



电子科技大学
University of Electronic Science and Technology of China



Outlier Detection of Trajectory

Yuanliang Zhang



Data Mining Lab, Big Data Research Center, UESTC
Email: junmshao@uestc.edu.cn
<http://staff.uestc.edu.cn/shaojunming>

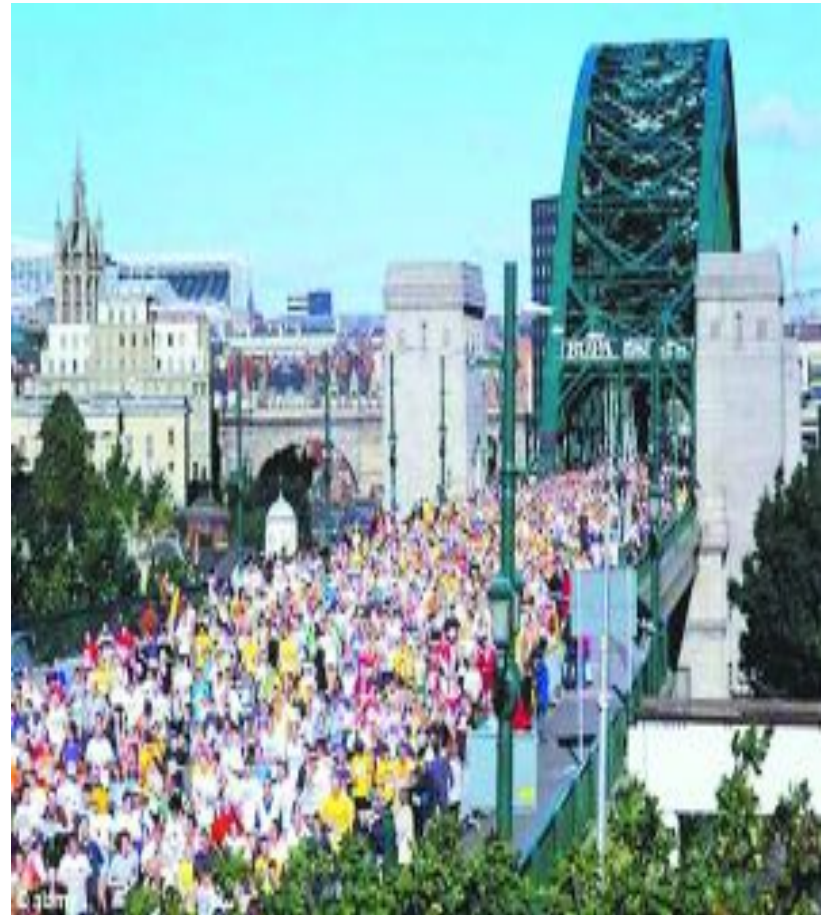
Outline

- Introduction
- Related Algorithms
- Conclusion & Discussion

Part 1 Introduction

Introduction

- Marathon



Introduction

- Two significant concepts
 - Trajectory
 - Trajectory Outlier(Anomaly)

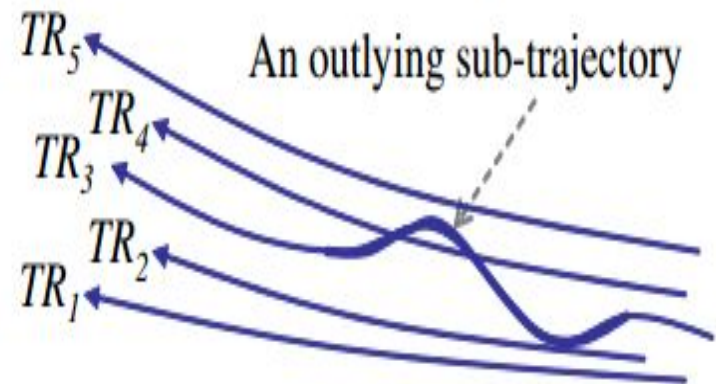
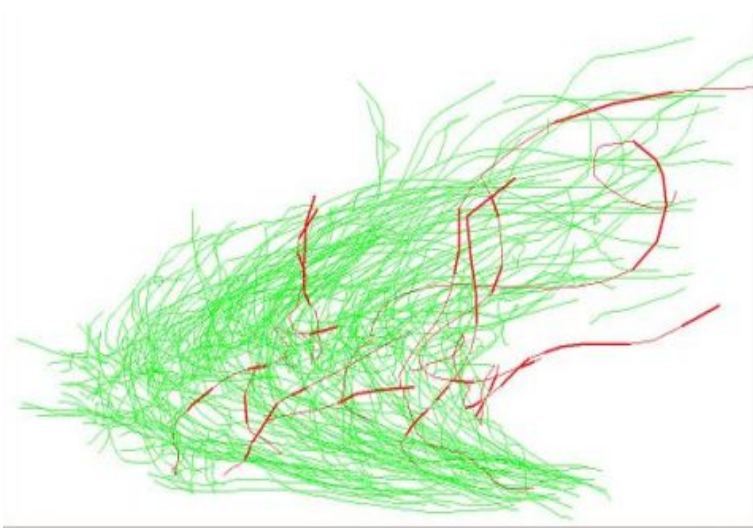
Introduction

- Trajectory:
 - A trace generated by a moving object
 - Moving objects: human, vehicles, animals, natural phenomena
 - GPS data
 - Location-based data : social network, credit card records
 - Represented by a series of chronologically ordered points
 - $p_1 \rightarrow p_2 \rightarrow \dots \rightarrow p_n$, where $p = (x, y, t)$



Introduction

- Trajectory outliers
 - A trajectory or a segment of trajectory
 - Condition 1: significantly different from other items in terms of some similarity metric
 - Condition 2: do not conform to an expected pattern

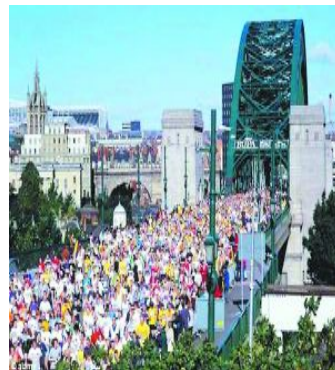


Introduction

- Marathon:
 - Which trajectories are outliers ?
 - Condition 1:



- Condition 2:



Introduction

- Applications
 - Transportation supervision
 - social network
 - Credit card fraud
 -

Part 2 Related Algorithms

Related Algorithm

- Distance-based
 - *Distance-based outliers: algorithms and applications* (Edwin M. Knorr et al. VLDB 2000)
- LOF
 - *LOF: Identifying Density-Based Local Outliers* (Markus M. Breunig et al. SIGMOD 2000)
- LOCI
 - *LOCI: Fast Outlier Detection Using the Local Correlation Integral* (Spiros Papadimitriou et al. ICDE 2003)

Related Algorithm

- TRAOD
 - *Trajectory Outlier Detection: A Partition-and-Detect Framework*
(Jae-Gil Lee et al. ICDE 2008)
- ROAM
 - *ROAM: Rule- and Motif-Based Anomaly Detection in Massive Moving Object Data Sets* (Xiaolei Li et al. SIAM 2007)
- TOD-SS
 - *Trajectory Outlier Detection Algorithm Based on Structural Features*
(Guan YUAN et al. Journal of Computational Information Systems 2011)

Distance-based

- Basic idea:
 - Summarize trajectory global information to a vector(point)
 - Distance function: weighted sum
 - DB(p,D)-outlier detection
- Summarization

$$P = \begin{bmatrix} P_{\text{start}} \\ P_{\text{end}} \\ P_{\text{heading}} \\ P_{\text{velocity}} \end{bmatrix}$$

where:

$$P_{\text{start}} = (x_{\text{start}}, y_{\text{start}}),$$

$$P_{\text{end}} = (x_{\text{end}}, y_{\text{end}}),$$

$$P_{\text{heading}} = (\text{avg}_{\text{heading}}, \text{max}_{\text{heading}}, \text{min}_{\text{heading}}),$$

$$P_{\text{velocity}} = (\text{avg}_{\text{velocity}}, \text{max}_{\text{velocity}}, \text{min}_{\text{velocity}}).$$

Distance-based

- Distance function:

- Weighted sum

$$D(P_1, P_2) =$$

$$\begin{bmatrix} D_{\text{start}}(P_1, P_2) \\ D_{\text{end}}(P_1, P_2) \\ D_{\text{heading}}(P_1, P_2) \\ D_{\text{velocity}}(P_1, P_2) \end{bmatrix} * [w_{\text{start}} \quad w_{\text{end}} \quad w_{\text{heading}} \quad w_{\text{velocity}}]$$

where $w_{\text{start}}, w_{\text{end}}, w_{\text{heading}}, w_{\text{start}}, w_{\text{velocity}}$ are user defined and application-dependent

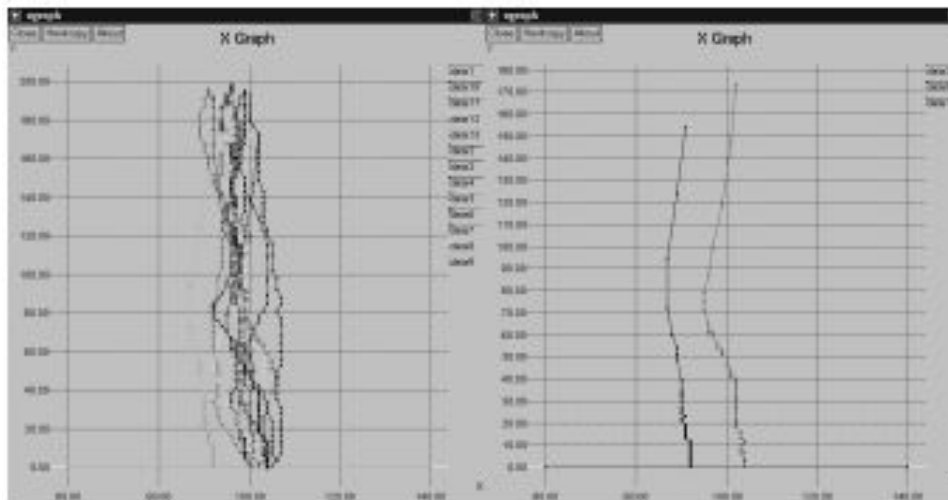
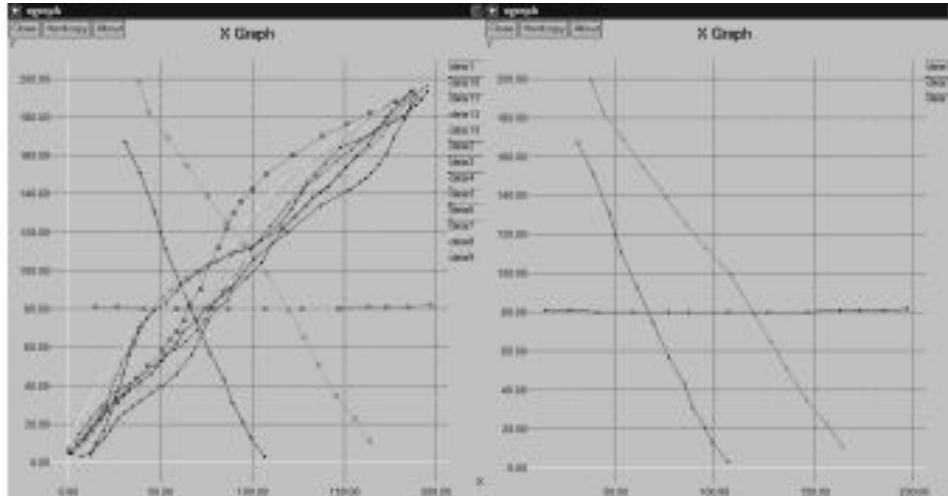
$D(P_1, P_2)$: Euclidean distance

- DB(p,D)-outlier

- An object O in a dataset T is a DB(p, D)-outlier if at least fraction p of the objects in T lies greater than distance D from O.

- P, D are user defined

Distance-based



Distance-based

- Drawbacks
 - Trajectory whole(global) information
 - Can't detect outlying portions
 - Doesn't consider different density distributions

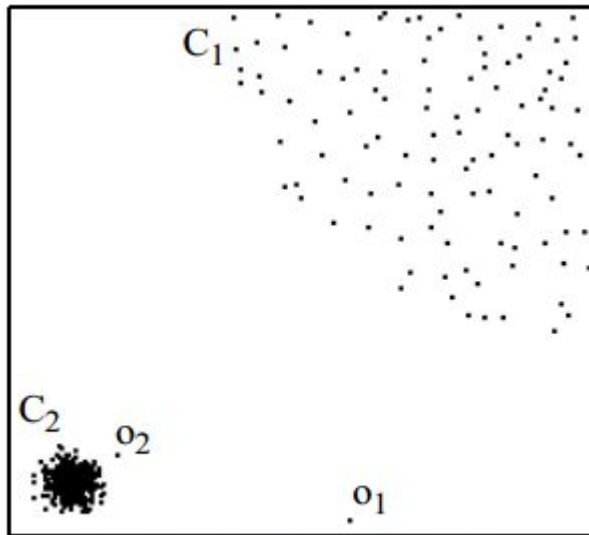


Figure 1: 2-*d* dataset DS1

LOF

- Basic idea:
 - Regard trajectory as point
 - local outlier factor (LOF) : degree of being an outlier
 - Depends on the local density of its neighborhood
- Preliminary
 - k-distance(p):
 - for at least k objects $o' \in D \setminus \{p\}$ it holds that $d(p, o') \leq d(p, o)$
 - for at most k-1 objects $o' \in D \setminus \{p\}$ it holds that $d(p, o') < d(p, o)$
 - k-distance neighborhood of an object p

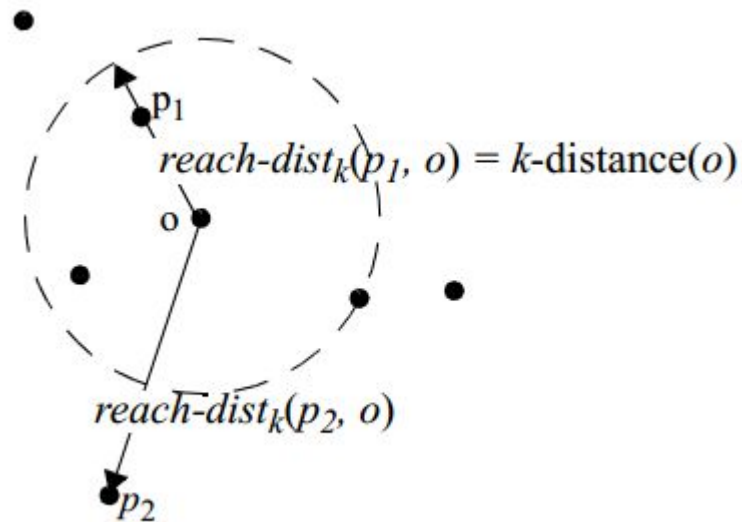
$$N_{k\text{-distance}(p)}(p) = \{ q \in D \setminus \{p\} \mid d(p, q) \leq k\text{-distance}(p) \}$$

LOF

- Preliminary

- Reachability distance of an object p w.r.t. object o

$$\text{reach-dist}_k(p, o) = \max \{ k\text{-distance}(o), d(p, o) \}$$



LOF

- local reachability density of an object p

$$lrd_{MinPts}(p) = 1 / \left(\frac{\sum_{o \in N_{MinPts}(p)} reach-dist_{MinPts}(p, o)}{|N_{MinPts}(p)|} \right)$$

- MinPts : a minimum number of objects, user defined
- (local) outlier factor of an object p

$$LOF_{MinPts}(p) = \frac{\sum_{o \in N_{MinPts}(p)} \frac{lrd_{MinPts}(o)}{lrd_{MinPts}(p)}}{|N_{MinPts}(p)|}$$

- For objects inside a cluster: LOF value is approximately 1
- For other objects: give tight lower and upper bounds on the LOF value

LOF

- Discussion
 - Assign to each object a degree of being an outlier
 - Density based
 - MinPts is non-trivial

LOCI

- Basic Idea
 - Density based
 - Multi-granularity deviation factor(MDEF)
 - Outlier detect:
 - MDEF value deviates significantly (more than three standard deviations) from the local averages

LOCI

- Preliminary

- **N(pi, r):** the set of r-neighbors of pi

$$N(p_i, r) \equiv \{p \in P \mid d(p, p_i) \leq r\}$$

- **n(pi, r):** the number of r-neighbors of pi

$$n(p_i, r) \equiv |N(p_i, r)|$$

- **n̂(pi, r, α):** average of n(p, αr) over the set of r neighbors of pi

$$\hat{n}(p_i, r, \alpha) \equiv \frac{\sum_{p \in N(p_i, r)} n(p, \alpha r)}{n(p_i, r)}$$

- **σn̂(pi, r, α):** standard deviation of n(p, αr) over the set of r-neighbors

$$\sigma_n^{\hat{}}(p_i, r, \alpha) \equiv \sqrt{\frac{\sum_{p \in N(p_i, r)} (n(p, \alpha r) - \hat{n}(p_i, r, \alpha))^2}{n(p_i, r)}}$$

LOCI

- MDEF (p_i, r, α): Multi-granularity deviation factor for point p_i at radius (or scale) r

$$MDEF(p_i, r, \alpha) = \frac{\hat{n}(p_i, r, \alpha) - n(p_i, \alpha r)}{\hat{n}(p_i, r, \alpha)} = 1 - \frac{n(p_i, \alpha r)}{\hat{n}(p_i, r, \alpha)}$$

- Normalized deviation

$$\sigma_{MDEF}(p_i, r, \alpha) = \frac{\sigma_{\hat{n}}(p_i, r, \alpha)}{\hat{n}(p_i, r, \alpha)}$$

LOCI

- LOCI outlier detection

$$MDEF(p_i, r, \alpha) > k_\sigma \sigma_{MDEF}(p_i, r, \alpha), \text{ where } k_\sigma = 3$$

```
// Pre-processing
Foreach  $p_i \in \mathbb{P}$ :
  Perform a range-search
  for  $N_i = \{p \in \mathbb{P} \mid d(p_i, p) \leq r_{max}\}$ 
  From  $N_i$ , construct a sorted list  $D_i$ 
  of the critical and  $\alpha$ -critical distances of  $p_i$ 
// Post-processing
Foreach  $p_i \in \mathbb{P}$ :
  For each radii  $r \in D_i$  (ascending):
    Update  $n(p_i, \alpha r)$  and  $\hat{n}(p_i, r, \alpha)$ 
    From  $n$  and  $\hat{n}$ , compute
       $MDEF(p_i, r, \alpha)$  and  $\sigma_{MDEF}(p_i, r, \alpha)$ 
    If  $MDEF(p_i, r, \alpha) > 3\sigma_{MDEF}(p_i, r, \alpha)$ ,
      flag  $p_i$ 
```

TRAOD

- Basic idea
 - Partition-and-Detect Framework
 - hybrid of the distance-based and density-based approaches

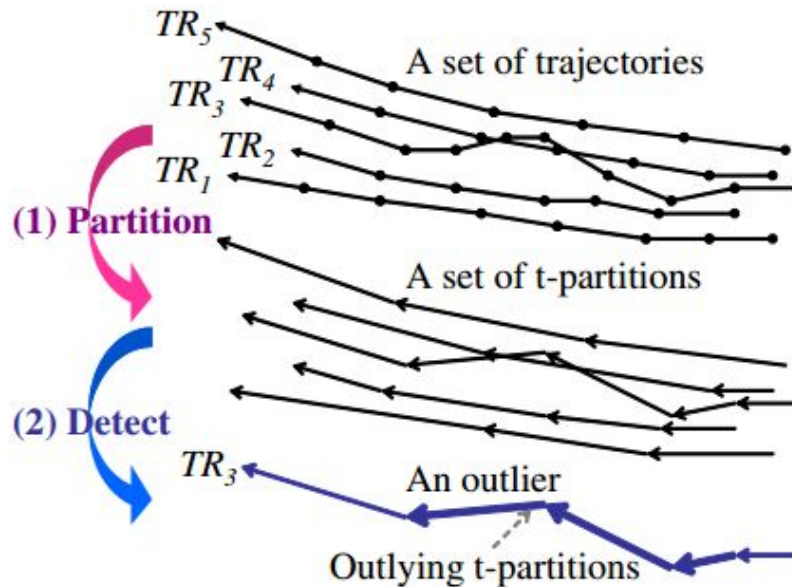


Fig. 2. An example of trajectory outlier detection in the partition-and-detect framework.

TRAOD

- Pre-definition
 - Trajectory
 - $TR_i = P_1 P_2 P_3 \cdots P_j \cdots P_{len_i}$, where P_j ($1 \leq j \leq len_i$) is a d-dimensional point
 - sub-trajectory
 - $P_{c_1} P_{c_2} \cdots P_{c_k}$ ($1 \leq c_1 < c_2 < \cdots < c_k \leq len_i$)
 - trajectory partition
 - line segment $p_i p_j$ ($i < j$)
 - t-partition
 - outlying trajectory partitions
 - t-partition does not have “enough” similar neighbors (i.e., close trajectories)

TRAOD

- Partition
 - MDL Principle
 - Distance between T-Partitions(distance function)
 - the perpendicular distance (d_{\perp})
 - the parallel distance (d_{\parallel})
 - the angle distance (d_{θ})

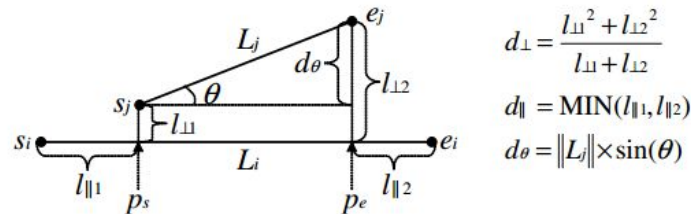


Fig. 4. Three components of the distance function for line segments.

$$\text{dist}(L_i, L_j) = w_{\perp} \cdot d_{\perp} + w_{\parallel} \cdot d_{\parallel} + (w_{\theta} \cdot d_{\theta}) \cdot L_i \cdot L_j$$

TRAOD

- Trajectory outlier detection

- Close

A trajectory TR_i is *close* to a t-partition $L_j \in P(TR_j)$ ($TR_i \neq TR_j$) if

$$\sum_{L_i \in CP(TR_i, L_j, D)} \text{len}(L_i) \geq \text{len}(L_j)$$

D is a parameter given by a user.

$P(TR_i)$: the set of all t -partitions of TR_i

$CP(TR_i, L_j, D)$: the set of TR_i 's t -partitions within the distance D from $L_j \in P(TR_j)$ ($TR_i \neq TR_j$)

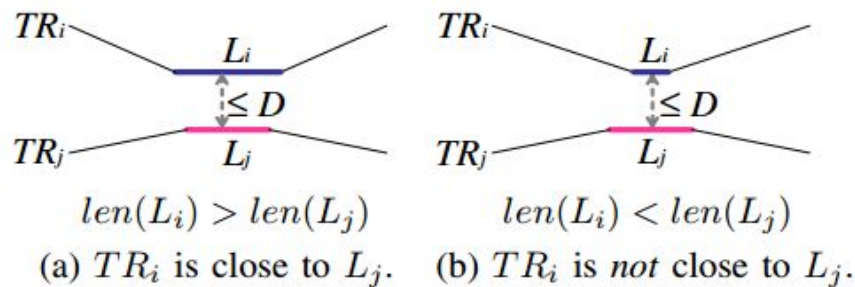


Fig. 3. The concept of the close trajectory.

TRAOD

- Trajectory outlier detection

- Outlying

A t-partition $L_i \in P(TR_i)$ is *outlying* if

$$|CTR(L_i, D)| \leq \lceil (1-p)|I| \rceil$$

$|I|$ indicates the total number of trajectories.

p is a parameter given by a user.

$CTR(L_i, D)$: the set of trajectories close to L_i

- Outlier

$$Ofrac(TR_i) = \frac{\sum_{L_i \in OP(TR_i, D, p)} len(L_i)}{\sum_{M_i \in P(TR_i)} len(M_i)} \geq F$$

- Where F is a parameter given by a user

TRAOD

- Incorporation of Density: adjusting coefficient

- Density

$$density(L_i) = | \bigcup_{TR_j \in I} CP(TR_j, L_i, \sigma) |$$

σ is the standard deviation of pairwise distances between t – partitions

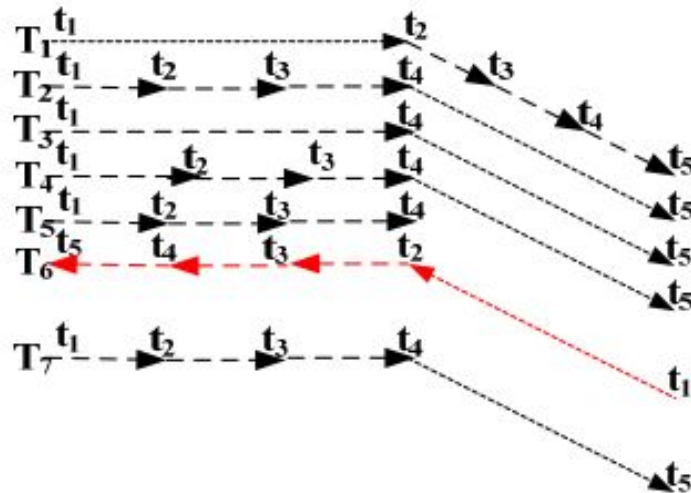
- adjusting coefficient

$$adj(L_i) = \frac{\sum_{L_j \in L} density(L_j) / |L|}{density(L_j)}, \text{ where } L = \bigcup_{TR_k \in I} P(TR_k)$$

$$adj(L_i) * |CTR(L_i, D)|$$

TRAOD

- Discussion
 - Detect outlying sub-trajectories
 - Without the temporal information



ROAM

- ROAM (Rule- and Motif-based Anomaly Detection in Moving Objects)
- Basic idea

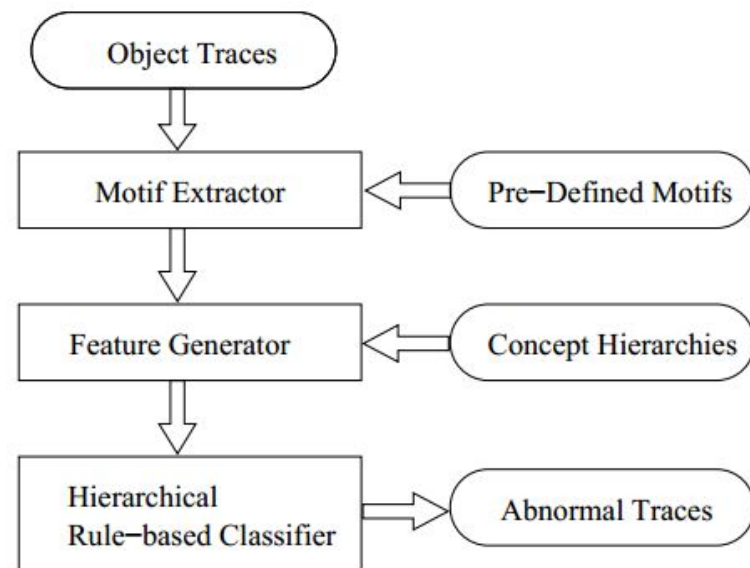
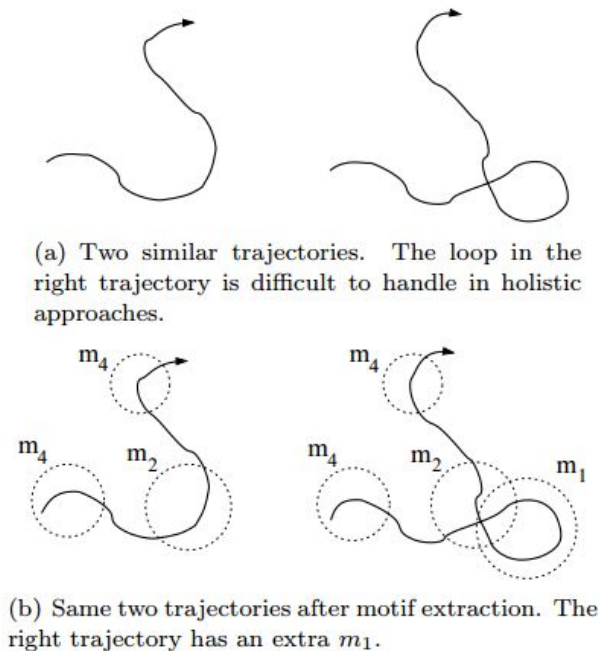


Figure 2: ROAM Framework

ROAM

- Motif Extractor

- Sliding window

- Window length w : user defined
- All windows are overlaid on top of each other

- Clustering

- group trajectory segments into representative sets

$$(m_i, t_{\text{start}}, t_{\text{end}}, l_{\text{start}}, l_{\text{end}})$$

- Motif expressions

- duration, top speed, avg speed, radius, and general location

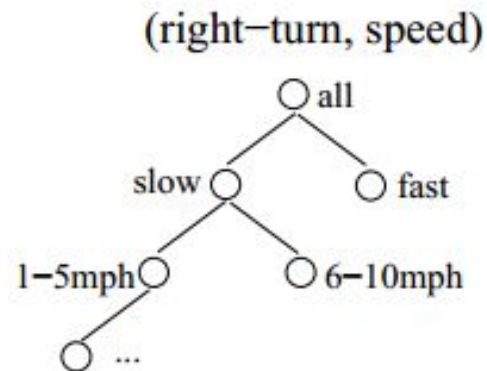
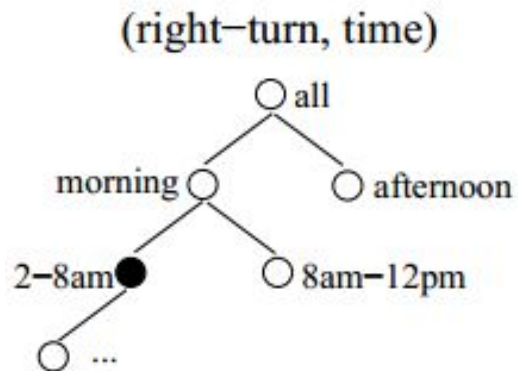
$$(m_i, V_1, V_2 \dots V_A)$$

ROAM

- Feature Generator

$$(m_i, a_j, v_k) \leftrightarrow f_x$$

- Feature generalization
 - tree-based data structure



ROAM

- Classification

- CHIP (Classification using Hierarchical Prediction Rules)

- Iteratively and greedily searches for the best available rule until all positive examples are covered
 - Tries to use high-level features

- A single rule r :

- $l_1 \wedge l_2 \wedge \dots \wedge l_n \rightarrow C$
 - l_i : a literal (or predicate) of the form (feature = value)
 - C : the class label
 - An example is “covered” by r if all the literals in r are satisfied in the example

ROAM

- Classification

- Foil Gain(l, r)

- A rule is learned one literal at a time. Literals are selected according to a weighted version of Foil Gain
 - Based on the positive and negative coverage of the rule before and after adding the literal

$$p_1 \left(\log_2 \frac{p_1}{p_1 + n_1} - \log_2 \frac{p_0}{p_0 + n_0} \right)$$

- p_0, n_0 : the number of positive and negative examples covered by rule r without literal l
 - p_1, n_1 : the number of positive and negative examples covered by rule $r \wedge l$

ROAM

- Classification

- $Exp_Gain(f, r)$

- The maximum Foil Gain achieved by any literal in any of its child features

$$Exp_Gain(f, r) = \max_{(l, f_i) \forall l, f_i \in Exp(f)} Foil_Gain(l, r)$$

- $Exp(f)$: the set of f 's children in F 's hierarchy
 - f is a feature and r is a rule

ROAM

Input: (1) Training set $\mathcal{D} = P \cup N$, where P and N are the positive and negative examples. (2) Initial feature set $\mathcal{F}_C \in \mathcal{F}'$.

Output: Set of classification rules R .

Method:

1. **while** not all of P is covered
2. initialize new rule r
3. **while true**
4. find literal l with highest $Foil_Gain(l, r)$
5. find feature f with highest $Exp_Gain(f, r)$
6. **if** both gains $<$ min_gain **then break**
7. **if** $Foil_Gain(l, r) > \beta \cdot Exp_Gain(f, r)$ **then**
8. add l to r
9. **else**
10. add feature f to \mathcal{F}_C
11. add r to R
12. **return** R

ROAM

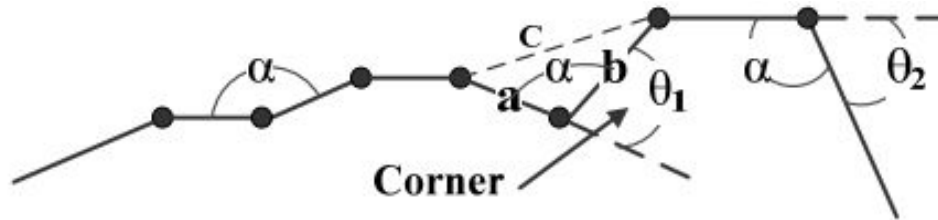
- Discussion
 - Detect outliers with multiple granularity of spatiotemporal features
 - Training set

TOD-SS

- Trajectory Outlier Detection algorithm based on Structural Similarity (TOD-SS)
- Basic idea
 - Partition trajectory : corner threshold
 - Abstract trajectory structure features
 - Outlier detection:
 - Segment outlier : SSIM(Structural Similarity) matrix
 - Trajectory outliers: the proportion of its segment outliers

TOD-SS

- The Partition of Trajectory



$$\alpha = \arccos\left(\frac{a^2 + b^2 - c^2}{2ab}\right) \quad \theta = \begin{cases} \pi - \alpha, & \text{if } (\vec{a} \times \vec{b} \geq 0) \\ \alpha - \pi, & \text{if } (\vec{a} \times \vec{b} < 0) \end{cases}$$

- Partition point: $|\theta| > \omega$, ω is user defined

TOD-SS

- Trajectory Structure

Trajectory Structure <Direction, Speed, Angle, Location>

- Direction

$$DirDist(L_i, L_j) = \begin{cases} \min(\|L_i\|, \|L_j\|) \times \sin(\varphi), & \text{if } (0 \leq \varphi \leq 90) \\ \min(\|L_i\|, \|L_j\|), & \text{if } (90 \leq \varphi \leq 180) \end{cases}$$

$\|L_j\|$ The length of L_j , and $j \in N \wedge 1 \leq j \leq |TS|$

- Speed:

$$SpeedDist(L_i, L_j) = \frac{1}{3} (S_{\max}(L_i, L_j) + S_{\text{avg}}(L_i, L_j) + S_{\min}(L_i, L_j))$$

- $S_{\max}(L_i, L_j) = |V_{\max}(L_i) - V_{\max}(L_j)|$
- S_{avg} and S_{\min} : the absolute divergence of average and minimum speed.

TOD-SS

- Trajectory Structure

- Angle:

$$AngleDist(L_i, L_j) = \frac{\sum_{1,1}^{P(L_i), P(L_j)} (|\theta_i - \theta_j|) / (|\theta_i| + |\theta_j|)}{P(L_i) + P(L_j)}$$

$P(x)$ The number points of a trajectory or a segment, and $x \in \{TD\} | x \in \{TS\}$

- Location

$$LocDist(L_i, L_j) = \max(h(L_i, L_j), h(L_j, L_i))$$

- $h(L_i, L_j) = \max_{a \in L_i} (\min_{b \in L_j} (dist(a, b)))$: Hausdorff distance

TOD-SS

- Structural Similarity

- Feature weights : user defined

- $W = \{W_D, W_S, W_A, W_L\}$

- $W_D \geq 0, W_S \geq 0, W_A \geq 0, W_L \geq 0; \quad W_D + W_S + W_A + W_L = 1.$

- Structural distance (SDIST)

$$SDIST(L_i, L_j) = (DirDist \times W_D + SpeedDist \times W_S + AngleDist \times W_A + LocDist \times W_L)$$

$$SSIM(L_i, L_j) = 1 - N(SDIST(L_i, L_j))$$

- Where $N(\dots)$ is a function of distance normalization

TOD-SS

- Structural Similarity
 - SSIM of a whole trajectory

$$MSSIM(TR_i, TR_j) = \frac{1}{\min(|L_i|, |L_j|)} \sum_{i=0, j=0}^{\min(|L_i|, |L_j|)} SSIM(L_i, L_j)$$

- Trajectory outlier detect
 - **ϵ -Neighbor:**
 - $SSIM(L_i, L_j) \geq \epsilon$, $L_j \in N_\epsilon(L_i)$, ϵ is user defined
 - **Segment Outlier:**
 - ϵ -neighbor count $|N_\epsilon(L_i)| < \sigma$,
 - σ is the neighbor threshold specified by users.

TOD-SS

- Trajectory outlier detect
 - **Trajectory Outlier:**
 - The proportion of outlier segments in $TR_i \geq$ given proportion threshold ξ
 - The total similarity (MSSIM) \leq threshold F

Table 1 Basic Symbols Notation.

Parameters	Description
TD, TS	Trajectory and Segment set
$\ TR_i\ $	The length of TR_i , and $(i \in N) \wedge (1 \leq i \leq TD)$
$\ L_j\ $	The length of L_j , and $j \in N \wedge 1 \leq j \leq TS $
$P(x)$	<u>The number points of a trajectory or a segment, and $x \in \{TD\} \cup \{TS\}$.</u>
$V_{flag}(L)$	Some kinds speed of L , where <i>flag</i> is <i>s, e, max, min, avg</i> , and represents start, end, maximum, minimum, average speed respectively.
$SO_{\epsilon, \sigma}(L)$	L is the ϵ, σ related segment outlier
$TO_{\zeta, F}(TR)$	TR is the ζ, F related trajectory outlier
$\theta, \omega, \epsilon, \sigma$	Corner, corner threshold, SSIM threshold, and neighbor threshold respectively.
ζ, F	Proportion threshold and structural similarity threshold



TOD-SS

Algorithm: Similarity Computing

Input: Trajectory set $I = \{TR_1, TR_2, \dots, TR_n\}$, ω (corner threshold), W (feature weights)

Output: Similarity matrixes M_L (Segments) and M (Trajectories)

*/*First phase: Trajectory Partition*/*

01: for each $TR_i \in I$ **do**

02: $TS \leftarrow$ Partition TR_i according to ω ; */*Second phase: Computation of Structural Similarity*/*

03: for each $L_i, L_j \in TS \wedge i \neq j \wedge L_i \cdot TR \neq L_j \cdot TR$ **do**

04: Initialize $SSIM(L_i, L_j)$;

05: $SSIM_Computer(L_i, L_j) * \{W\}$; */* multiply by weighs */*
end for

06: Generate Segment Matrix M_L ;

07: for each $TR_i, TR_j \in I \wedge i \neq j$ **do** */* compute similarity matrix*/*

08: Initialize $MSSIM(TR_i, TR_j)$;

09: Calculate $MSSIM(TR_i, TR_j)$;
end for

10: Generate Trajectory Matrix M ;

End.

Algorithm: Outlier Detection

Input: Similarity matrixes M_L (Segments), M (Trajectories),
 ε (neighbor threshold), σ, ζ, F

Output: Segment Outliers, Trajectory Outliers

01: for each $Segment\ Pair(L_i, L_j) \in M_L$ **do**

02: if $SSIM(L_i, L_j) \geq \varepsilon$ **then**

03: $N_\varepsilon(L_i) \leftarrow L_j$; */*set the ε -neighbor of L_i */*

04: for each $L_i \in TS$ **do**

05: if $|N_\varepsilon(L_i)| \leq \sigma$ **then**

06: Set $SO_{\varepsilon, \sigma}(L_i) \leftarrow L_i$; */*mark L_i as an outlier*/*

07: for each $TR_i \in I$ **do**

08: if $Count(L\ is\ an\ outlier \wedge L \in TR_i) / TR_i \cdot SegmentCount > \zeta \wedge$
 $MSSIM(TR_i) < F$ **then**

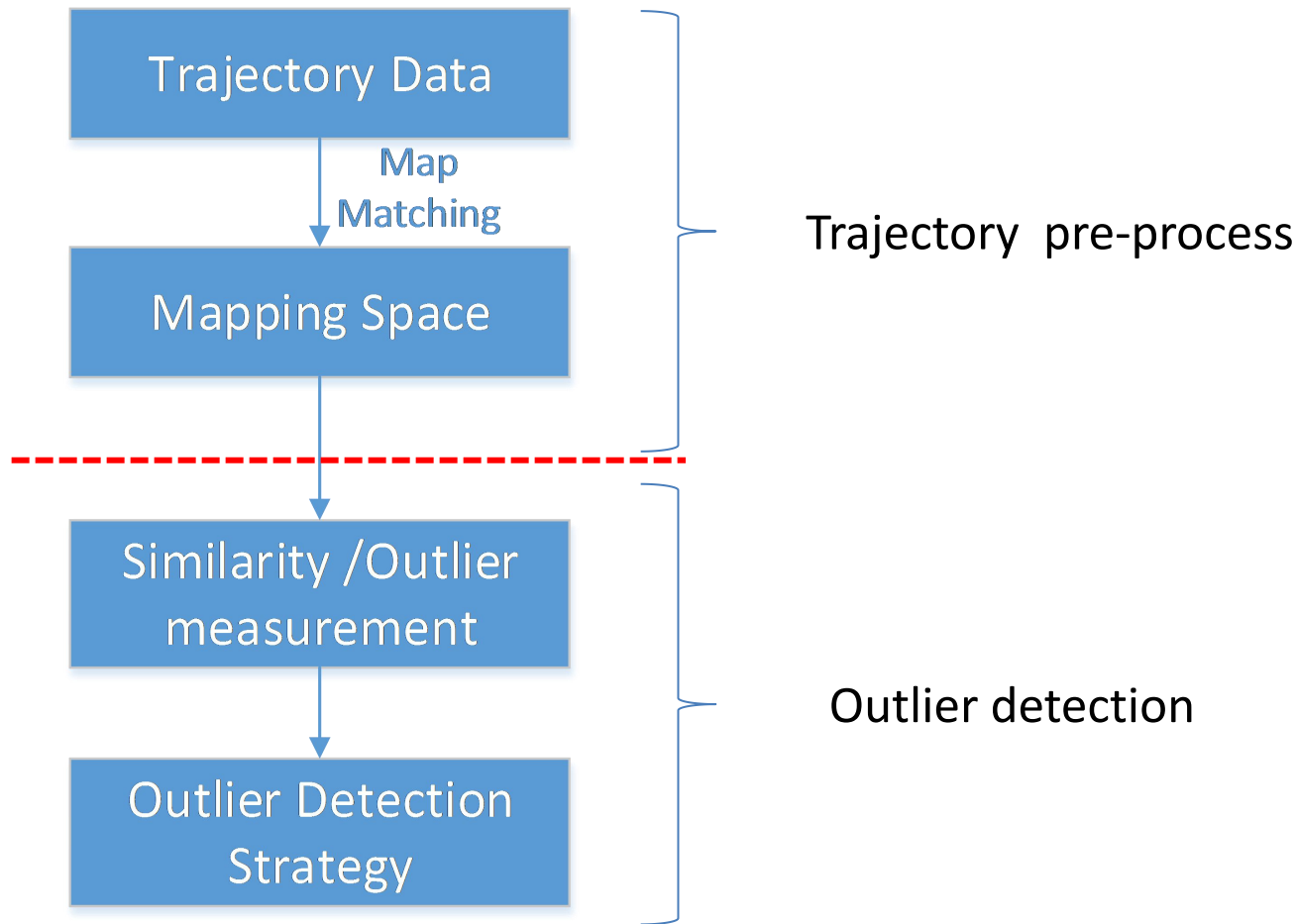
09: Set $TO_{\varepsilon, F}(TR_i) \leftarrow TR_i$; */*mark TR_i as an outlier*/*

End.

Part 3 Conclusion & Discussion

Conclusion & Discussion

- Basic Framework

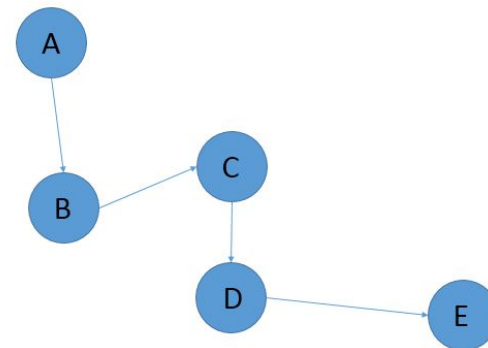


Conclusion & Discussion

- Trajectory pre-process
 - Raw data: $P(L, t)$
 - Space mapping:
 - Trajectory \rightarrow Vector(Point)
 - Merit: can utilize common anomaly detection approaches
 - Defect: may exist information loss
 - Trajectory \rightarrow Line segment
 - Merit: original space; intuitive
 - Defect : partition strategy ; temporal property

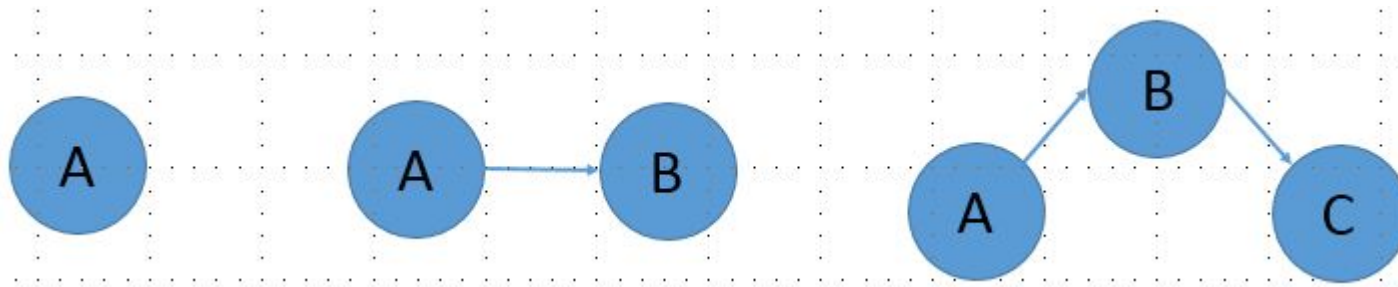
Conclusion & Discussion

- Trajectory outlier detection
 - Doesn't have a uniform definition
 - Personal idea:
 - Location distribution
 - Temporal information
 - Semantic meanings
 - Context information
- Future work
 - Trajectory → Region of interest
 - Pattern mining



Conclusion & Discussion

- Future work
 - Pattern mining



- Global temporal distribution
- Similarity measurement
- Outlier detection