

Advanced Kernel Clustering Technique To Space Related Networking Action Encapsulation

Abhinaya. N. A¹, PG Student¹, Rebecca. R², Asst. Professor,
Department of Computer Science and Engineering,

¹PG student, Parisutham Institute of Technology and Science, Thanjavur, TamilNadu, India.

²Asst. Professor, Parisutham Institute of Technology and Science, Thanjavur, TamilNadu, India.

Abstract—Data characterization is an important concept in data mining for identifying a compact representation of a dataset. In the spatial network activity summarization (SNAS), we are supposed to given a spatial network and a bulk of activities (e.g., pedestrian accidents reports, crime reports) and the aim is to find k shortest paths that summarize the activities. SNAS is more important for the applications where supervision occurs along linear paths such as roadways, rail tracks, etc. SNAS is computationally a try because of the large number of k subsets of shortest paths in a spatial network. Past work was centered on either geometry or sub graph-based techniques (e.g., only one single path), and cannot summarize events using multiple paths. This paper recommends a K-Main Routes (KMR) approach that develops k shortest paths to summarize events. KMR generalizes K-means for network space but uses shortest paths than ellipses to summarize activities. To enhance performance, KMR uses network Voronoi, separate and join, and trim strategies. Experimental results on synthetic and real data show that KMR with our performance-tuning decisions obtains substantial computational savings without losing summary path coverage.

I. INTRODUCTION

SPATIAL network activity summarization (SNAS) has its important applications in areas where monitoring occurs along linear paths in the network. This gives examples of such applications where the generative model of monitoring is inherently linear. For illustration, the transportation planners and engineers may have to find road segments/stretches that have harm for pedestrians and require redesign [1]; crime technicians will look for center of offence occurring along certain streets to lead law enforcement [2]; and hydrologists will attempt to summarize environmental changes on the water resources to make understand the action of river networks and lakes [3]. Casually, the SNAS problem can be defined as follows: Given a spatial network, a group of activities and their locations (e.g., placed on a node or an edge), and a given number of paths k , find a set of k shortest distance paths that enlarges the sum of activities on the paths (counting activities that are on overlapping paths only once) and a dividing of activities across the paths. Based on the area, an activity will be the location of a pedestrian fatality, a car accident, a train accident occurrences, etc. SNAS considers that each path is a shortest distance path because in applications such as locomotive transportation planning, the intension is usually to help the public to reach at their destination as soon as possible. Consider the input

consisting of eight nodes, seven edges (with edge weights of 1 for simplicity), eleven activities, and $k = 2$, mentioning that two routes and groups are desired. Table 1 clearly tells the shortest distance paths for the spatial network and their respective activity ranges (i.e., the sum of events on every shortest path). The result has two shortest paths and two groups of events. The shortest paths are the descriptions for every group and each shortest path extends the activity coverage for the group it mentions. For example, path (A, B, C) is the representative for the group consists of activities 1, 2, 3, 4, and 5, and path (D, E, F) is the representative for the group consists of activities 6, 7, 8, 9, 10, and 11.

Detecting the set of k shortest distance paths that extends the count of activities on selected route is computationally a try. This is because to the thing that if k shortest paths are selected among all shortest paths in a spatial network, there are a large number of chances for large k , i.e., (nk) , where n is the no of shortest distance route. This is due to various subsets of k shortest paths can overlap or have the same shortest routes. For disjoint paths, the consequences would be relatively less computationally challenging.

1) A Framework for Data Characterization

Data classification is a major concept in data mining that limits the techniques for detecting a together description or depiction of a dataset. The process usually includes defining a bulk of groups, finding a representative for every group, and investigating a statistic for each group (e.g., sum, mean, standard deviation). These notions vary based on the category of the data being classified. Table 2 presents a classification framework for three categories of data. An assumption of the first, relational table characterization, is the GROUP BY clause in SQL which is used to group rows in a table having common values to announce SQL aggregation functions such as mean and standard deviation. The group definition in this type of section is a division of rows and the group representation is having different values of attributes such as age-group, salary group, etc. The second category is spatial Euclidean characterization that adds heat maps and hotspot location analysis. Heat maps provide a graphical format of data in which individual values in a matrix are varied according to colors. The group definition for heat maps might add a set of pixels, and the group representation might be a subset of these pixels. Hotspots are a special kind of classification pattern where objects in hotspot regions have high resemblance in contrast to one another and are dissimilar to

all the objects that are outside the hotspot [6]. These spatial classifications are based on spatial point areas where the group definition is a classification of space and the groups could be represented by points, hexagons, pentagons, polygons, ellipses, or line-strings.

2) An Illustration on Application Domain: Preventing the Pedestrian Fatalities

To consider the occurrences of SNAS, we concentrate on the problem of pedestrian accidents. As per the recent policy report, more than 47,700 pedestrians were murdered in the United States from 2000 and 2009 [1]. More than 688,000 pedestrians were met with accidents over the same time period, which equal to a pedestrian is being hit by a vehicle each 7 minutes. Pedestrian accidents have risen in many of the places, including 15 of the country's largest metropolitan areas, even the overall pedestrian deaths have fallen [1]. Domain experts are trying to attribute pedestrian accidents largely to the maker of streets, which have been engineered for fastening traffic with little or no provision for people on foot, in wheel of chairs or on bicycles [1]. Day to day activities have moved away from main streets directions higher speed based on the special importance on traffic movement. This has occurred in most of the fatal pedestrian crashes happening on this wide, high capability and high-speed. They lack sidewalks such as platforms, crosswalks (or have crosswalks spaced too far apart), street lights, and school and public bus shelters [4].

II. RELATED WORK

Classifying the activities by grouping is a great research field in the area of data mining. Past approaches have generally been geometry-based [8]–[12] or network-based [13]. In geometry-based characterization, classification of spatial data is on the basis of grouping same points distributed in planar space where distance is measured with the help of Euclidean distance and not the network distance. Those approaches concentrates on the discovery of the geometry (e.g., circle, ellipse) of high concentration regions [5] and include K-Means [8], K-medoid [9], [10], P-median [11] and Nearest Neighbor Hierarchical Clustering [12]. These methods does not mind the upcoming spatial network; they group spatial objects that are very close in the form of Euclidean distance but not close in terms of the network path distance. Thus, they may be unsuccessful to group the activities that occur on the same street. In network-based characterization, the spatial objects are gathered using network (e.g., road) distance. Present methods of network-based characterization such as Mean Streets [7], Maximal Sub graph Finding (MSGF) [14], and Clumping [15] group activities over multiple paths, a single path/sub graph, or no paths at all. Mean Streets [13] finds new streets or routes with abnormal high activity levels. It is not designed to categorize the activities

over k paths because the number of high offence streets returned is always relatively small. MSGF [14] identifies the maximal sub graph (e.g., a single path, $k = 1$) with the consideration of a user defined length. The Network-Based Variable-Distance Clumping Method (NT-VCM) [22] is an illustration of the clumping approach [15]. NT-VCM groups activities that are within a certain shortest path distance of each other on the network; in order to run NT-VCM, a distance threshold is required [16].

Network distance might also be applied to few of the geometry-based approaches such as K-Means [8]. Fig. 3(c) shows an example of the ellipses output by K-Means using network distance. The left ellipse joins the activities 1, 2, 3, 6, 7, 8, and 11 where the right ellipse joins the activities 4, 5, 9, and 10. Though it is generalized with network distances, these methods result point-based or ellipsoid-based groups, not route. The main aim of network based K-Means (e.g., reduces the within-cluster sum of squares) is vary from that of SNAS (e.g., enlarges activity coverage), which explains their uniqueness in output.

III. BASIC CONCEPT

We discussed our basic concepts as follows:

Definition 1. A **spatial network** $G = (N, E)$ contains of a node set N and an edge set E , where every element u in N is related with a pair of real numbers (x, y) mentioning the spatial location of the node in a Euclidean plane [26]. Edge set E is a subset of the cross product $N \times N$. Each element $e = (u, v)$ in E is an edge that merges node u to node v . An illustration of a spatial network is shown in Fig. 1(a). In the figure, circles indicate nodes and lines indicate edges. A roadway route network is an example of a spatial network where nodes indicate street intersections and edges indicates streets.

Definition 2. An **activity set** A is a collection of activities. An **activity** $a \in A$ is an object of interest merged with only one edge $e \in E$ or one node $n \in N$. In Fig. 1(a), activities are mentioned as squares. In transportation planning, an activity may be the location area of a pedestrian accident; in offence analysis, an activity may be the location area of a theft; and in disaster reply an activity might be the location of a request for relief supplying purpose.

This work concentrates on characterizing discrete activity events (e.g., pedestrian fatalities, crime reports) related with a point on a network. This does not affect that all events must compulsorily be related with a point in a street. Further, the other network properties such as GPS strategies and traffic concentration of road networks are not considered. The main function used in SNAS is based on extending the activity coverage of summary paths, not on reducing the distance of events to summary paths.

SOURCE	SINK	SHORTEST PATH	ACTIVITY COVERAGE	SOURCE	SINK	SHORTEST PATH	ACTIVITY COVERAGE
A	B	$\langle A, B \rangle$	3	C	E	$\langle C, B, E \rangle$	2
A	C	$\langle A, B, C \rangle$	5	C	F	$\langle C, B, E, F \rangle$	4
A	D	$\langle A, B, E, D \rangle$	6	C	G	$\langle C, B, E, G \rangle$	3
A	E	$\langle A, B, E \rangle$	3	C	H	$\langle C, B, E, G, H \rangle$	3
A	F	$\langle A, B, E, F \rangle$	5	D	E	$\langle D, E \rangle$	3
A	G	$\langle A, B, E, G \rangle$	4	D	F	$\langle D, E, F \rangle$	5
A	H	$\langle A, B, E, G, H \rangle$	4	D	G	$\langle D, E, G \rangle$	4
B	C	$\langle B, C \rangle$	2	D	H	$\langle D, E, G, H \rangle$	4
B	D	$\langle B, E, D \rangle$	3	E	F	$\langle E, F \rangle$	2
B	E	$\langle B, E \rangle$	0	E	G	$\langle E, G \rangle$	1
B	F	$\langle B, E, F \rangle$	2	E	H	$\langle E, G, H \rangle$	1
B	G	$\langle B, E, G \rangle$	1	F	G	$\langle F, E, G \rangle$	3
B	H	$\langle B, E, G, H \rangle$	1	F	H	$\langle F, E, G, H \rangle$	3
C	D	$\langle C, B, E, D \rangle$	5	G	H	$\langle G, H \rangle$	0

Table 1: Description of shortest path(Activity coverage refers to the no of activities on a path)

Data Genre	Group Definition (Partitioning Criteria)	Group Representation Choices	Statistic
Relational Table (a set of rows)	a partition of rows	Distinct values of attributes (e.g., age-group)	sum, count, mean, etc.
Spatial (Euclidean Space)	a partition of space	points, polygons, ellipses, line-strings	sum, count, mean, etc.
Spatial Network (Neighbor Relationship)	a partition of a graph	node, path, tree, subgraph	sum, count, mean, etc.

Table 2: Describes the characterization framework for various data genres

IV. SPATIAL NETWORK ACTIVITY SUMMARIZATION

This section mentions the computational model of SNAS. It also explains the K-Main Routes (KMR) algorithm and its performance-tuning planning Network Voronoi activity allotment, Divide and join Summary route RE computation, and passive Node trimming.

1. Structure of SNAS

In SNAS, the best solution may not be different. Among the best solutions there is some path that starts and ends at active nodes. These approaches are formally shown.

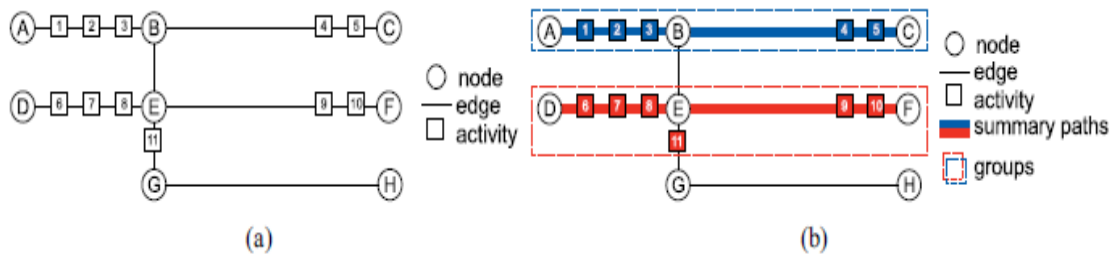


Fig.1(a)Input and(b) output of spatial network activity summarization(SNAS).

2. K-Main Routes Algorithm

Algorithm 1 represents the pseudo code for the past K-Main Routes (KMR) approach. The basic model of KMR seems that of K-Means [8] in terms of selecting initial nodes, forming k groups, and modernizing the representative of every group till the assignments no longer changes. Line 1 of Algorithm 1 all shortest path P that begin and end with active nodes (inactive node trimming).Then, k paths from P

Definition 3. An **active edge** is an edge $e \in E$ that has 1 or more events. An **active node** is a node u connect by an active edge or a node that has one or more events, or both. An **inactive node** is a node that is disconnected by any active edges. Edges A,B and B,C in the Fig. 1(a) are active because they each have at least one active edge and nodes A, B, C, D, E, F, and finally G are all active because they are all grouped by active edges. By opposite, Node H is an inactive because it is disconnected by any active edges.

are chosen as beginning summary paths, which are the "initiator" for KMR (line 2). This algorithm then further executes in two main phases. First, it creates k groups by

allotting each event to its nearest summary path (line 4). Then, it modernizes the summary path of each group by measuring the shortest route that enlarges activity coverage (line 5). Allocating and updating repeats till the summary

paths will never change and the final conclusion summary paths and joined that are returned (line 8).

Let's see the algorithm below:

Algorithm 1 K-Main Routes approach(KMR) Algorithm

Input:

- 1) a spatial network $G = (N, E)$,
- 2) a set of activities A ,
- 3) a number of routes k ,
- 4) mode1 $\in \{naive, NOVA_TKDE\}$,
- 5) mode2 $\in \{naive, D-SPARE_TKDE\}$

Output:

A summary path is a set of size k and a dividing of activities through these summary paths, where the utilization is to enlarge the event coverage of each summary path for the group.

Algorithm:

- 1: $P \leftarrow$ shortest paths between active nodes of G
- 2: $P \leftarrow k$ summary paths $\in P$; stable Groups \leftarrow false;
- 3: **while** not stable Groups **do**
- 4: **Phase 1:** current Groups \leftarrow Assign Activities-
 To Summary Paths($G, A, k, P, mode1$)
- 5: **Phase 2:** $P \leftarrow$ Recompute Summary Paths
 ($G, A, k, current Groups, mode2$)
- 6: **if** $P' = P$ **then** stable Groups \leftarrow true
- 7: $P' \leftarrow P$
- 8: **return** currentGroups

This is the algorithm used for finding the shortest path

V. HYPOTHETICAL ANALYSIS

In this section, we represent a proof of NP-Completeness for spatial network activity summarization. We also represent the proofs of exact validness for our new performance-tuning approaches.

1. Proof of NP-Completeness

For easy way, we start by defining a generalized, conclusion version of SNAS where all the set of routes may be random and show this problem to be NP-complete. We then present a proof of rough sketch indicating the conclusion version of SNAS with shortest paths is also NP-complete.

VI. EXPERIMENTAL EVALUATION

The aim of our test was to analyze KMR with and without performance-tuning. Scalability was analyzed by differing and monitoring the effect of four workload criteria: edges, events, paths, and active node ratio.

Definition 4. The **active node ratio**, ANR, is the ratio of active nodes to all edges.

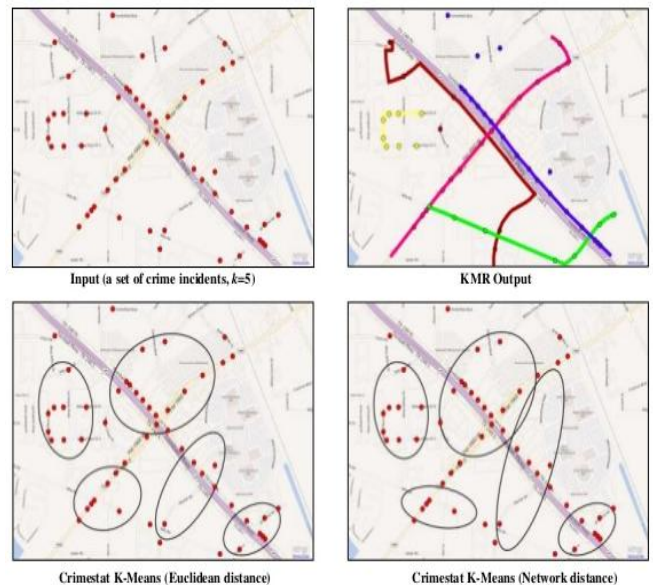
The active edges ratio in Fig. 1(a) is 7/8, that is the total no of active edges 7 is divided by the total no of edges, 8.

For every work done test, we ran seven versions of KMR:

- KMR with NOVA_TKDE only (KMR_V)
- KMR with D-SPARE_TKDE only (KMR_D)
- KMR with both NOVA_TKDE and D-SPARE_TKDE (KMR_{VD})
- KMR with inactive node pruning only (KMR_I) [23]
- KMR with both inactive node pruning and NOVA_TKDE (KMR_{IV})
- KMR with both inactive node pruning and D-SPARE_TKDE (KMR_{ID})
- KMR with all three performance-tuning conclusion (KMR_{IVD})

All tests were done on a Mac Pro with a 2 x Xeon Quad Core 2.26 GHz processor and 16 GB RAM. In all test, the summary paths were beginning to the top- k edge-disconnected shortest paths.

VII. PERFORMANCE EVALUATION



VIII. CONCLUSION

This work found out the defects of spatial network activity characterization in association to important application domains such as preventing pedestrian accidents and crime analysis. We proposed a K-Main Routes (KMR) algorithm that creates a set of k shortest paths to group activities. KMR uses passive node trimming, Network Voronoi activity Assignment (NOVA_TKDE) and Divide and join Summary route Recompile (D-SPARE_TKDE) to improve its performance and scalability. We presented a study on comparing KMR with other characterization encapsulation approaches on pedestrian accident data. Experimental evaluation using both synthetic and real-world data sets indicate that the performance-tuning decisions utilized by KMR obtained substantial computational savings without reducing the coverage of the resulting summary paths.

ACKNOWLEDGEMENT

My sincere thanks to my guide Mrs.R.Rebecca Asst.Professor, HOD, Department of Computer Science and Engineering, Parisutham Institute of Technology and Science, Thanjavur for her help and guidance.

REFERENCES

- [1] M. Ernst, M. Lang, and S. Davis. (2011). Dangerous by design: Solving the epidemic of preventable pedestrian deaths. Transportation for America: Surface Transportation Policy Partnership. Washington, DC, USA. [Online]. Available: <http://trid.trb.org/view.aspx?id=1148931>
- [2] J. Eck, S. Chainey, J. Cameron, M. Leitner, and R. Wilson. (2005, Aug.). Mapping Crime: Understanding Hot Spots, U.S. Department of Justice. Washington, DC, USA [Online]. Available: <https://www.ncjrs.gov/pdffiles1/nij/209393.pdf>
- [3] D. Matthews, S. Effler, C. Driscoll, S. O'Donnell, and C. Matthews, "Electron budgets for the hypolimnion of a recovering urban lake, 1989-2004: Response to changes in organic carbon deposition and availability of electron acceptors," *Limnol. Oceanogr.*, vol. 53, no. 2, pp. 743-759, 2008.
- [4] Huffington Post. (2013, Apr.). Hungary: Snowstorm Strands Thousands in Their Cars [Online]. Available: http://www.huffingtonpost.com/huff-wires/20130315/eu-europe-snow/?utm_hp_ref=travel&ir=travelM1
- [5] R. Wronski. (2013, Apr. 9). Metra argues for delay of 'fail-safe' rail system. Chicago Tribune [Online]. Available: <http://www.chicagotribune.com/news/local/ct-met-metra-collision-prevention-20130409,0,3710438.story>
- [6] S. Shekhar, M. Evans, J. Kang, and P. Mohan, "Identifying patterns in spatial information: A survey of methods," *WIREs Data Mining Knowl. Discov.*, vol. 1, no. 3, pp. 193-214, Apr. 2011.
- [7] Fatality Analysis Reporting System (FARS). Encyclopedia, National Highway Traffic Safety Administration (NHTSA) [Online]. Available: <http://www.nhtsa.gov/FARS>
- [8] J. MacQueen et al., "Some methods for classification and analysis of multivariate observations," in *Proc. 5th Berkeley Symp. Math. Statist. Probab.*, vol. 1. Berkeley, CA, USA, 1967, p. 14.
- [9] L. Kaufman and P. Rousseeuw, *Finding Groups in Data: An Introduction to Cluster Analysis*. Wiley Online Library, 1990.
- [10] R. Ng and J. Han, "Efficient and effective clustering methods for spatial data mining," in *Proc. 20th Int. Conf. VLDB*, San Francisco, CA, USA, 1994, pp. 144-155.
- [11] M. Resende and R. Werneck, "A hybrid heuristic for the p-median problem," *J. Heuristics*, vol. 10, no. 1, pp. 59-88, Jan. 2004.
- [12] R. D'Andrade, "U-statistic hierarchical clustering," *Psychometrika*, vol. 43, no. 1, pp. 59-67, Mar. 1978.
- [13] M. Celik, S. Shekhar, B. George, J. Rogers, and J. Shine, "Discovering and quantifying mean streets: A summary of results," Univ. Minnesota, Minneapolis, MN, USA, Tech. Rep. 07-025, 2007.
- [14] K. Buchin et al., "Detecting hotspots in geographic networks," in *Proc. Adv. GIScience*, Berlin, Germany, 2009, pp. 217-231.
- [15] S. Roach, *The Theory of Random Clumping*. London, U.K.: Methuen, 1968.
- [16] A. Okabe, K. Okunuki, and S. Shiode, "The SANET toolbox: New methods for network spatial analysis," *Trans. GIS*, vol. 10, no. 4, pp. 535-550, Jul. 2006.

AUTHOR DETAILS



N.A. Abhinaya Completed B.E., (CSE) at Periyar Maniammai University, Thanjavur in 2013. Now pursuing M.E., (CSE) at Parisutham Institute of Technology and Science, Thanjavur.



R.Rebecca Completed B.E., at Karunya University, Coimbatore. Completed M.Tech., in PRIST University and doing Ph.d in PRIST University, Vallam, Thanjavur. Working as Asst. Prof in Computer Science And Engineering Department in Parisutham Institute Technology Science, Thanjavur.