

Trajectory Outlier Detection

Usman Gohar

About me

- B.S in Computer Engineering
- Data Science Intern – Progos Tech (2017)
- M.Sc in Computer Science (University of Minnesota Duluth)
- Currently – Software Developer (Data Science) OATI (Open Access Technology International)
- Towards Data Science Medium Contributor

[Medium.com/@usman.gohar](https://medium.com/@usman.gohar)

linkedin.com/in/usman-gohar

Outline

- Introduction
- DBSCAN
- Methodology
- Results
- Conclusion & Q/A

Goal

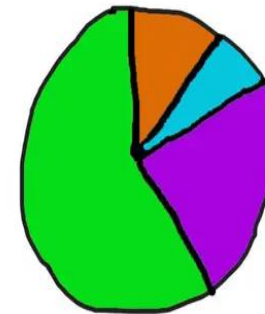
THOUGHTS YOU HAVE ON
THE FIRST DAY OF A NEW JOB:



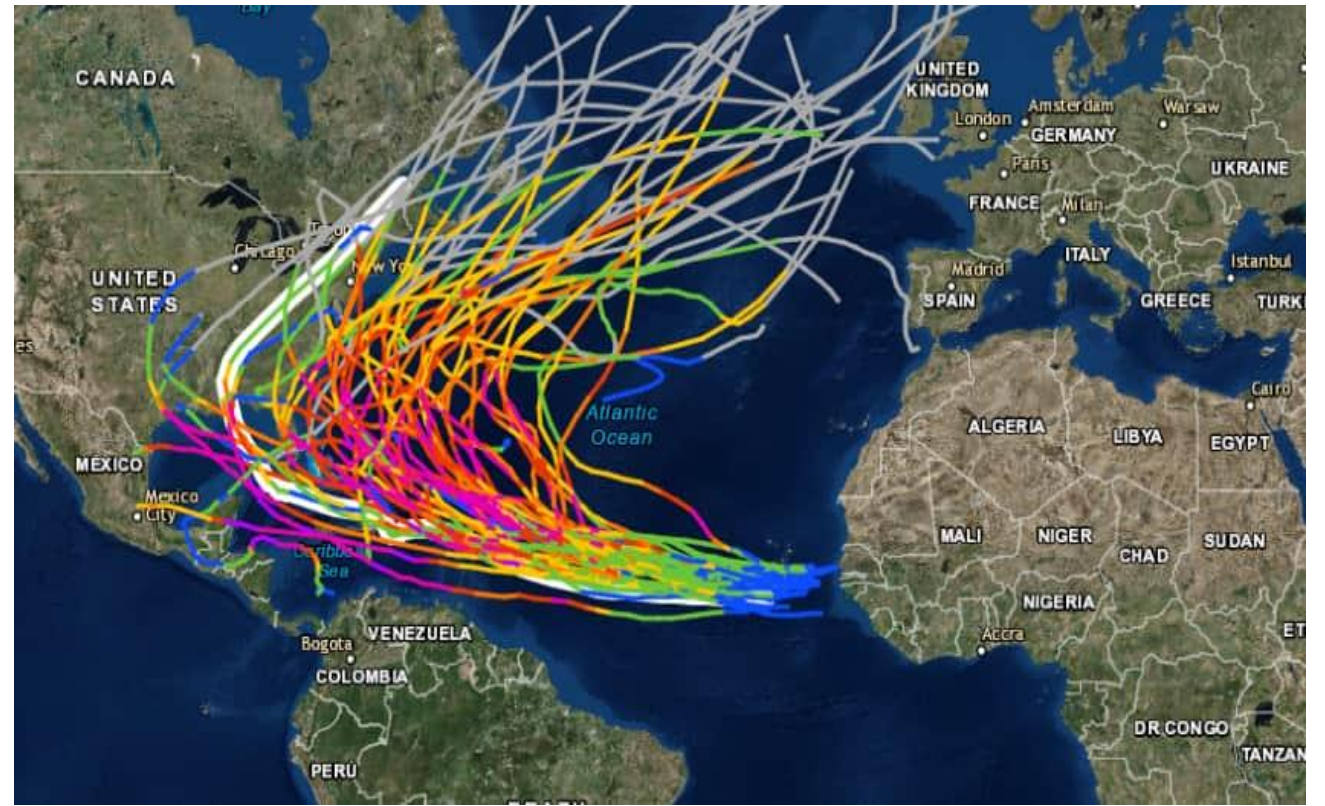
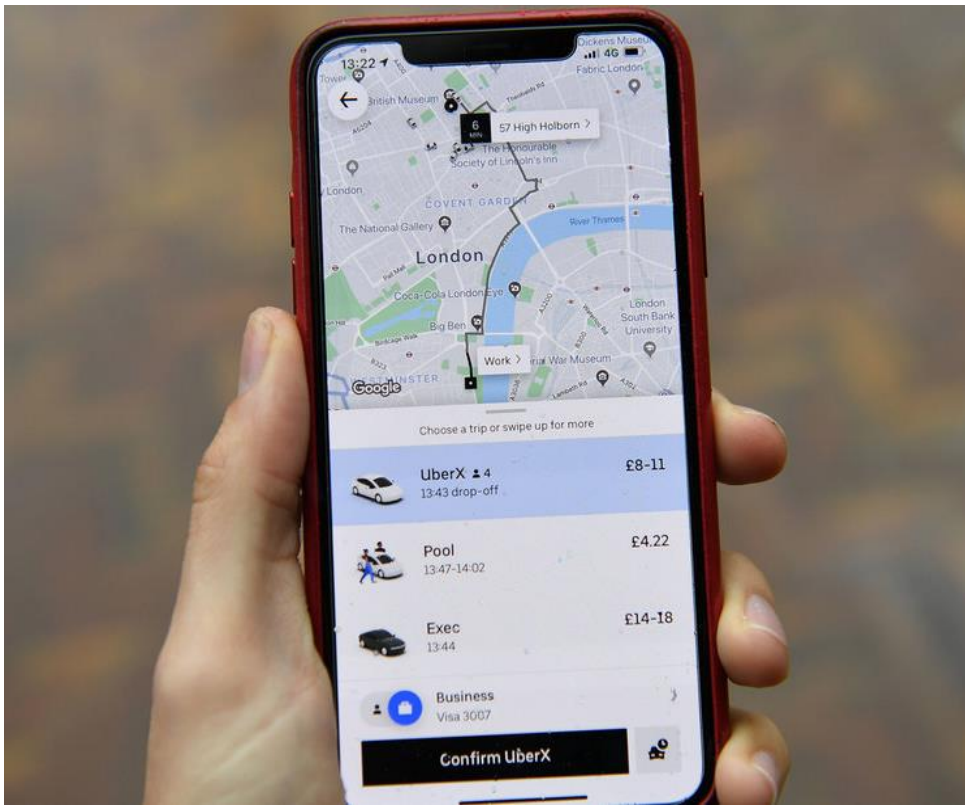
- MY BOSS IS GOING TO REALIZE I WAS A HUGE MISTAKE.
- MY BOSS IS GOING TO REALIZE I WAS A HUGE MISTAKE.



YOUR THOUGHTS WHEN SOMEONE SAYS THAT
YOU WOULD BE GOOD FOR A JOB/ROLE/TEAM:



- WHAT?
- WHY?
- HAVE YOU MET ME?
- MAYBE THEY'RE JUST TRYING TO BE NICE

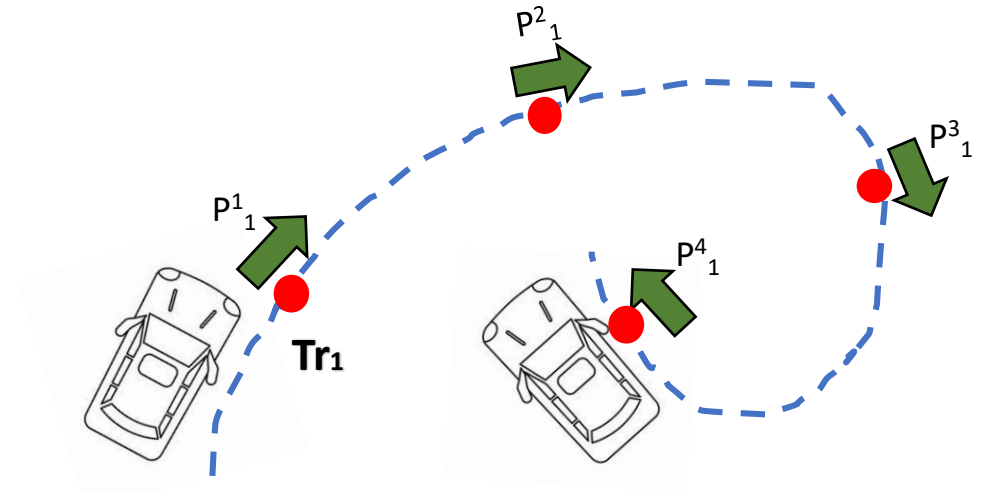


What is a Trajectory?

- Informally, a trajectory is a path an object takes over time
- Or even simply, a record of the movement of an object

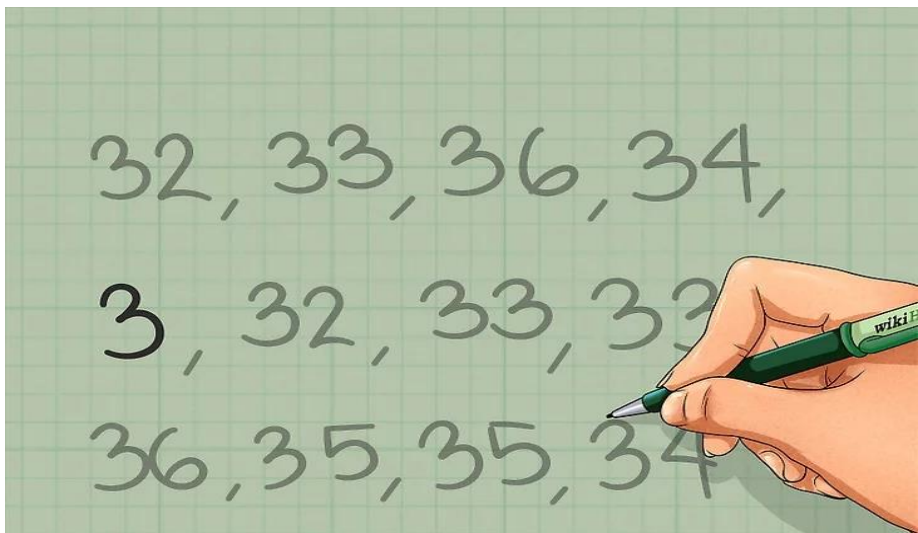
Meet a Trajectory

- A single multidimensional point P_i^j generated from a moving object O , using a tracking device, at time-bin t_j , is called a trajectory point of the trajectory T_{ri}
- The trajectory of a moving object O is then defined as a sequence of such trajectory points produced at time-bins $\{t_1, t_2, \dots, t_j\}$ denoted as $T_{ri} = \{P_i^1, P_i^2, \dots, P_i^j\}$ as shown on the right
- A time-bin is the most smallest unit of time interval for a trajectory which can consist of a single or multiple points that record its movements



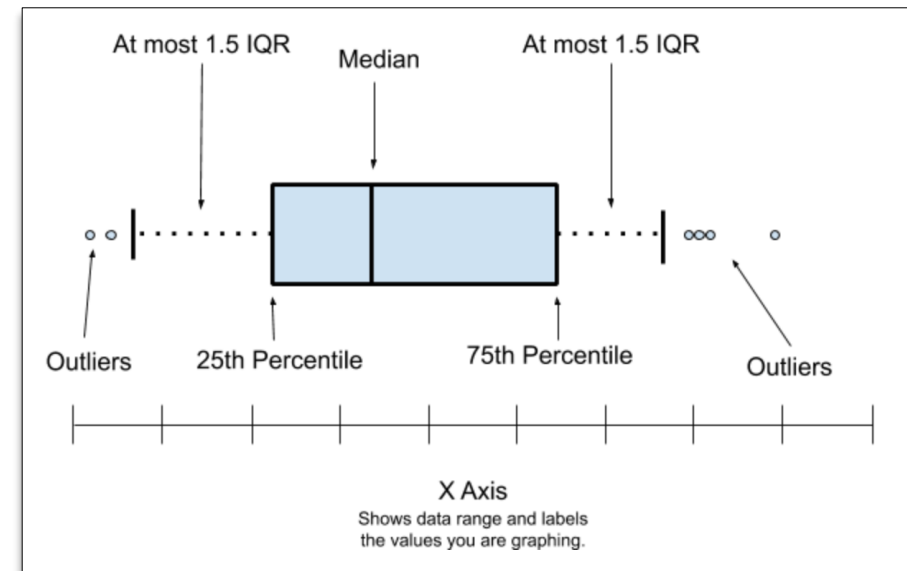
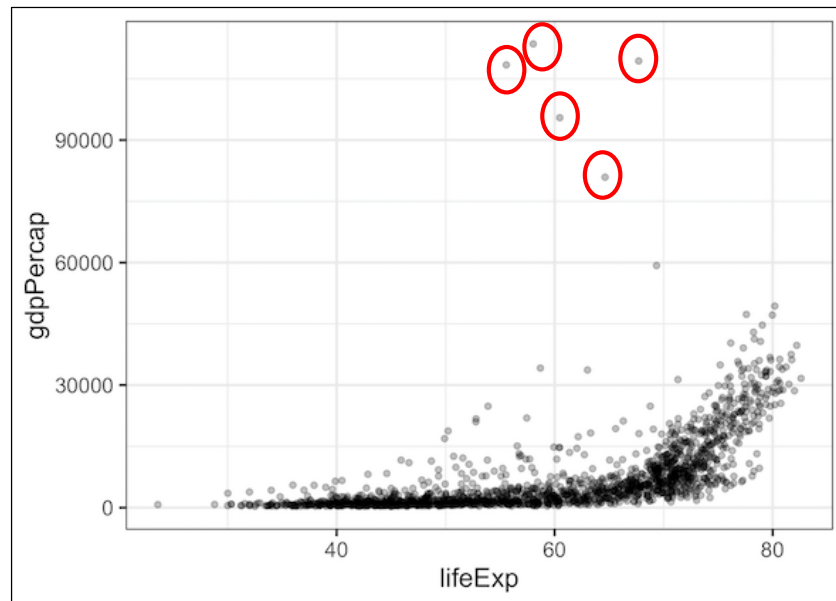
Meet the Outlier

- In statistics, it is simply an observation that is different from other observations (Wiki)
- It could be textual data, image etc.



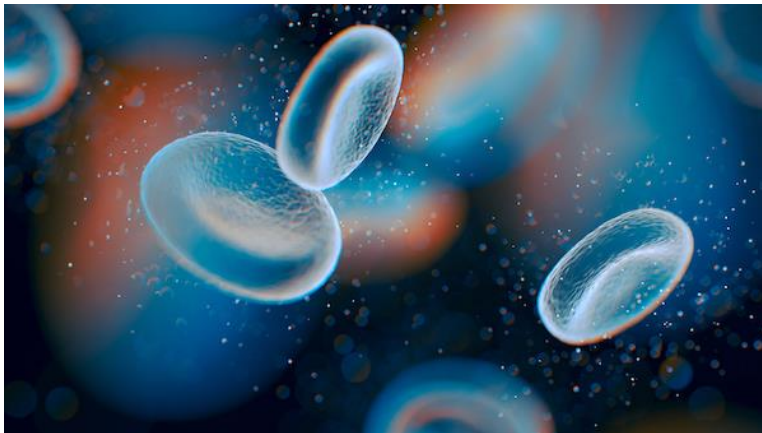
Find the Outlier

Formally, any point in the data that does not conform to the overall underlying pattern in the data



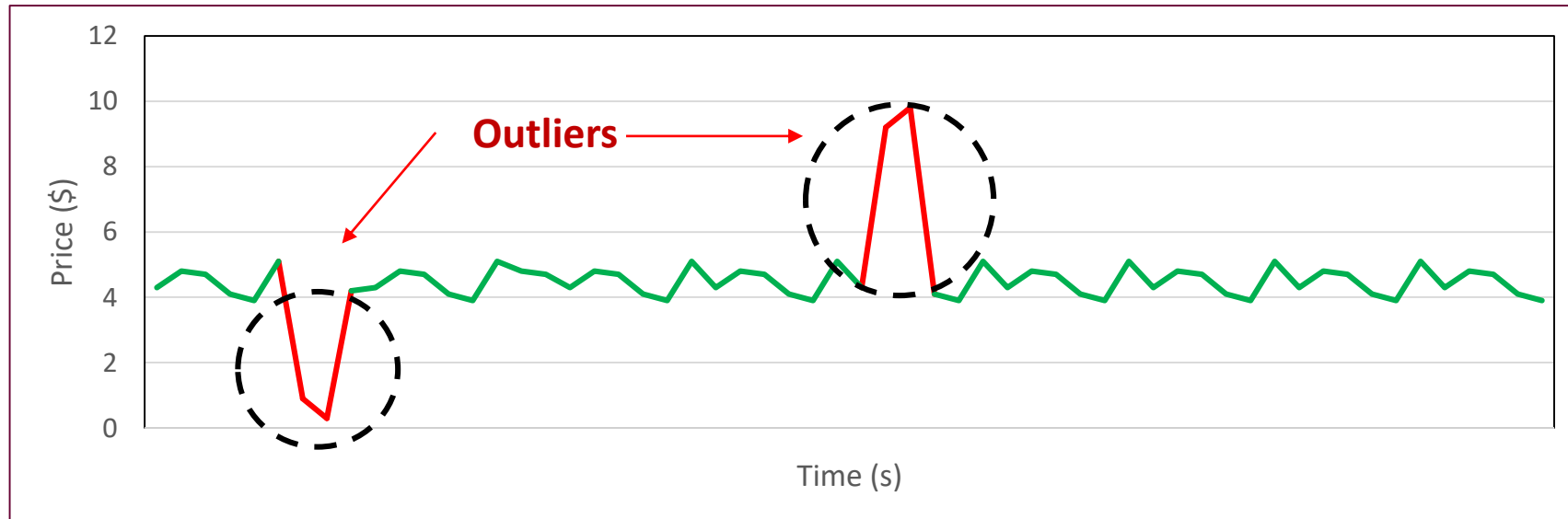
Good or Bad?

- Credit Card Fraud Detection
- Cancerous Cell Detection through imagery
- Interesting information about the data (Machine failing)
- Sway your prediction model e.g. in regression
- Reduce accuracy
- Loss of data



What are Trajectory Outliers?

- Outliers are data points that behave inconsistently compared to the general pattern of data points in the all the trajectories in the window
- Example below shows two outliers in a stock price “trajectory”



Applications

- Real time outlier detection is a time-critical process
- Applications include traffic management, stock-price monitoring



Source: Google Images

**Imperative that
outliers are detected
in a timely fashion!!**



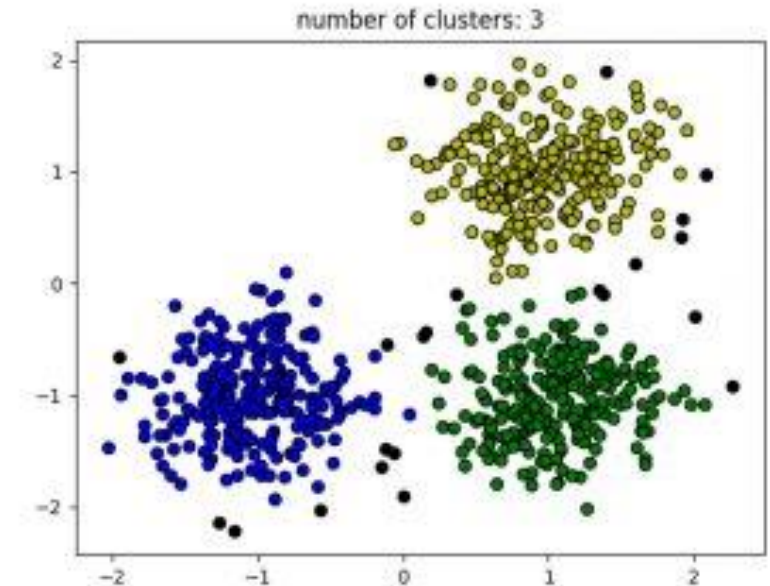
Source: Google Images

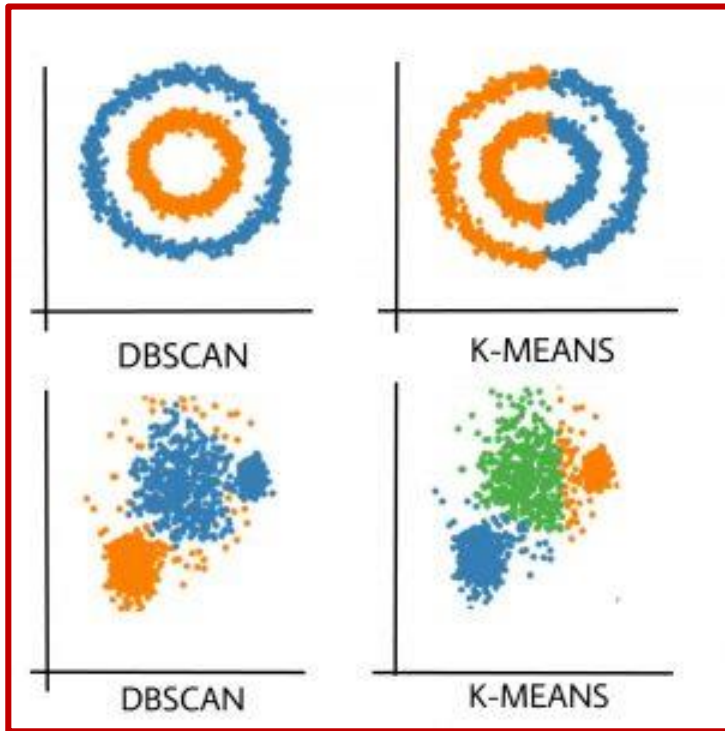
Outline

- Introduction
- **Literature Review**
- Methodology
- Results
- Conclusion & Future Work

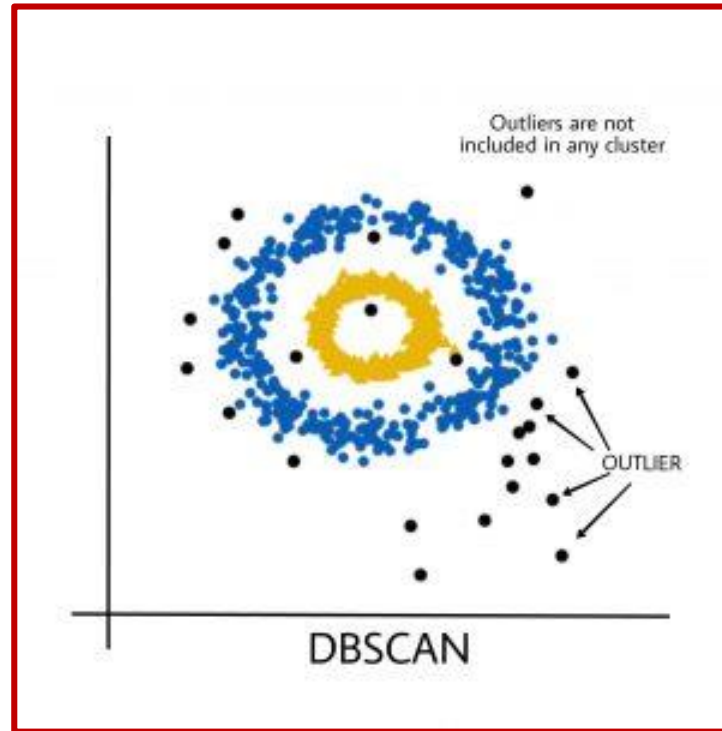
DBSCAN

- Clustering algorithm (similar to K-means)
- Clusters points based on density
- Algorithm as two parameters:
 - 1) Eps (Distance threshold between two points)
 - 2) MinPoints (Equal to number of Dimensions + 1)

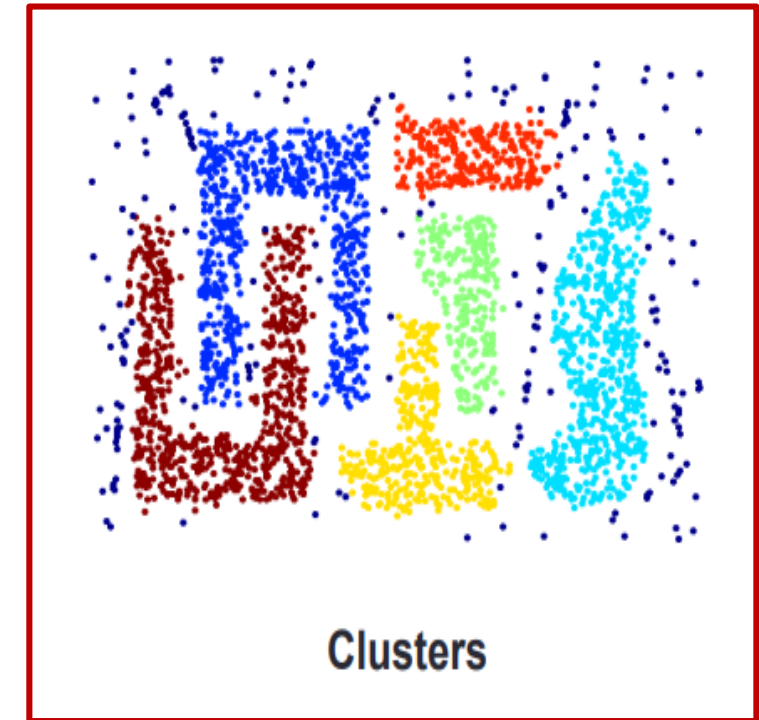




1) K-Means form spherical clusters only



2) K-Means skewed by outliers



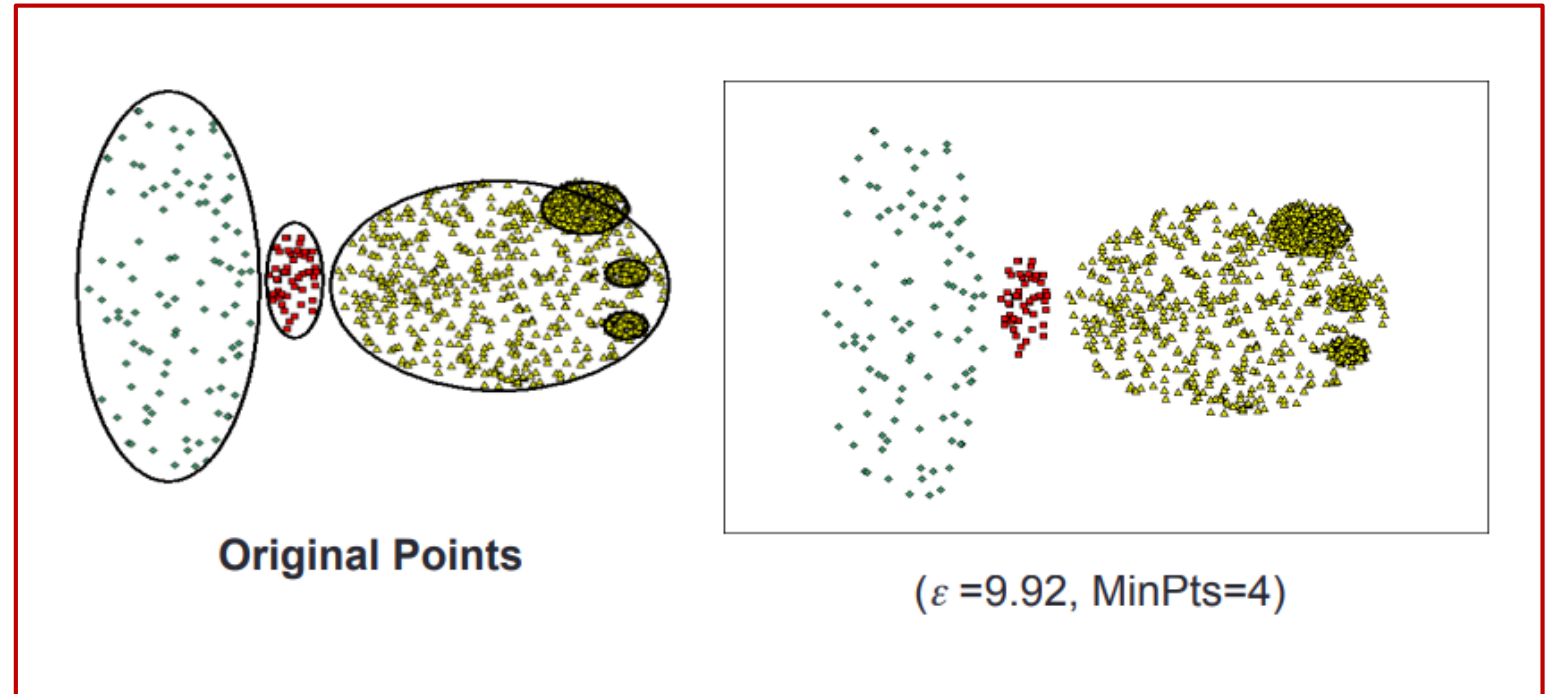
3) Can handle clusters of various shapes

DBSCAN vs K-Means

DBSCAN

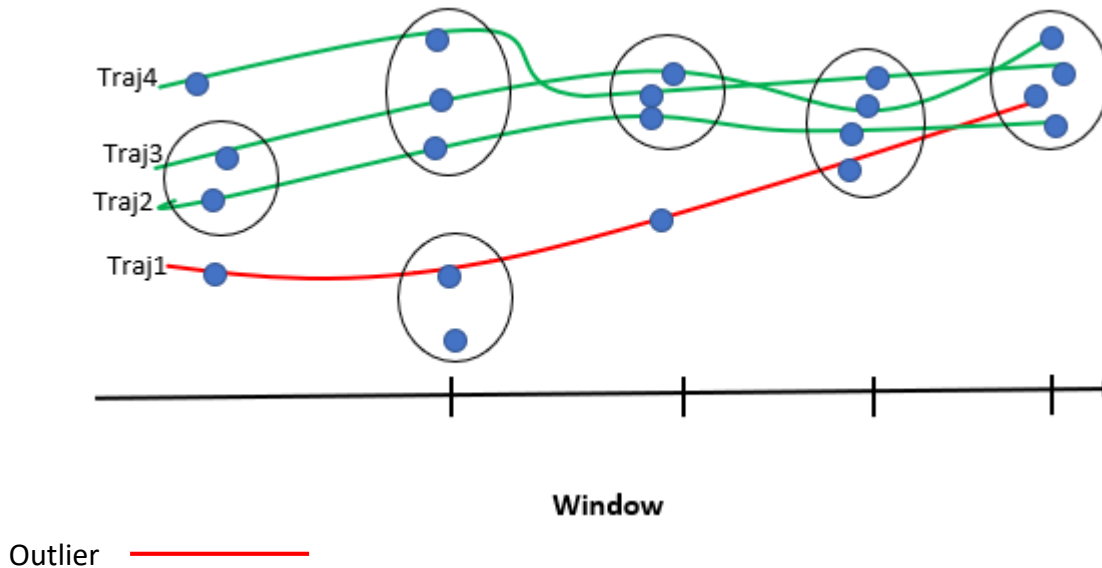
Disadvantages

- Cannot handle varying densities
- Sensitive to parameters
- Can't apply directly to trajectories



Objectives

Outlier Detection over Massive-Scale Trajectory Streams



Parameters:

- Distance Threshold d
- Neighbor Count Threshold k
- Time-Bin Count thr

Definition:

Given a distance threshold d , neighbor count threshold k and time-bin count threshold thr , a trajectory is an outlier in the window if it has fewer than k trajectory neighbors i.e.

$$|TN(Tri, d, thrj)| < k$$

Haversine Approximation

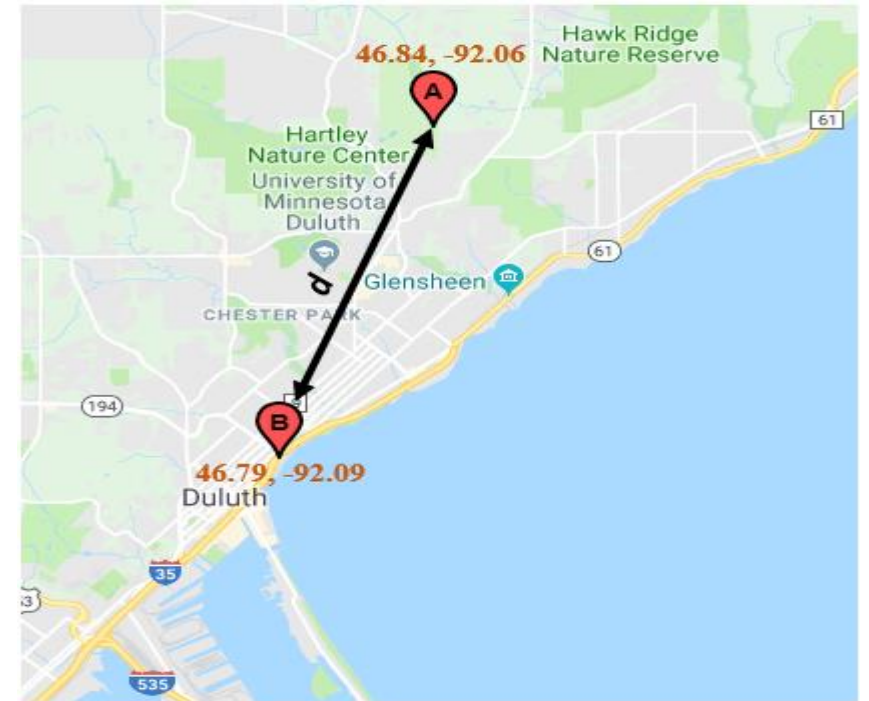
- Calculates the shortest distance between two points on a sphere given the latitude and the longitude

$$a = \sin^2\left(\frac{\Delta\phi}{2}\right) + \sin^2\left(\frac{\Delta\gamma}{2}\right) * \cos(\theta_1) * \cos(\phi_2)$$

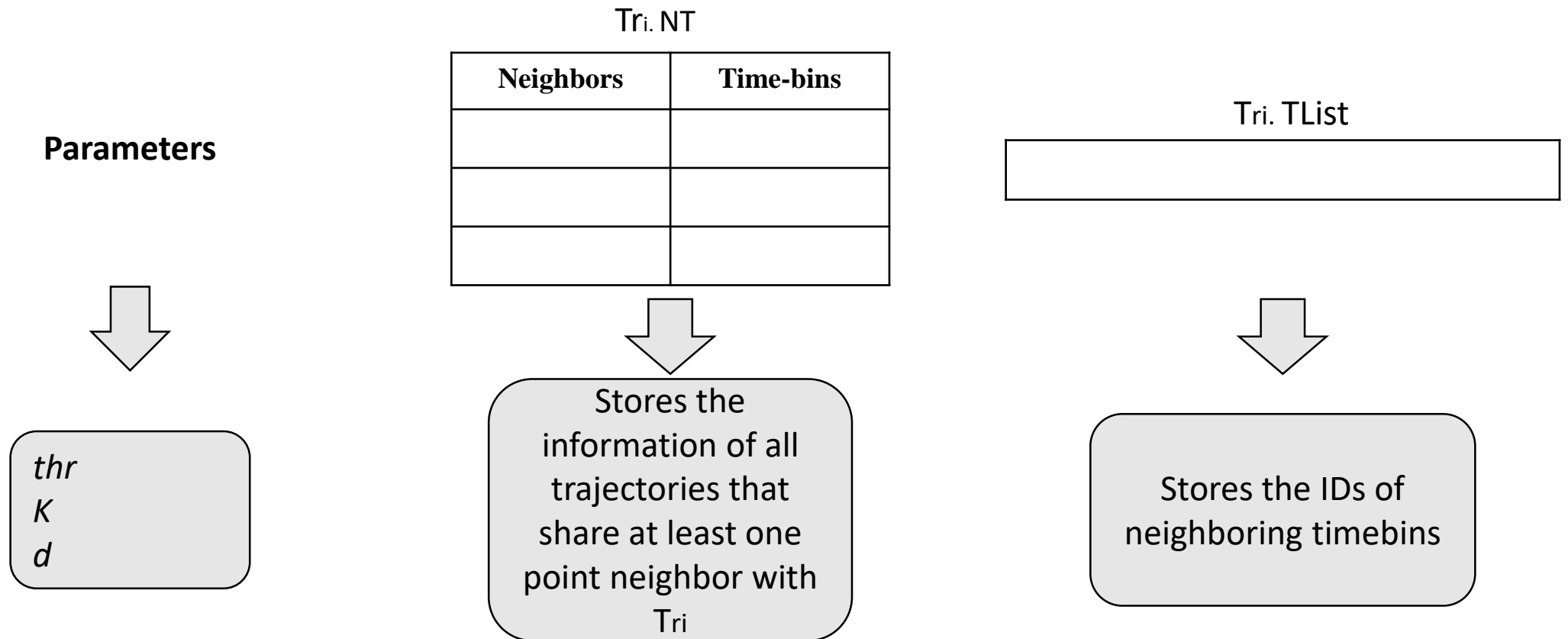
$$c = 2 * \text{atan2}(\sqrt{a}, \sqrt{1-a})$$

$$d = R * c$$

where ϕ , γ , R is the latitude, longitude and the radius of the query search respectively



Data Structures of ODMTS



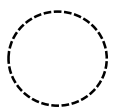
Range Query (Neighbors)

Find all neighbors of each point in each time-bin

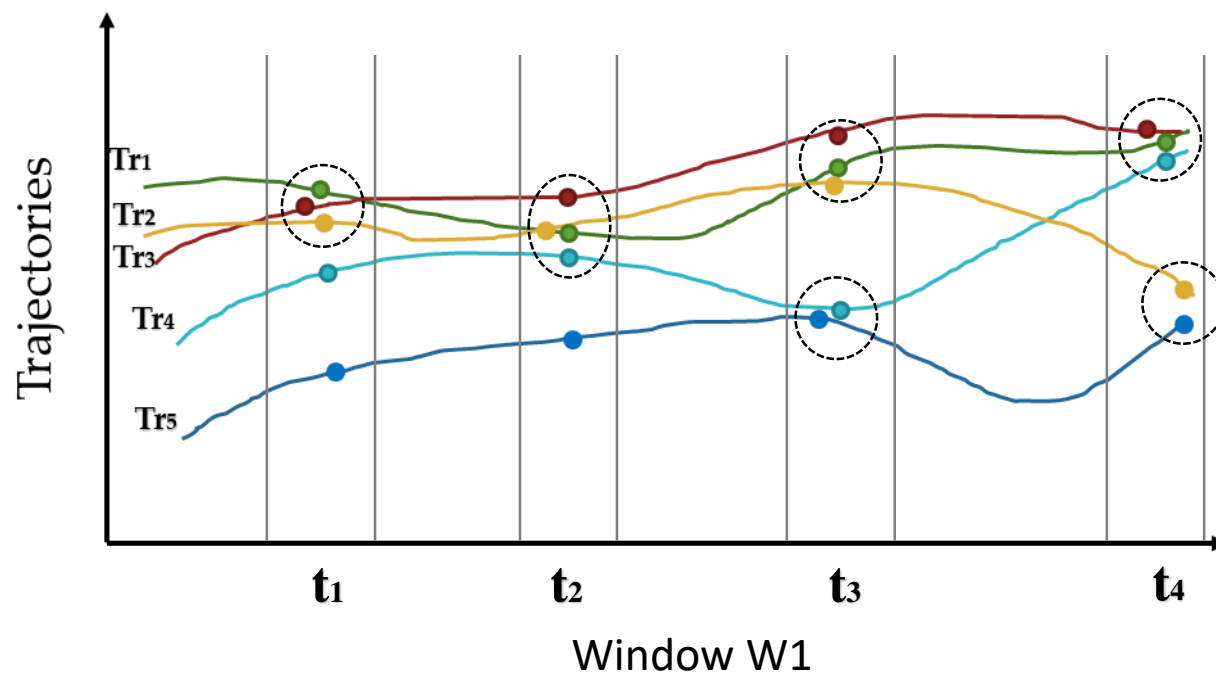
$D = 500\text{m}$

$\text{Thr} = 2$

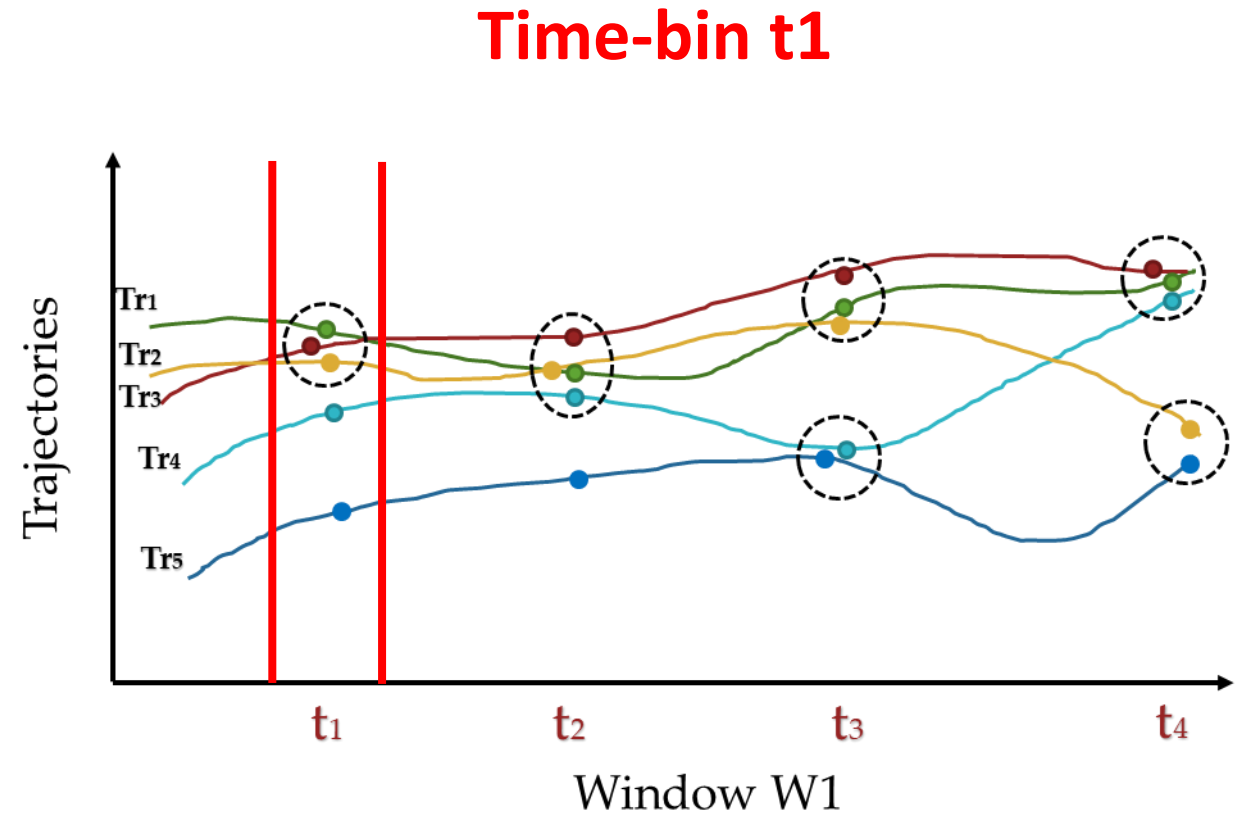
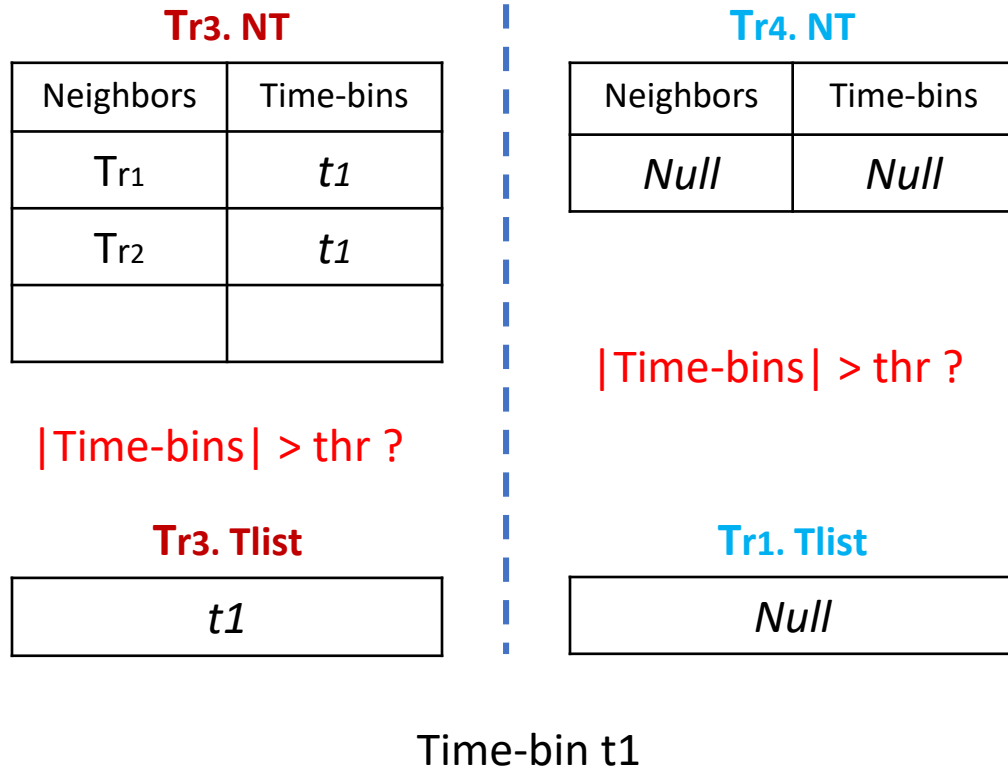
$K = 3$



Point Neighbor
($d < \epsilon$)



ODMTS Example



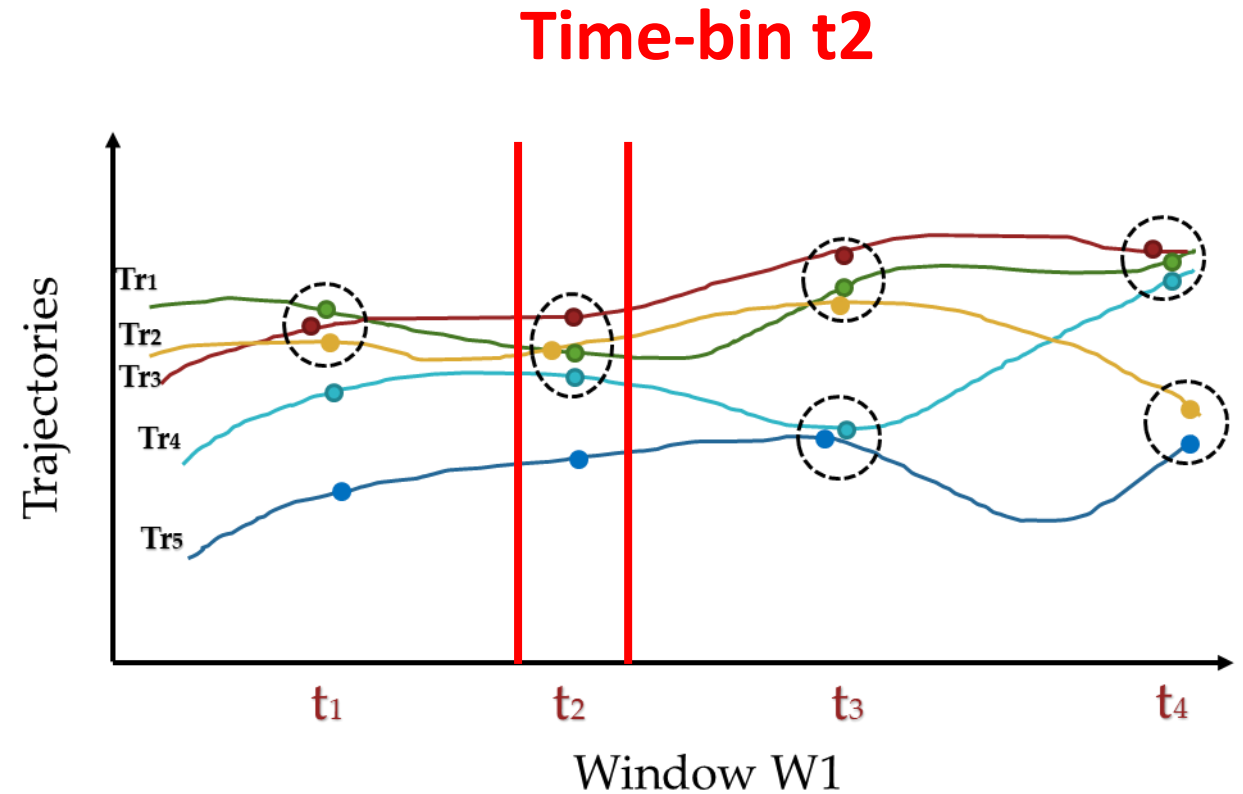
ODMTS Example

Tr3. NT		Tr4. NT	
Neighbors	Time-bins	Neighbors	Time-bins
Tr1	$t1, t2$	Tr1	$t2$
Tr2	$t1, t2$	Tr2	$t2$
Tr4	$t2$	Tr3	$t2$

Tr3. Tlist		Tr4. Tlist	
$t1, t2$		$t2$	

Time-bin $t2$

$t1, t2$



ODMTS Example

Tr3. NT

Neighbors	Time-bins
Tr ₁	t_1, t_2, t_3
Tr ₂	t_1, t_2, t_3
Tr ₄	t_2

$|\text{Time-bins}| > \text{thr} ?$

Tr3. Tlist

t_1, t_2, t_3

Tr4. NT

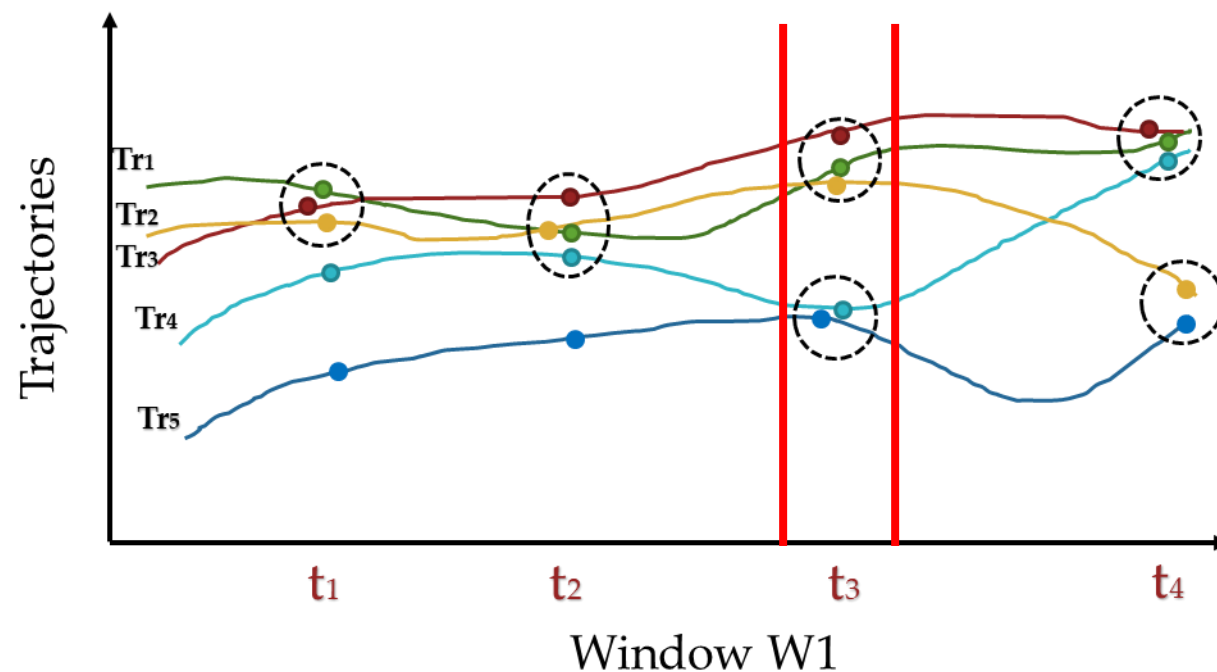
Neighbors	Time-bins
Tr ₁	t_2
Tr ₂	t_2
Tr ₃	t_2
Tr ₅	t_3

Tr4. Tlist

t_2

Time-bin t_3

Time-bin t_3



ODMTS Example

Tr3. NT

Neighbors	Time-bins
Tr ₁	t_1, t_2, t_3, t_4
Tr ₂	t_1, t_2, t_3
Tr ₄	t_2, t_4

$|\text{Time-bins}| > \text{thr} ?$

Tr3. Tlist

t_1, t_2, t_3, t_4

Tr4. NT

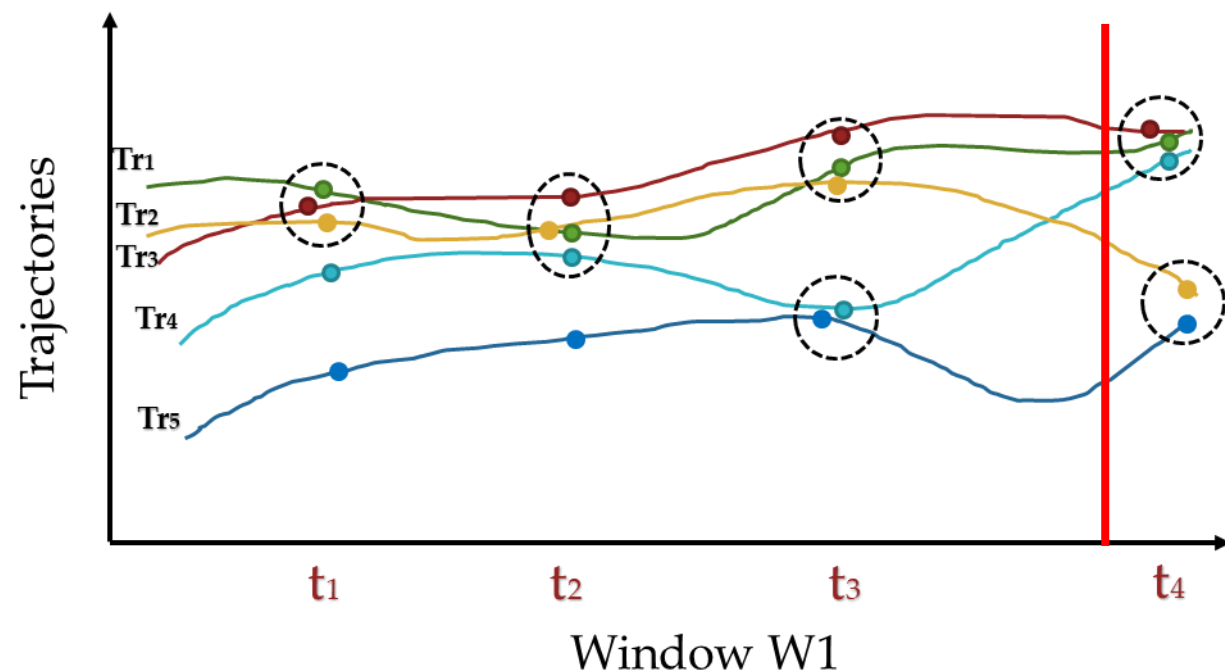
Neighbors	Time-bins
Tr ₁	t_2, t_4
Tr ₂	t_2
Tr ₃	t_2, t_4
Tr ₅	t_3

Tr4. Tlist

t_2, t_4

Time-bin t_4

Time-bin t_4



ODMTS Example

- Check each entry of $Tr_i.NT$:

Are $|time-bins| \geq Thr_j$?

- If at least K entries fulfill the above condition, Trajectory is inlier.
 - Else otherwise
-
- Hence, Tr_3 is an inlier and Tr_4 is an outlier

Outlier Test

Tr3. NT		Tr4. NT	
Neighbors	Time-bins	Neighbors	Time-bins
Tr1	$t1, t2, t3, t4$	Tr1	$t2, t4$
Tr2	$t1, t2, t3$	Tr2	$t2$
Tr4	$t2, t4$	Tr3	$t2, t4$
		Tr5	$t3$

Tr3. Tlist	Tr4. Tlist
$t1, t2, t3, t4$	$t2, t4$

Time-bin $t4$

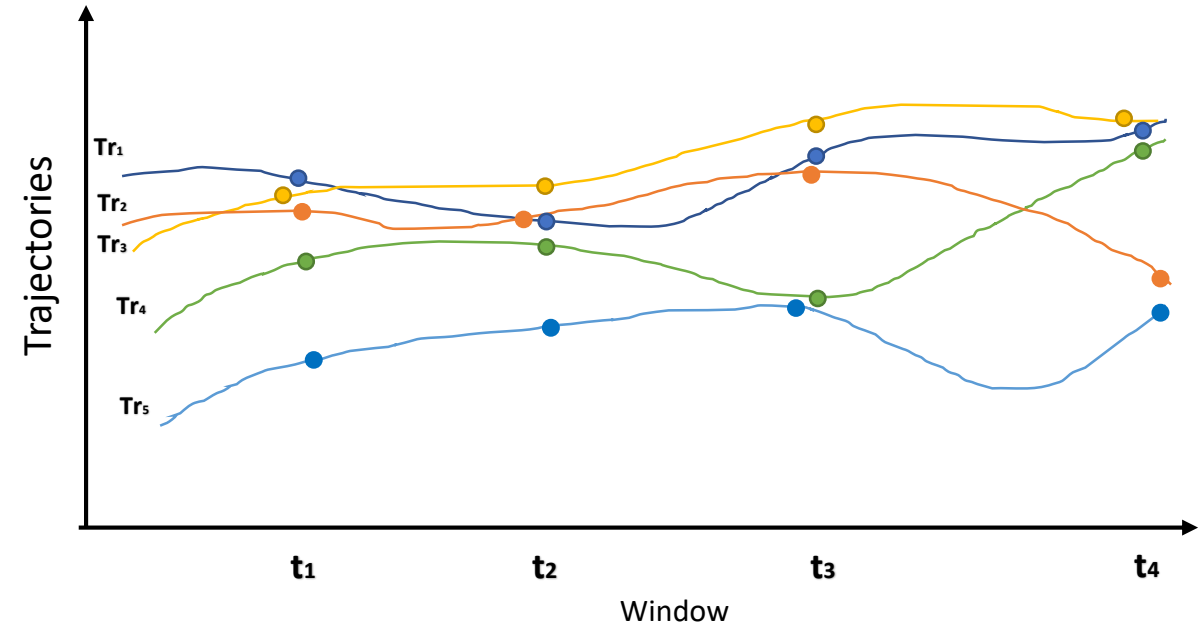
Pseudocode of ODMTS

Algorithm 1 Distance Based Outlier Detection (ODMTS)

Input Set of Trajectories, parameters: d , k , thr_i

Output Trajectory Outliers

```
1: for each  $TR_i$  do
2:   for each  $TR_k$  do
3:     if  $dist(p_i^j, p_k^j) < d$  then
4:        $TR_i.NT.insert(TR_k)$ 
5:     end if
6:   end for
7:   if  $TR_i.Count(t_{bin}) > thr$  then
8:      $TR_i.Tlist.insert(t_{bin})$ 
9:   end if
10:  if  $TR_i.size < k$  then
11:     $TR_i$  is an Outlier
12:  end if
13: end for
```



Experimental Setup

Hardware

- 40 cores
- 512 GB RAM

Performance Measures

- Total Execution Time

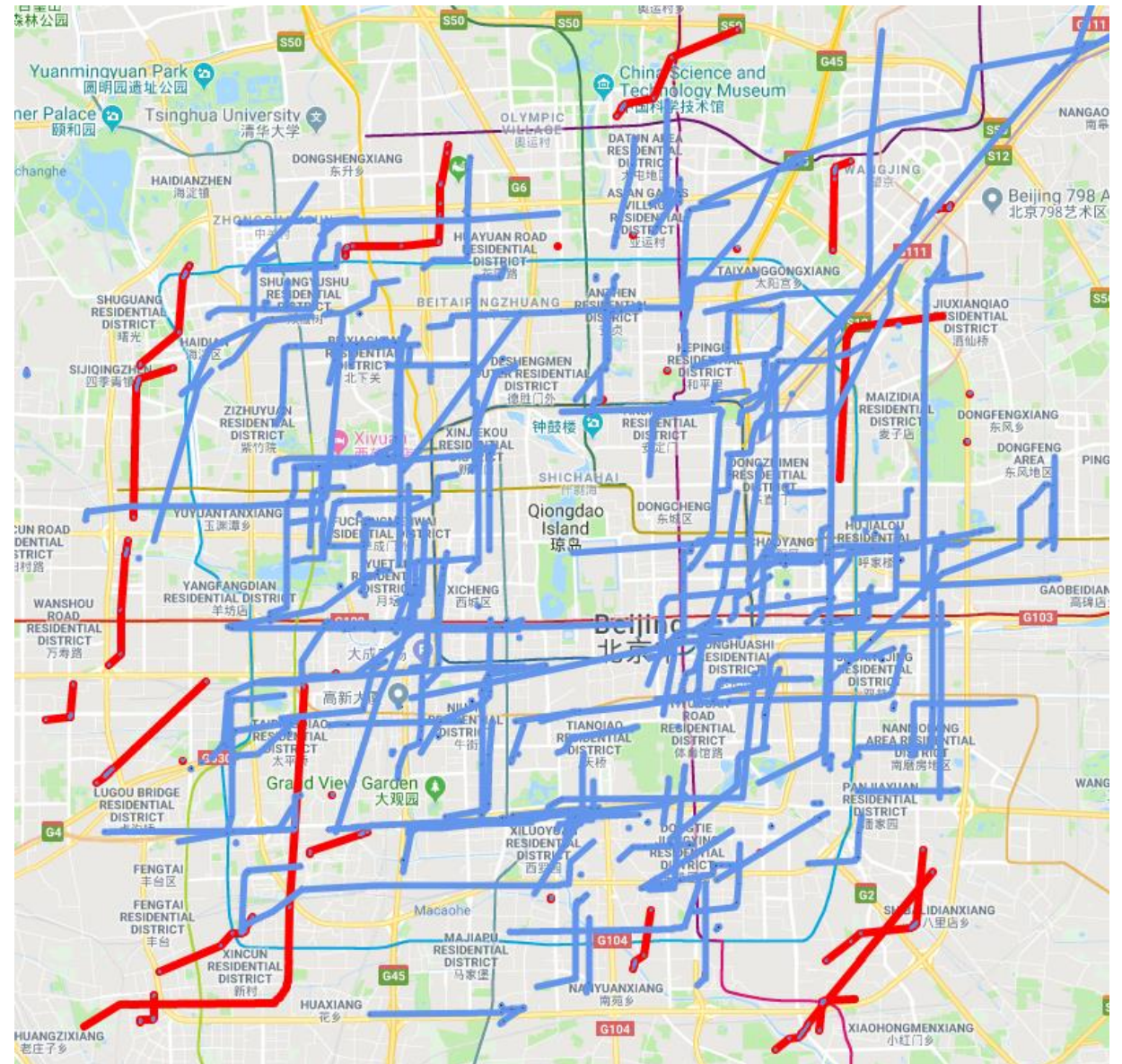
Software

- Python

Datasets

Dataset	No. of Trajectories	No. of Points (million)	No. of Attr.	Duration	Distance	Area
T-Drive	10357	15	2	7 days	9 million	Beijing
Geolife	17621	23.6	7	5 years	1.2 million	Beijing
Porto	1710671	-	9	1 year	-	Porto

Result of ODMTS



Outline

- Introduction
- Literature Review
- **Methodology**
- **Results**
- Conclusion & Future Work

Approach 1: Use Spatial Data Structures for Speed

Goal & Motivation:

- Use k-d trees and r-trees to improve running performance of ODMTS
- Compare the performances for both the trees
- Good for multi-dimensional points as in our case (latitude, longitude geo-coordinates)
- Our initial experiments revealed search query to be the slowest
- Suitable for range query search in large datasets



Key Idea Behind Approach 1

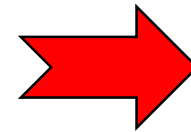
- Evaluated using python's cProfiler

Algorithm 1 Distance Based Outlier Detection (ODMTS)

Input Set of Trajectories, parameters: d, k, thr_i

Output Trajectory Outliers

```
1: for each  $TR_i$  do
2:   for each  $TR_k$  do
3:     if  $dist(p_i^j, p_k^j) < d$  then
4:        $TR_i.NT.insert(TR_k)$ 
5:     end if
6:   end for
7:   if  $TR_i.Count(t_{bin}) > thr$  then
8:      $TR_i.Tlist.insert(t_{bin})$ 
9:   end if
10:  if  $TR_i.size < k$  then
11:     $TR_i$  is an Outlier
12:  end if
13: end for
```



Pseudocode

- No change in pseudocode except insertion and search query
- Insert all geo-coordinates in tree
- Use a range query to retrieve neighbors in a radius r

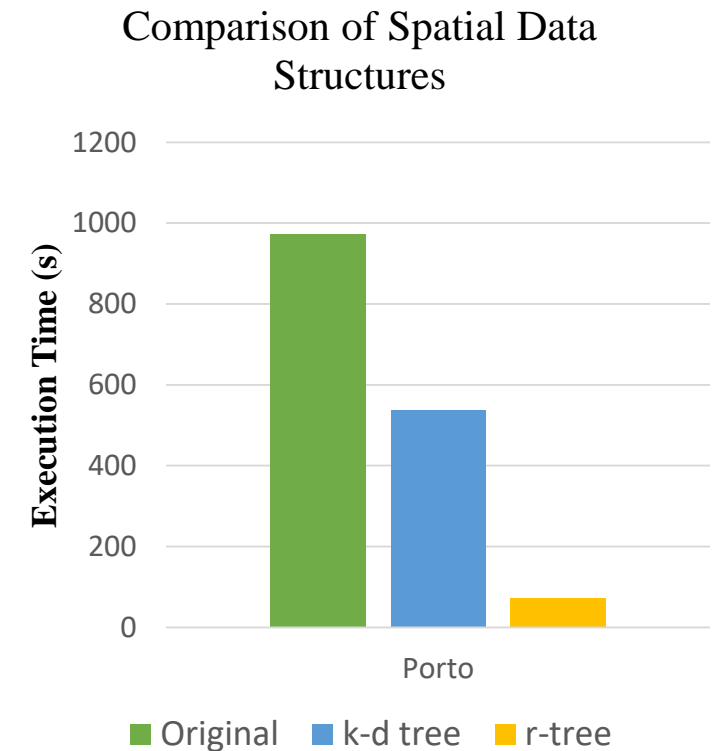
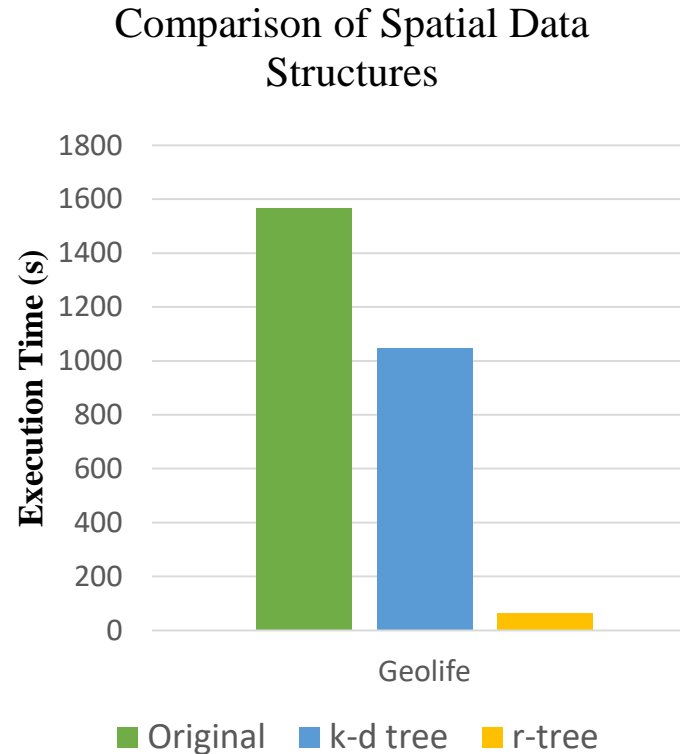
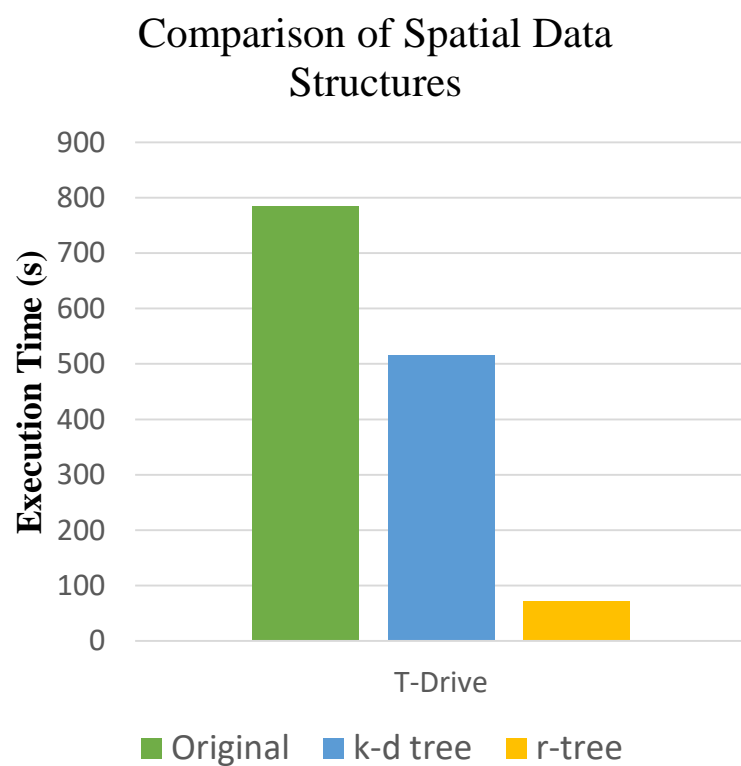
Algorithm 1 Distance Based Outlier Detection (ODMTS)

Input Set of Trajectories, parameters: d, k, thr_i

Output Trajectory Outliers

```
1: for each  $TR_i$  do
2:   for each  $TR_k$  do
3:     if  $dist(p_i^j, p_k^j) < d$  then
4:        $TR_i.NT.insert(TR_k)$ 
5:     end if
6:   end for
7:   if  $TR_i.Count(t_{bin}) > thr$  then
8:      $TR_i.Tlist.insert(t_{bin})$ 
9:   end if
10:  if  $TR_i.size < k$  then
11:     $TR_i$  is an Outlier
12:  end if
13: end for
```

Results of Approach 1



The experiments showed that using an R-tree improved the execution time performance of the ODMTS algorithm by 10x!!

Approach 2: Parallelization Strategy

Goal & Motivation:

- Propose parallelization strategy for ODMTS
- Study scalability in terms of execution time
- Study memory usage vs performance gain
- Our initial experiments revealed search query to be the slowest
- Harness the computing power of modern machines



Pseudocode

- Divide the workload among cores
- Single tree, passed to each core
- Each core performs search query on a subset of the dataset

Algorithm 2 Parallel - Distance Based Outlier Detection (ODMTS)

Input Set of Trajectories, parameters: d, k, thr_i

Output Trajectory Outliers

```
1: Insert all trajectory points in K-D/R Tree
2: Divide the Trajectory Dataset into Count(cores)
3: Insert divided Datasets in separate  $TR_iList$ 
4: for each core do
5:   Run a Ball-Point Query Search on all Trajectories in  $TR_iList$ 
6:   Insert query search result in  $Neighbor_iList$ 
7:   for each in  $Neighbor_iList$  do
8:     if  $Count(Neighbor_iList[TR_i]) < k$  then
9:        $TR_i$  is an Outlier
10:    end if
11:  end for
12: end for
```

Experimental Setup

Hardware

- Akka
- 40 cores
- 512 GB RAM

Performance Measures

- Total Execution Time

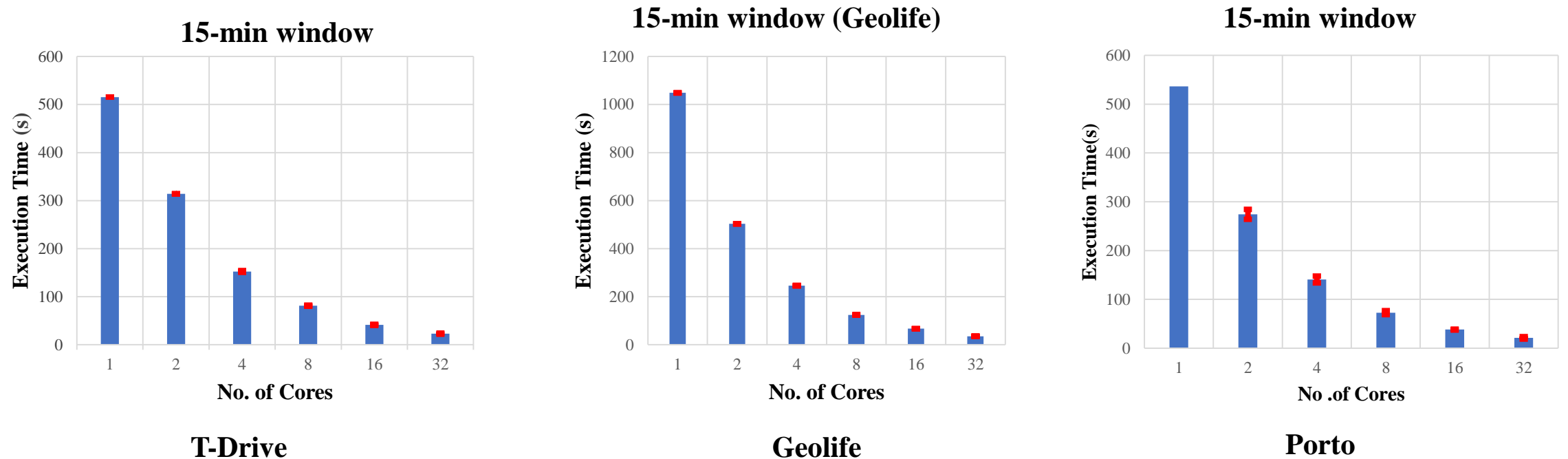
Software

- Python
- Multiprocessing

Datasets

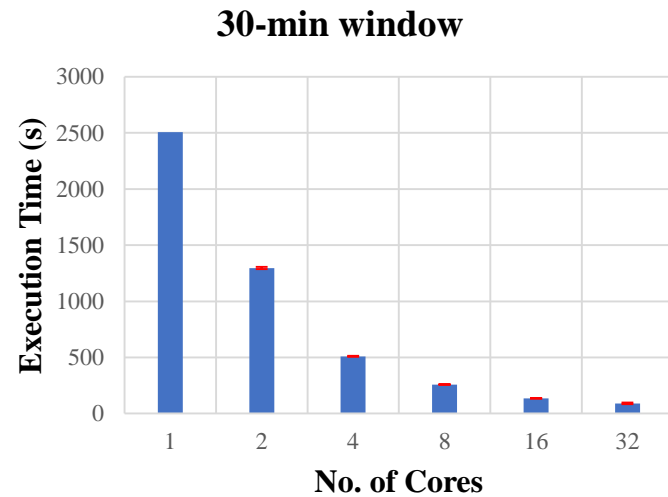
Dataset	No. of Trajectories	No. of Points (million)	No. of Attr.	Duration	Distance	Area
T-Drive	10357	15	2	7 days	9 million	Beijing
Geolife	17621	23.6	7	5 years	1.2 million	Beijing
Porto	1710671	-	9	1 year	-	Porto

Results of Parallelization Strategy

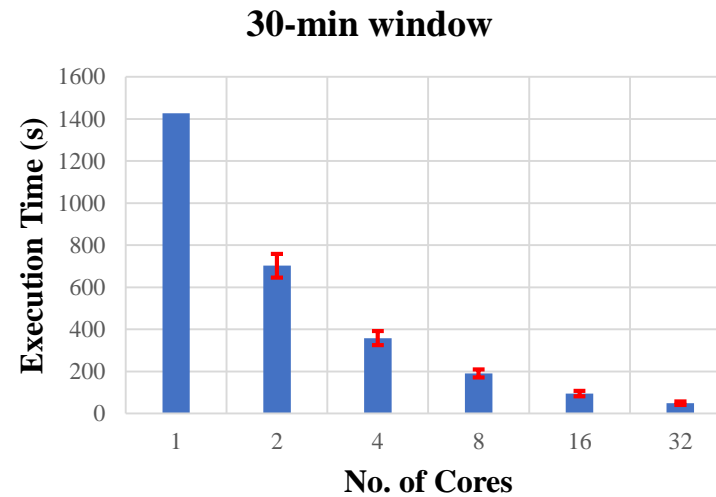


The experiments showed a linear decrease in execution time for k-d trees as the number of cores are increased!

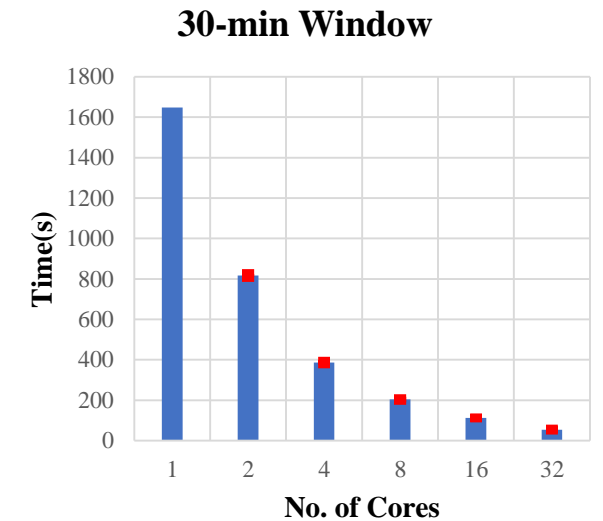
Results of Parallelization Strategy



T-Drive



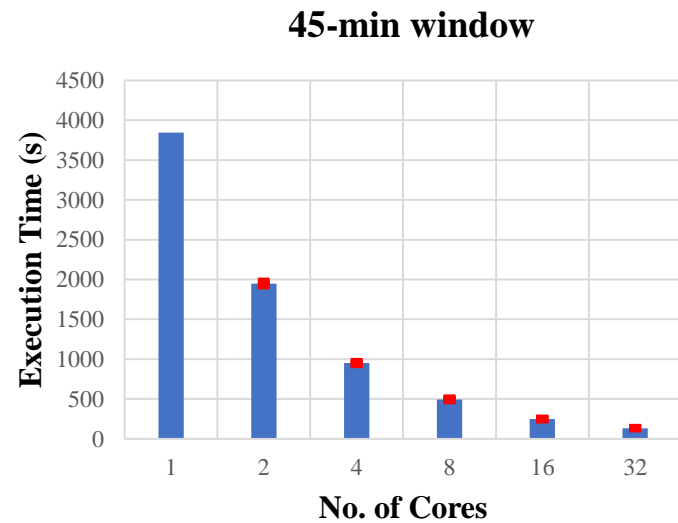
Geolife



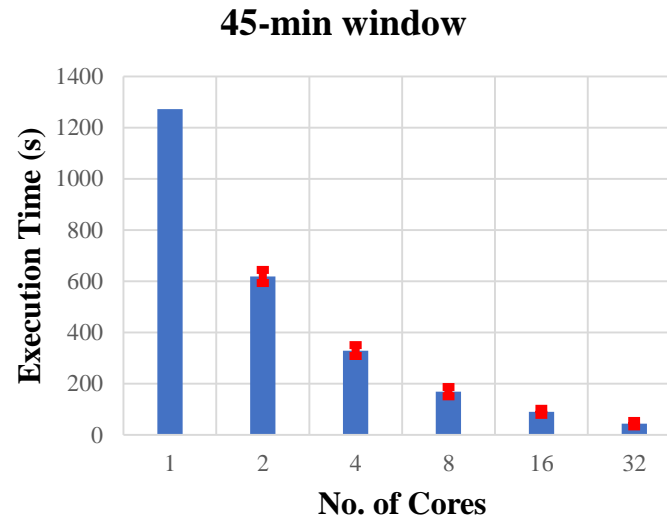
Porto

Workload is perfectly balanced among cores and the algorithm scales!!

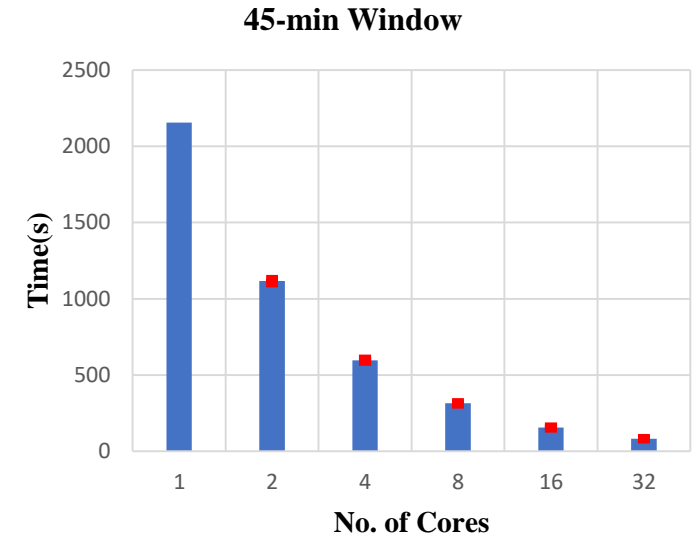
Results of Parallelization Strategy



T-Drive



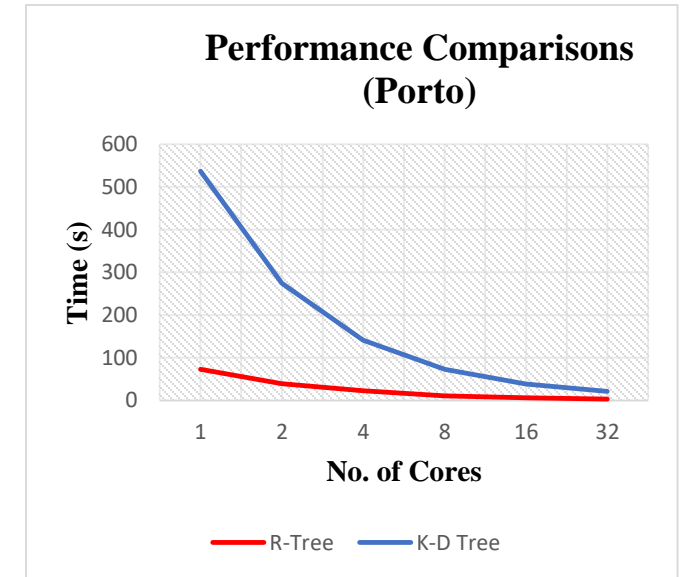
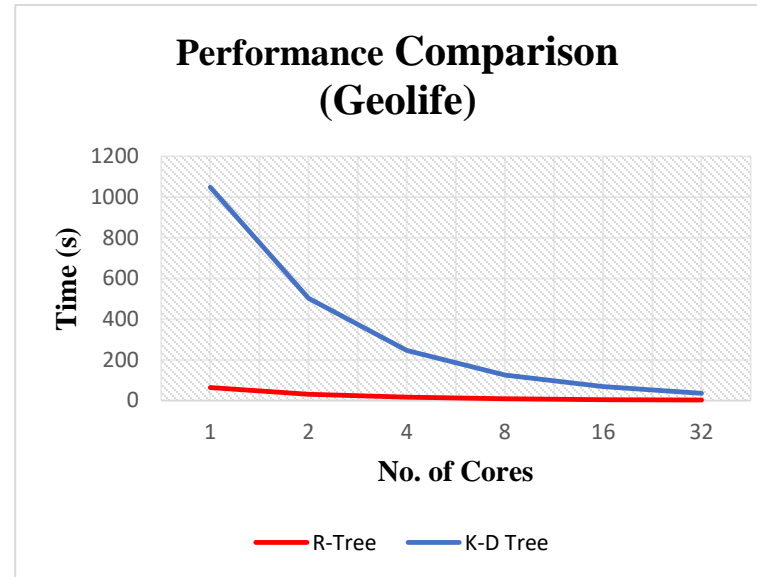
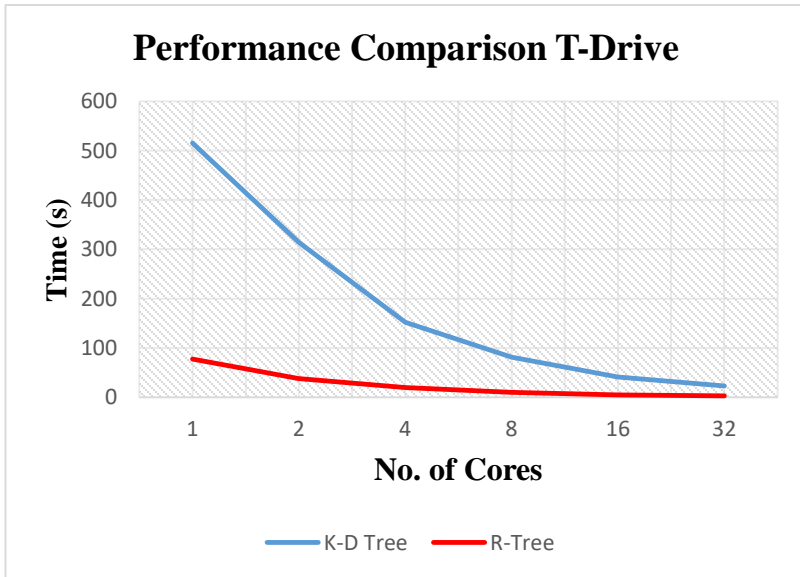
Geolife



Porto

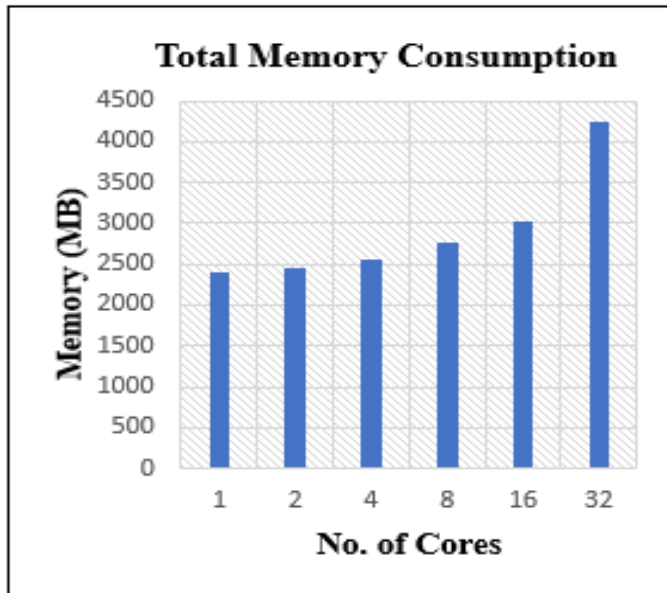
Scales linearly with the number of cores. In our case 32X decrease!

Results of Parallelization Strategy

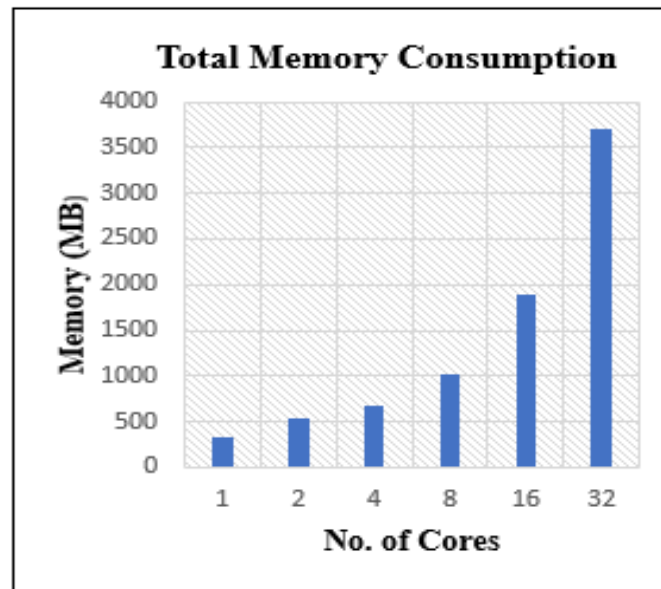


However, R-trees do not scale linearly in terms of execution times as the number of cores are increased, unlike k-d trees!!

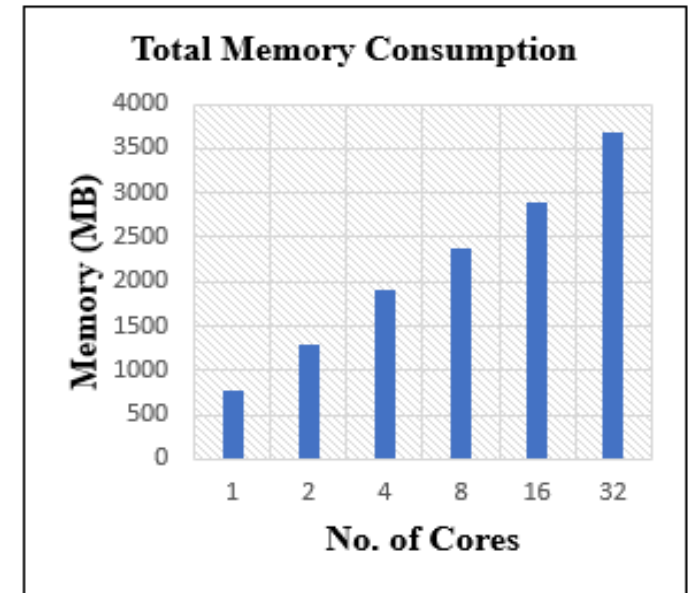
Results of Parallelization Strategy



T-Drive



Geolife



Porto

Total Memory Consumption shows an almost linear increase as number of cores are increased. This is because each core keeps a separate copy of the tree whilst having the same number of trajectory points!!

Approach 3: PDMTS

Goal & Motivation:

- Propose partition based streaming algorithm PDMTS
- Aim to detect trajectories that are significantly different but for a small period of the overall time of the trajectory
- Introduce windowing and temporal comparisons against other trajectories in the window
- Some outliers only hidden as sub-trajectories and not as a whole



Pseudocode

- Modifying TRAOD and Outlier Detection over Massive-Scale Trajectory Stream
- Detect outliers that are significantly different from the rest of the dataset but only for a small period of the trajectory
- Introduce windowing and temporal comparisons against other trajectories in the window
- First partition trajectory as in TRAOD then detect outliers using ODMTS
- Evaluate using Precision, Recall and F-Score

Algorithm 3 Streaming TRAOD (S-TRAOD)

Input Set of Trajectories, parameters: d, k, thr_i

Output Trajectory Outliers

```
1: —* Partitioning Phase *—
2: for each  $TR_i$  do
3:   Partition  $TR_i$  at coarse granularity using MDL ( $L_i$ )
4: end for
5: —* Detection Phase *—
6: for each partition  $L_i$  do
7:   for each partition  $L_j$  do
8:     if  $dist(p_i^j, p_k^j) < d$  then
9:        $TR_i.NT.insert(TR_k)$ 
10:    end if
11:  end for
12:  if  $TR_i.Count(t_{bin}) > thr$  then
13:     $TR_i.Tlist.insert(t_{bin})$ 
14:  end if
15:  if  $TR_i.size < k$  then
16:     $TR_i$  is an Outlier
17:  end if
18: end for
```

Experimental Setup

Hardware

- Akka
- 40 cores
- 512 GB RAM

Performance Measures

- Precision
- Recall
- F-Score

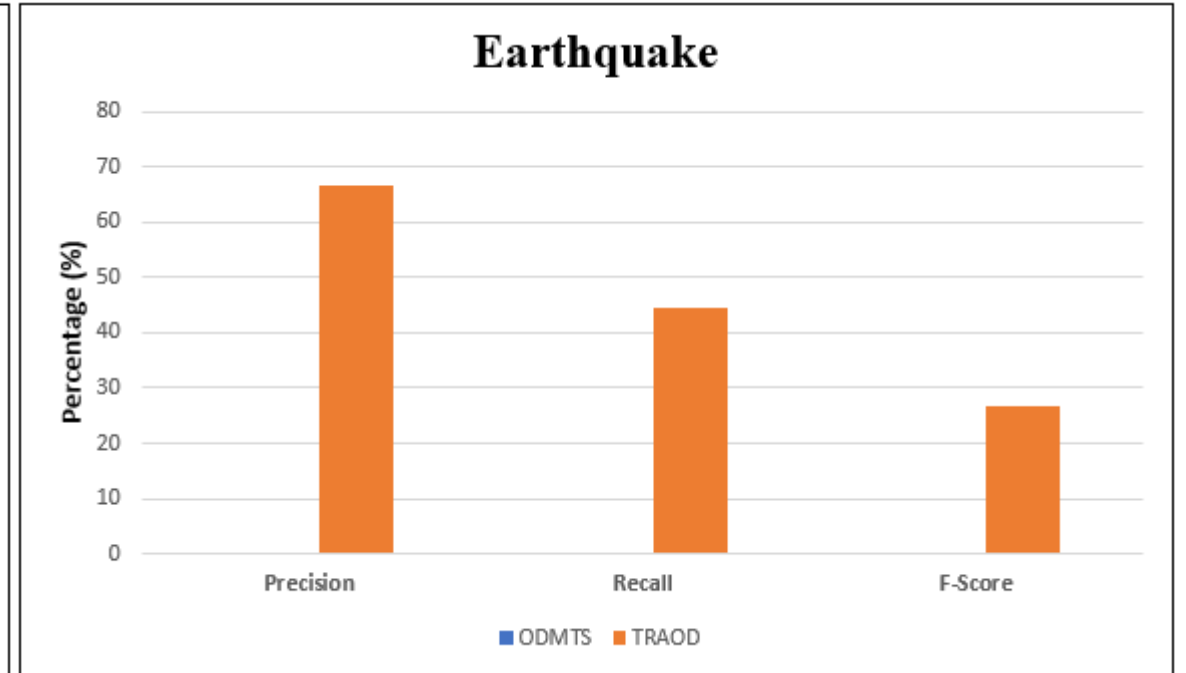
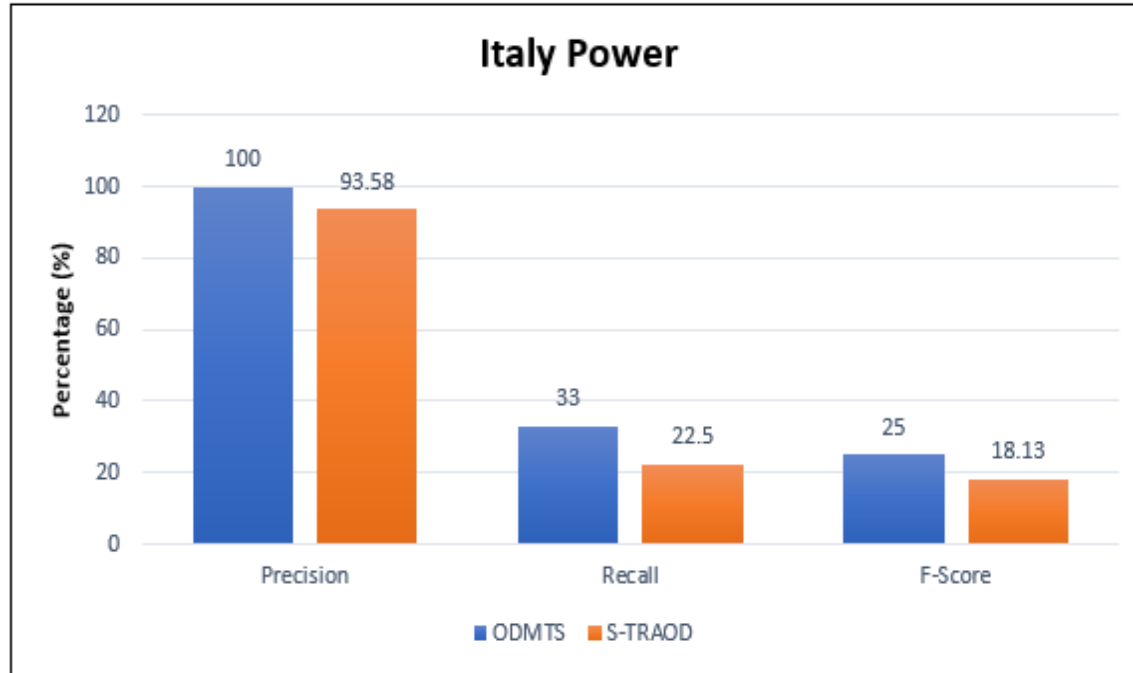
Software

- Python

Datasets

Dataset	No. of Trajectories	No. of Points (million)	No. of Attr.	Duration	Area
Earthquake	322	512	1	36 years	California
Italy Power	67	24	1	1 year	Italy

Results of Approach 3



ODMTS fails to detect any outlier on Earthquake dataset!!
Performs similarly on Italy Power!!

Outline

- Introduction
- Literature Review
- Methodology
- Results
- **Conclusion**

Conclusions for Approach 1

- Slowest block of our pseudocode was the range query for the neighbor search.
- To address this, we used k-d trees and r-trees, which are spatial data structures, to improve range queries
- The experiments showed that using an R-tree improved the execution time performance of the ODMTS algorithm by 10x compared to without them

References

- Edwin M. Knorr, Raymond T. Ng, and Vladimir Tucakov. 2000. Distance-based outliers: Algorithms and applications. *The VLDB Journal* 8, 3–4, 237–253
- Jae-Gil Lee, Jiawei Han, and Xiaolei Li. 2008. Trajectory outlier detection: A partition-and-detect framework. In *Proceedings of ICDE*. IEEE, 140–149
- Xiaolei Li, Jiawei Han, Sangkyum Kim, and Hector Gonzalez. 2007. ROAM: Rule- and motif-based anomaly detection in massive moving object data sets. In *Proceedings of SDM*. SIAM, 273–284
- Wei Liu, Yu Zheng, and Sanjay Chawla. 2011. Discovering spatio-temporal causal interactions in traffic data streams. In *Proceedings of SIGKDD*. ACM, 1010–1018
- Jing Yuan, Yu Zheng, Xing Xie, and Guangzhong Sun. 2013. T-Drive: Enhancing driving directions with taxi drivers' intelligence. *Transactions on Knowledge and Data Engineering* 25, 1, 220–232
- Yu Zheng, Xing Xie, and Wei-Ying Ma. 2010. GeoLife: A collaborative social networking service among user, location and trajectory. *IEEE Data Engineering Bulletin* 33, 2, 32–40
- Yingyi Bu, Lei Chen, Ada Wai-Chee Fu, and Dawei Liu. 2009. Efficient anomaly monitoring over moving object trajectory streams. In *Proceedings of SIGKDD*. ACM, 159–168

QUESTIONS?