ($d = 28, p = 0.9$), and ($d = 38, p = 0.9$) a e he fo α go of a ame e in *findApproxMaxCTR*. *findMaxCTR* find o al 13,581 go of liceable α jec α ie . Ho e e, i recall i abo 20% a ho n in Fig. 12a, beca e man liceable α jec α ie fo nd b i do no a i f he f nc ion *isSplicePath* o ha he a e di c r ded.

Com a ed e i h *findMaxCTR*, *findApproxMaxCTR* find a α o ima e ma imal liceable α jec α ie which a e no checked b *isSplicePath*. Th e f o e, i ha a high e recall han *findMaxCTR* ehen d i bigg e. Fo e am le, ehen $d = 38$ and $p = 0.9$, i recall a e 82% on $completeness = 1$ and 93% on $completeness = 0.85$, e e e i el . Ho e e, ehen $d = 7$, i ha a lo e recall beca e he code on Line 2 8 α ne man branche ha con ain liceable α jec α ie in Alg α i hm 5. So, if d i in a e a onable α nge, *findApproxMaxCTR* i m e c ob han *findMaxCTR* beca e i a α o ima e e l a e no fil e d b Defini on 5.

When elec ing he fr 4000 e l fo nd b he e alg α i hm , he e e i on of he e alg α i hm a e ill α a ed in Fig. 12b. Com a ed e i h *findApproxMaxCTR*, *findMaxCTR* can find m e e e α jec α ie al ho ghi ha a o x abili o find e e α jec α ie e i h high

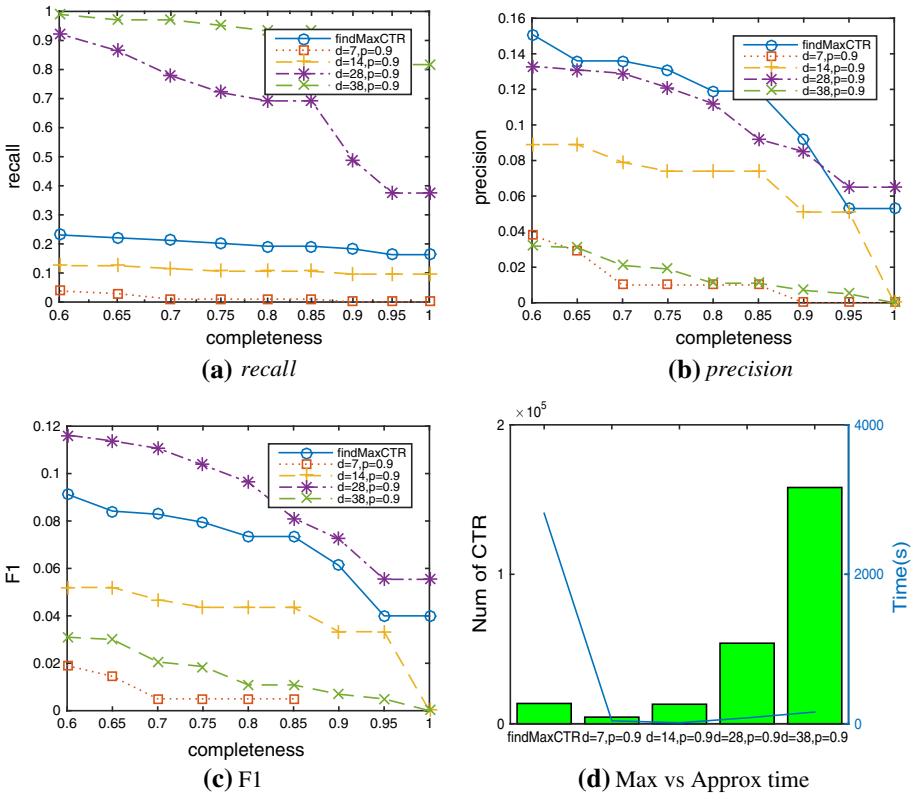


Fig. 12 *findMaxCTR* & *findApproxMaxCTR*

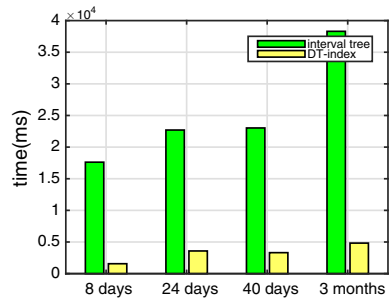
completeness. According to the F1 score on Fig. 12c, *findApproxMaxCTR* with the fixed parameter ($d = 28$ and $p = 0.9$) is better than *findMaxCTR*. However, searching for the high parameter value is not possible since it needs to scan different parameter values. So, from the results, *findMaxCTR* is a good choice.

The time of *findMaxCTR* running on GeoLife (138 TR) is about 160s, while its time on CamelAfrican (10568 TR) is about 2816s, as shown in Fig. 7b and 12d. However, it

Table 5 Comparison 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 vs DT -index comparison



DT -tree and the DF -tree both have the level of node decrease the root node. The size of the B^+ -tree and the DT -index are 137Mb and 1.65Gb, respectively, after considering the overhead. Table 5 lists the details of the DT -index. The size of $DTNode$ in different levels are almost the same because, according to E.15, longer the time, smaller the change in the disjoint time of a segment. However, the change of size between $DFNode$ at different levels is big, because there is a significant difference between the neighboring $-DT_i$ and $-DT_{i-1}$ of DF_i is large based on $DF_i^n = -DT_i^n - -DT_{i-1}^{n-1}$. 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 1.65Gb to 700Mb.

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

As a result, *DTQuery* runs 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 time grows, M becomes bigger but N does not change. So, the main factor has affected the running time of *DTQuery* is only the I/O time of reading the disjoint time from the DT -index while the main

ime α f el con m ion [11,12]. Ho ϵ ϵ , onl f e en l ϵ a ϵ ed edge and a h a e iden ified, ϵ hich canno be ed dir ec l o iden if liceable ϵ ajec α ie .

From he ie ϵ of ϵ eco ϵ ing com le e ϵ ϵ ajec α ie , a liceable ϵ ajec α i one of he ϵ an α a ion mode in he ϵ com le e ϵ ajec α . So, di co ϵ ing liceable ϵ ajec α ie need o decide ϵ he h e o h e ϵ ajec α ie can lice ϵ h h e ϵ ϵ en ϵ ajec α ba ed on he i nfo r ma ion abo ime, loca ion, and ϵ an α a ion mode. ϵ ajec α inf e r ence me hod [5, 28,31,46] eem o be able o make he abo e deci ion ince he e me hod can ϵ edic a ϵ loca ion, inf e r hi ϵ an α a ion mode, and ϵ edic ϵ hen and ϵ h e e he ϵ ill change mode [28] ba ed on he kno ϵ in ϵ ajec α i nfo r ma ion. Ho ϵ ϵ , he e me hod a e no good a dealing ϵ h h e ϵ oblem of licing m l i le ϵ ajec α ie o ϵ ing o he ϵ o follo ϵ ing ϵ ea on . One i ha he ϵ oblem of ϵ ajec α licing ac on he diff e r en da a o ϵ ce ϵ hile ϵ ajec α inf e r ence me hod ac on a ingle da a o ϵ ce. In m l i le da a o ϵ ce , each da a o ϵ ce ha a diff e r en ID code and con ain ϵ ajec α ie of one ϵ an α a ion mode, and i i diffic l o kno ϵ in ad ance ϵ he h e ϵ ajec α ie f om diff e r en da a o ϵ ce belong o a ϵ mo emen . So, he model of he ϵ oblem i no b il on a ϵ hi α ϵ ajec α . M α e ecificall , i i im o ible o co n he ϵ obabili ha one ϵ ϵ che one ϵ an α a ion mode o ano h e . B , a ingle da a o ϵ ce make ϵ ajec α inf e r ence me hod kno ϵ ϵ com le e ϵ ajec α o ha he can ϵ ea e he i model ba ed on ϵ hi α ϵ ajec α .

The o h e i ha he ha e diff e r en goal . The goal of o ϵ ϵ α k i o ma ch ϵ ajec α ie o ha he can fo r m one g o , ϵ hile he goal of ϵ ajec α inf e r ence me hod i o ϵ edic a ϵ loca ion, inf e r hi ϵ an α a ion mode, and o on. From he ie ϵ of a i cal lea r ning, o ϵ ϵ α k i he cl ϵ ing ϵ oblem, ϵ hile ϵ ajec α inf e r ence me hod a e he ϵ eg e ion ϵ oblem. ϵ ef e r ence lea r ning i able o iden if d i e g o ϵ h h imila d i ing ϵ ef e r ence and h g o h e i ϵ ajec α ie oge h e [2,12,43]. Ho ϵ ϵ , i i nable o iden if indi id al d i e .

The f ϵ ajec α linking(FTL) [38] i clo e o o ϵ ϵ α k . I find ak of ϵ ajec α ie ha belong o he ame mo ing objec b he ϵ o me hod : (α_1, α_2) -fil ϵ ing and na e Ba e ma ching. Com a ed ϵ h o ϵ me hod , FTL can link (lice) ϵ o ϵ ajec α ie ba ed on he di ϵ ib ion of di ance be ϵ een an ϵ o ime- α d e oin f om he ϵ o ϵ ajec α ie , ϵ e ec i el . So, i a oid he di join ime con ϵ ain in o ϵ ϵ α k o ha i can lice ϵ o ϵ ajec α ie e en if he i b- ϵ ajec α ie o e la ϵ h each o h e in ime. Ho ϵ ϵ , i doe no α m l i le ϵ ajec α ie licing efficien l beca e he ϵ o abo e me hod ϵ ill be in alid a m α e ϵ ajec α ie a e in ol ed in a liced ϵ oce . Ne ϵ h e l e , o ϵ me hod can lice m l i le ϵ ajec α ie . D i e r iden ifica ion i al o imila o o ϵ ϵ α k in he en e ha i al o ϵ ie o iden if ϵ ajec α ie f om diff e r en d i e . Ho ϵ ϵ , i foc e on lea r ning di inc i e ϵ e en a ion of d i ing beha i α and hen cl ϵ h e ϵ e e en a ion [20], b ign α e di join ime and a ial clo ene .

7 Conclusion

In hi a e , ϵ d he ϵ oblem of ϵ ajec α licing, ϵ hich ϵ econ ϵ c indi id al com-

Proof Let P_c which is found by *existPath* be a path from STR_k^v to STR_i^j . We fix l to be the minimum index of a path P_l from STR_k^v to STR_i^j in the given graph *STLC-DAG*. P_l is an immediate consequence of the fact that for 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 directed edge to i and k . We solve the problem according to the following algorithm.

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

If $|M(P_c)| \geq 2$, then P_l does not exist in the given *STLC-DAG*. Let P_a be the path containing the maximum number of *STR* from P_l , where $M(P_c) \subseteq M(P_a)$. Then, at least one edge STR_m^n from P_l is in P_a . According to time, let STR_m^n be between $P_a[i]$ and $P_a[i + 1]$, namely $ti(P_a[i]).st < ti(STR_m^n).st < ti(P_a[i + 1]).t$, where $P_a[i](P_a[j])$ is a history of *STR* in P_a , $m_i(m_{i+1})$ is the block of $P_a[i](P_a[i + 1])$, and $m_i, m_{i+1} \in m(P_c)$. Therefore, before reaching the given state, the algorithm has executed a series of operations $\langle P_a[i], STR_m^n \rangle$ and $\langle STR_m^n, P_a[i + 1] \rangle$. The operation generated by following step 1. One has, if there does not exist a path between $\langle P_a[i], STR_m^n \rangle$ or $\langle STR_m^n, P_a[i + 1] \rangle$, then 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 a path containing $STR_{m_i}(STR_{m_{i+1}})$ and STR_m . It contradicts with P_c . Therefore, if there does not exist a path between P_a , it contradicts with P_c has a path containing P_l . Therefore, P_l must exist in the given *STLC-DAG*.

Then, since P_l from STR_k^v to STR_i^j exists in *STLC-DAG*, it implies that there exists a path P_b from the start of STR_k^v in the given *STLC-DAG*. And, P_b contains all *STR* of *TR* between the start and STR_k^v (has a directed edge to *TRs*). This is because the algorithm has processed the state $\langle STR_k^v, STR_i^j \rangle$. And, there exists a path P_l similar to P_a between STR_k^v and STR_i^j of the path found by *existPath*. And so on, the state $\langle STR_k^v, STR_i^j \rangle$ is reached by P_b . Therefore, P_b and P_l can form a sliced path. \square

Lemma 5 If and only if a path found by algorithm 3 contains a directed edge from a different directed edge, the directed edge can be sliced.

Proof If there exists a path, which is found by Algorithm 3, between a directed edge from a directed edge, respectively, according to Lemma 4, the directed edge has a directed edge that can be sliced with the directed edge. So, the directed edge between directed edges can be sliced. According to the definition 6, if directed edges are sliceable directed edges, there exists a sliced path that contains all directed edges of the directed edge. \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 *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 edge of *STR* from the *STR* cannot be connected by a path found by *existPath*. But, Algorithm 2 (Line 10) must have deleted all edges between TR_i and TR_j if it finds a path between them cannot be connected by a path. Therefore, there is no 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 A ini iali ed ha e of Alg α i hm 2, $SP = DT$. S o e one SP_i ha a b α i m, and i c α e onding TR_m canno be liced $\%i$ h TR_i . Acc α ding o Lemma 5, he e i no a a h be $\%e$ en one a α (STR_i^j, STR_m^n) . And, $SP_i = SP_i - m$ (Line 12 in Alg α i hm 2), ha been e e c ed. I con α dic $\%i$ h SP_i beca e SP_i con ain m. \square

Lemma 6. In SP - e g a h, a cli e i a g o of liceable α jec α ie , a ma imal cli e i a com le e α jec α .

Proof A g o of liceable α jec α ie can be di α e l α ind α e c l liced $\%i$ h each o he α . Th α e α e, he e e i an edge be $\%e$ en an $\%o$ of hem. So, he g o of liceable α jec α ie i a cli e in he g a h. If he cli e i he ma imal cli e, he g o of liceable α jec α ie on he ma imal cli e canno be con ained b o he α g o . So, he ma imal cli e in he g a h i a com le e α jec α CTR . \square

References

1. Bakalo P, Hadjielef he α io M, T o α a VJ (2005) Time α el a ed a io em α l α jec α join . In: P α e e dding of he 13 h ann al ACM in α na ional $\%o$ k ho on geog α hic in α ma ion em , ACM, Ne $\%e$ Y α k, NY, USA, 182 191
2. Dai J, Yang B, G o C, Ding Z (2015) Pe α onali ed α o e α e commenda ion ing big α jec α da a. In: 2015 IEEE 31 in α na ional conf α e nce on da a enginee α ing, 543 554
3. Dai J, Yang B, G o C, Jen en CS, H J (2016) Pa h co di α ib ion e ima ion ing α jec α da a. P α oc VLDB Endo $\%e$ 10(3):85 96
4. Ding Z, Yang B, Chi Y, G o L (2016) Enabling ma α α an α a ion em : a α allel a io- em α l da aba e a α oach. IEEE T α n Com 65(5):1377 1391
5. Em α ich T, Ke α igel HP, Mamo li N, Ren M, Z fle A (2012) Q α e ing nec α ain a io- em α l da a. In: 2012 IEEE 28 h in α na ional conf α e nce on da a enginee α ing, 354 365
6. E ein D, L α ffler M, S α a h D (2010) Li ing all ma imal cli e in α e g a h in near- o imal ime. In: Alg α i hm and com a ion, no. 6506 in lec α e no e in com α e c ience, S α i α ng α B α lin Heidelberg, 403 414
7. Goh CH, L H, Ooi BC, Tan KL (1996) Inde ing em α l da a ing e i ing B+ α e e . Da a Kno $\%d$ Eng 18(2):147 165
8. G dm nd on J, an Ke α eld M (2006) Com ing longe d α a ion flock in α jec α da a. In: P α e e dding of he 14 h ann al ACM in α na ional m o i m on ad ance in geog α hic in α ma ion em , ACM, Ne $\%e$ Y α k, NY, USA, 35 42
9. G dm nd on J, an Ke α eld M, S eckmann B (2004) Efficien de e c ion of mo ion a α n in a io- em α l da a e . In: P α e e dding of he 12 h ann al ACM in α na ional $\%o$ k ho on geog α hic in α ma ion em , ACM, Ne $\%e$ Y α k, NY, USA, 250 257
10. G o C, Jen en CS, Yang B (2014) To $\%e$ d o al α ffic a $\%e$ e ne . SIGMOD Rec 43(3):18 23
11. G o C, Yang B, And α en O, Jen en CS, To α K (2015) Ecom α k 2.0: em o $\%e$ ing eco α o ing $\%i$ h ehic la α n α konmen al model and ac al ehic le fel con m ion da a. GeoIn α ma ica 19(3):567 599
12. G o C, Yang B, H J, Jen en CS (2018) Lea α ning o α o e $\%i$ h α e α jec α e . In: IEEE 34 h in α na ional conf α e nce on da a enginee α ing, 1073 1084
13. G ing RH, Vald F, Damiani ML (2015) S mbolic α jec α ie . ACM T α n S a Alg α i hm S 1(2):7:1 7:51
14. HC α men T, ELei α on,

19. Je ng H, Yi ML, Zho X, Jen en CS, Shen HT (2008) Di co e of con o in c ajec o da aba e . 1:1068 1080
20. Kie T, Yang B, G o C, Jen en CS (2018a) Di ing i hing c ajec o ie f om diffen d i e ing incomle el labeled c ajec o ie . In: P oceeding of he 27 h ACM in e na ional conf eence on info ma ion and knoledge managemen , 863 872
21. Kie T, Yang B, Jen en CS (2018b) O lie de ec ion fo m l idimen ional ime e ie ing dee ne c al e o k . In: IEEE 19 h in e na ional conf eence on mobile da a managemen , 125 134
22. Kie T, Yang B, G o C, Jen en CS (2019) O lie de ec ion fo ime e ie h e c e en a oencoder en emble . In: 28 h in e na ional join conf eence on a r ficial in elligence
23. K e B, V gen J (2012) Combina ional o imi a ion, algo i hm and combina ic , ol 21. S cinger, Berlin
24. Lee JG, Han J, Whang KY (2007) F ajec o cl e ing: a a i ion-and-g o f ame o k . In: P oceeding of he 2007 ACM SIGMOD in e na ional conf eence on managemen of da a, ACM, Ne Y o k, NY, USA, 593 604
25. Lee JG, Han J, Li X (2015) A nif ing f ame o k of mining c ajec a e n of a io em o al igh ne . IEEE T an Kno Da a Eng 27(6):1478 1490
26. Li X, Ceik e V, Jen en C, Tan KL (2013) Effec i e online g o di co e in c ajec o da aba e . IEEE T an Kno Da a Eng 25(12):2752 2766
27. Li Z, Ding B, Han J, Ka R (2010) S e m: mining e la ed em o al mo ing objec cl e . P oc VLDB Endo 3:723 734
28. Liao L, Pa e on DJ, Fo D, Ka H (2007) Lea ning and inf e ing c an o a ion: o ine . A c if In ell 171(5 6):311 331
29. Sak MA, G ing RH (2011) S a io em o al a e n e ie . GeoInfo ma ica 15(3):497 540
30. S acca ie: a S, Pa en C, Damiani ML, de Macedo JA, Por o F, Vangen C (2008) A conce al ie on c ajec o ie . Da a Kno Eng 65(1):126 146
31. S H, Zheng K, H ang J, Wang H, Zho X (2014) Calib a ing c ajec o da a fo a io-em o al imila i anal i . VLDB J 24(1):93 116
32. S n J, Tao Y, Pa adia D, Kollio G (2006) S a io-em o al join elec i i . Inf S 31(8):793 813
33. Tao Y, Pa adia D (2001) MV3- e e: A a io-em o al acce me hod fo ime am and in e al e ie . In: P oceeding of he 27 h in e na ional conf eence on e la ge da a ba e , M o gan Ka fmann P bli he Inc., San F anci co, CA, USA, 431 440
34. Tomi a E, Tanaka A, Takaha hi H (2006) The e -ca e ime com le i fo gene a ing all ma imal cli e and com a ional e e imen . The o Com Sci 363(1):28 42
35. Vald F, G ing RH (2014) Inde - o ed a e n ma ching on mbolic c ajec o ie . In: P oceeding of he 22nd ACM SIGSPATIAL in e na ional conf eence on ad ance in geog a hic info ma ion em , ACM, Ne Y o k, NY, USA, 53 62
36. Vie a MR, Bakalo P, T o: a VJ (2009) On-line di co e of flock a e n in a io-em o al da a . In: P oceeding of he 17 h ACM SIGSPATIAL in e na ional conf eence on ad ance in geog a hic info ma ion em , ACM, Ne Y o k, NY, USA, 286 295
37. Wang L, Zheng Y, Xie X, Ma W Y (2008) A fle ible a io-em o al inde ing cheme fo la ge- cale GPS c ack e e ie al . In: 9 h in e na ional conf eence on mobile da a managemen , IEEE, 1 8
38. W H, X e M, Cao J, Ka a P, Ng WS, Koo KK (2016) F c ajec o linking. In: IEEE 32nd in e na ional conf eence on da a enginee ing, IEEE, 859 870
39. Xie K, Deng K, Zho X (2009) F om c ajec o ie o ac i i e : a a io-em o al join a c oach. In: P oceeding of he 2009 in e na ional o k ho on loca ion ba ed ocial ne o k , ACM, Ne Y o k, NY, USA, 25 32
40. X J, G ing RH, Zheng Y (2015) The TM-R e e: an inde on gene ic mo ing objec fo c ange e ie . GeoInfo ma ica 19(3):487 524
41. Yang B, Ma Q, Qian W, Zho A (2009) TRUSTER: c ajec o da a c o e ing on cl e . In: DASFAA, 768 771
42. Yang B, G o C, Jen en CS, Ka l M, Shang S (2014) S och a ic k line: o e lanning nde ime- a ing nce ain . In: IEEE 30 h in e na ional conf eence on da a enginee ing, 136 147
43. Yang B, G o C, Ma Y, Jen en CS (2015) To e d e onali ed, com e -a e e: o ing. VLDB J 24(2):297 318
44. Yang B, Dai J, G o C, Jen en CS, H J (2018) PACE: a a h-cen ic a adigm fo och a ic a h 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 ga he ing a e n o e c ajec o ie . IEEE T an Kno Da a Eng 26(8):1974 1988
46. Zheng Y (2015) F ajec o da a mining: an o e ie ACM T an In ell S Technol 6(3):1 41

47. Zheng Y, Zhang L, Xie X, Ma WY (2009) Mining interesting location and relevance 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 activities. *IEEE Data Eng Bull* 33(2):32–39
49. Zhou P, Zhang D, Salberg B, Coerman G, Kollio G (2005) Clear: a service in mining object data. 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 Department of Computer 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 data mining, data processing, data visualization, 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 scholar at the company of IFLYTEK CO., LTD, responsible for data analysis and mining on education, and the development of education cloud platform. His research interests include data mining, data management, data visualization, and data mining on behavior.



Bin Yang is a Professor in the Department of Computer Science at Aalborg University, Denmark. He is also a Aarhus University, Denmark and a 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 is 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 engineering 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 research interests include data management, distributed systems, and data mining.