

AN AUTOMATIC SEMANTIC MAP GENERATION METHOD USING TRAJECTORY DATA

Yu Miao¹, Xuehua Tang^{2*}, Zhongyuan Wang³

¹ School of Computer Science, Wuhan University, Wuhan, China - 2015302580003@whu.edu.cn

² School of Remote Sensing and Information Engineering, Wuhan University, Wuhan, China - tangxuehua@whu.edu.cn

³ School of Computer Science, Wuhan University, Wuhan, China - wzy_hope@163.com

KEY WORDS: Semantic map generation, Semantic pattern, Semantic recognition, Stay point, Trajectory data

ABSTRACT:

It's easily to obtain the geometric information of terrain features in a timely manner using advanced surveying and mapping methods, but it is impossible to obtain their semantic information with low latency due to the rapid development of cities. The popularity of GPS-enabled devices and technologies provide us a large number of personal location information. Moreover, it is possible to extract the personal or group behavior pattern due to the regularity of human behavior. Those conditions make it possible to extract and identify human behavior patterns from their trajectory data. In this paper, we present an automatic semantic map generation method that extract semantic patterns and take advantage of them to tagging spatial objects in an unknown region based on known semantic patterns. We study the regularity of trajectory data and build the semantic pattern based on the regularity of human behavior. Most importantly, we use known semantic patterns to identify the semantics of the stay points in the unknown region, and use this method to realize the semantic recognition of the stay points. Results of the experiments show the effectiveness of our proposed method.

1. INTRODUCTION

Due to the rapid development of cities, advanced surveying and mapping methods are only able to obtain the geometric information of terrain features in a timely manner, but it is impossible to obtain their semantic information with low latency. How to obtain map timely semantic information has become the bottleneck of map construction.

The popularity of sensor-enabled mobile phones and GPS tracking technology provides huge amount of human trajectory data that include lots of implicit information of human activities. Abundant researchers have engaged in the analysis and mining of trajectory data and shown that there is a particular spatiotemporal regularity in individual behavior patterns (Gonzalez et al., 2008; Rossi et al., 2015; Lv et al., 2016; Zhang et al., 2019). Trajectory data provides new opportunities for gaining knowledge about individual movement behavior and social interactions (Kwan, 2004, 2013; Raubal et al., 2004; Griffith et al., 2013; Palmer et al., 2013). Based on the regularity of personal behavior, large amount of mobility trajectory data about individuals has been generated and collected, which has been applied in the field of traffic, geographic, health, and social science (Shoval et al., 2011; Wesolowski et al., 2012; Richardson et al., 2013; Shen et al., 2013). Therefore, it is feasible to explore the semantics of maps by using the behavior characteristics contained in personal trajectory data. Existing map construction based on personal trajectory data is mainly applied to indoor maps (Zhang et al., 2019; Zhou et al., 2019; Liu et al., 2019), but has not been applied to the semantic recognition of outdoor maps. In response to this need, this paper proposes an automatic semantic map generation method based on mobile phone trajectory data. The new method uses the trajectory data sets of the semantic region of the known map to extract the individual or group behavior semantic pattern. Based on the semantic pattern, we match the known patterns with the trajectory data of the unknown regions to identify the map semantics of these regions. The new method includes three major steps: 1) Investigate and define the semantic patterns using trajectory data based on the spatiotemporal semantic information of stay points; 2) Extract the stay points from the trajectory data in unknown semantic regions; 3) Matching the known semantic patterns with the stay point

clusters in the regions without semantics to obtain semantic annotation of map.

According to this process, the remainder of this paper is organized as follows: Section2 introduces the method for the extraction of the stay point and the semantics of stay point. Section3 describes the definition and extraction of the semantic pattern-based trajectory data. Section4 describes a method of automatic map semantic recognition based trajectory data. Section5 shows the experimental evaluation on the real trajectory datasets. Section6 summarize our main work and contribution.

2. STAY POINT CLUSTER AND SEMANTICS EXTRACTION

To build individual or group behavior semantic patterns, the first step is to extract the stay points from the original trajectory data, and then match the stay points with the spatial objects of map to extract the semantic information.

2.1 Stay Point Cluster

The first stage transforms trajectory data into a series of stay point clusters, which are locations where the users stay for a period of time rather than just passing by. Semantic information of trajectory data mainly concentrated in the stay points where users stopped for a certain purpose. Extracting stay points is essential to get the semantic patterns.

For each trajectory $T = \{t_1, t_2, \dots, t_i, \dots, t_n\}$, where t_i is the i -th point in the T trajectory, which contains the latitude, the longitude and time information. We use a spatiotemporal threshold method to get the stay point. Time threshold THR_time is the minimum duration of staying point cluster. Distance threshold $THR_distance$ is the maximum distance between all trajectory points in the staying point cluster and the initial point of the point cluster.

The steps to extract staying point clusters from trajectory data are as follow: 1) Use the next trajectory point as the initial point for the stay point cluster; 2) Add the subsequent trajectory points to the stay point cluster until the distance between the next point with the initial point is greater than the distance threshold $THR_distance$; 3) Regard the point cluster as a staying point cluster if the duration is greater than the time threshold

THR_time. The algorithm is shown as Algorithm 1 for the pseudo-code.

Algorithm 1 Stay point cluster detection

INPUT: Trajectory data T , time threshold THR_time, distance threshold THR_distance;
OUTPUT: Stay point cluster set $C = \{c_1, c_2, \dots, c_i, \dots, c_n\}$

- 1: Initialize empty point set $c = \{\}$ and stay point cluster set $C = \{\}$
- 2: for all p in T do
- 3: if c is empty or distance $(p_{i-1}, p_i) < \text{THR_distance}$ then
- 4: add p_i to c
- 5: else
- 6: if time_duration(c) $> \text{THR_time}$ then
- 7: add c to C
- 8: end if
- 9: $c = \{\}$
- 10: end for

2.2 Stay point semantics extraction

To extract the stay point semantics, we firstly match the stay point clusters to the spatial objects of known map. This step is most fundamental for both semantic pattern extraction and automatic semantic identification. Since the extraction result of staying points is a set of point clusters, to extract the semantic information of the staying points, we must first match each point in the set of point clusters with the map. Considering the large number of data at the staying points, we adopted the method of hierarchical matching in order to ensure the efficiency and accuracy of matching at the same time. This hierarchical matching has three stages. Firstly, the Geohash matching is used to preliminary index the spatial objects in a certain range around the stay point and narrow the list of candidates. Then we use the bounding rectangle of the spatial objects to match the remaining points. Finally, Pnpoly algorithm is applied to implement the accurate matching. Based on the matching results of individual staying points, the matching result of each staying point cluster is calculated according to the length of staying.

2.2.1 Geohash matching: In order to filter spatial objects in a certain range around the stay point and shorten the list of candidates, we adopt Geohash method to our matching method. Geohash method uses bisection along longitude and latitude to grid the earth and use 0 or 1 to represent the region generated by each division. When bisection is conducted along the longitude direction, the code of the left region is 0, and the code of the right region is 1; when bisection is conducted along the latitude direction, the code of the lower region is 0, and the code of the upper region is 1. Every five divisions divided origin region into 32 grids with 5bit binary code. Grids are finally represented by base32 code, a number of 0-9 or a lowercase English character except a, i, l, o, converted from the original binary code. Continuous gridding can get higher accuracy and longer Geohash code. Figure 1 shows the procedure of Geohash code.

However, the stay points may not be in the center of grids generated by Geohash method. Geohash method only make sure that the stay points are within the grid, which may cause false result. For instance, as shown in Figure 2, when two points are close to each other physically, Geohash code gets opposite result. Geohash code shows that C is closer to A than B while A and B are divided into two grids. In order to deal with this problem, we simultaneously use the grid where the staying point is located and the eight adjacent grids for subsequent matching at the next step instead of single grid.

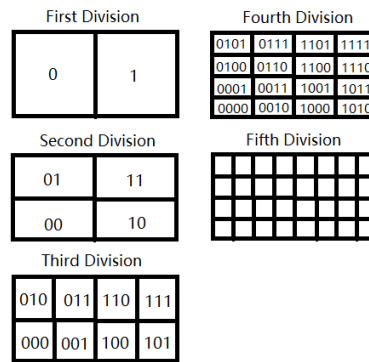


Figure 1. Procedure of Geohash code

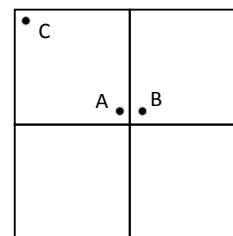


Figure 2. False result of Geohash Method

2.2.2 Bounding rectangle matching: Bounding rectangle is a rectangle of the spatial object, which represents the range of latitude and longitude of a spatial object. The bounding rectangle could be built by the maximum and minimum value of the objects' coordinates. Based on the result of the Geohash index, we adopt bounding rectangles for further matching. When the stay point falls into the rectangle of the spatial object, it is most likely to fall into this object. Bounding rectangle matching further narrows the search scope.

2.2.3 Pnpoly algorithm matching: The bounding rectangle is the approximate range of the spatial object. If a stay point falls within the rectangles of multiple objects, more accurate judgment is required. For example, as shown in Figure 3, a stay point will be attached to an object surrounding it through bounding match while it is out of the object. In order to solve this problem, we employ Pnpoly algorithm, which is a classic algorithm for judging the point within a polygon, to implement the final decision of stay point matching.

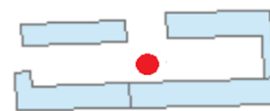


Figure 3. Wrong judgment of bounding matching

2.2.4 Stay point cluster matching: After the matching of individual staying points is completed, the stay point cluster matching is calculated according to the staying time. The spatial object where users stay for longest time is regarded as the matching result to this stay point cluster.

Algorithm 2 for the pseudo-code shows the whole procedure of matching. Based on the stay point cluster set $C = \{c_1, c_2, \dots, c_i, \dots, c_n\}$ and the spatial object set BS, there are four steps to match the stay point cluster with the object of C. In the first stage, we use Geohash to filter out the candidate object set CB within a certain range around the trajectory point p_j from the

spatial object sets BS. The second stage is to further match by the bounding rectangles of the spatial objects. Then the stay point is precisely matched to an object using Pnpoly algorithm. The last stage is to determine the attributions and semantics of staying point clusters according to the length of staying and get the stay point clusters by semantics. Experimental results show that even in the case of mobile GPS signal position shift and dense objects, our method still has a high accuracy. As the amount of candidate objects obtained by Geohash is often only a few hundred at most, the hierarchical matching method get both high efficiency and high accuracy.

Algorithm 2 Place matching method

INPUT: Stay point cluster set $C = \{c_1, c_2, \dots, c_i, \dots, c_n\}$, Spatial object set BS
OUTPUT: Stay point cluster set with corresponding objects $C' = \{c'_1, c'_2, \dots, c'_i, \dots, c'_n\}$
 1: **for** all stay point cluster c_i in C **do**
 2: **for** all trajectory point p_j in c_i **do**
 3: Detect all candidate objects CB for trajectory point p_j from BS by geohash
 4: **for** all candidate object cb_k in CB **do**
 5: **if** trajectory point p_j is within the bounding box of candidate object cb_k **then**
 6: **if** trajectory point p_j is in the region of candidate object cb_k **then**
 7: candidate object cb_k is the corresponding object to trajectory point p_j
 8: **end if**
 9: **end if**
 10: **end for**
 11: **end for**
 12: Choose the object user stay for the longest time in as the corresponding one to stay point cluster c_i , and add object id to it.
 13: **end for**

3. SEMANTIC PATTERN DEFINITION AND EXTRACTION

3.1 Semantic pattern

In order to represent the spatiotemporal information of the stay points from the trajectory data, the semantic pattern is defined as a place preference matrix M as:

$$M = \begin{bmatrix} X_1 \\ \vdots \\ X_m \end{bmatrix} \quad (1)$$

Where X = place preference vector of a certain type of person in a certain time interval

Assuming the time interval is one hour, there are 24 periods corresponding to 24 hours. Then, M includes 24 place preference vectors (X_1, \dots, X_{24}) . Considering typical daily activity such as dining usually has a time period for about 15 minutes, we adopt 15 minutes as the length of time interval, where $m=96$.

To achieve the integration of different matrix, we integrated the original vector into the total stay time and translated the stay time for each object to the stay probability. Therefore, the place preference of the i -th period X_i is defined as:

$$X_i = (x_{i1}, x_{i2}, \dots, x_{in}, T_i) \quad (2)$$

Where n = the number of semantic types

T_i = the sum of the length of stay time of the i -th object

x_{ij} = the probability of staying in the type j -th object
 x_{ij} is defined as:

$$x_{ij} = \frac{t_{ij}}{T_i} \quad (3)$$

Where t_{ij} = the staying time for the j -th object of the i -th period

3.2 Semantic pattern fusion

Due to the data collection, etc., a person or a class of people may has multiple behavior semantic patterns. In order to solve this problem, we define pattern fusion to integrate multiple patterns. Based on the definition of semantic pattern, the essence of pattern fusion is the fusion of two preference matrices. Assumed that $X_i = [x_{i1}, x_{i2}, \dots, x_{in}, T_{xi}]$ and $Y = [y_{i1}, y_{i2}, \dots, y_{in}, T_{yi}]$ are two place preference vectors for the same time interval i , the fusion of the two place preference vectors $Z = [z_{i1}, z_{i2}, \dots, z_{in}, T_{zi}]$ is defined as:

$$z_{ij} = \frac{x_{ij} * T_{xi} + y_{ij} * T_{yi}}{T_{xi} + T_{yi}}, 1 \leq j \leq n \quad (4)$$

Where $T_{zi} = T_{xi} + T_{yi}$

3.3 Semantic pattern extraction

As mentioned before, using the results of staying point cluster matching with the spatial objects of known map, we can extract the semantic information of the trajectory and construct its semantic pattern matrix. To obtain the semantic pattern, we assign semantic feature to stay point clusters according to the spatial objects' information.

Firstly, we set the location preference vector $X_j (1 \leq j \leq m)$ for each time period in the semantic model $M = [X_1, X_2, \dots, X_j, \dots, X_m]$ to the initial value $[0, \dots, 0]$. Corresponding to the semantic model period, the staying point cluster is divided into multiple time slices $\{p_1, p_2, \dots, p_j, \dots, p_l\}$ and each time slice is assigned a location preference vector $Y_j = [0, \dots, 1, \dots, 0, t]$, where t is the duration of the time slice p_j , the value 1 represents the object type of the staying point cluster. The final semantic pattern can be obtained by fusing the location preference vector Y of all time slices with the location preference vector X of the corresponding time period in the semantic pattern. See Algorithm 3 for the pseudo code.

Algorithm 3 Semantic pattern extraction

INPUT: Stay point cluster set with corresponding objects $C' = \{c'_1, c'_2, \dots, c'_i, \dots, c'_n\}$

OUTPUT: Semantic pattern $M = [X_1, X_2, \dots, X_j, \dots, X_m]$

1: Initialize semantic pattern $M = [X_1, X_2, \dots, X_j, \dots, X_m]$ with a certain time interval, where $X_j = [0, 0, \dots, 0, 0], 1 \leq j \leq m$
 2: **for** all stay point cluster c'_i in stay point cluster set C' **do**
 3: Add semantic mark to stay point cluster c'_i according to its corresponding object
 4: Divide the time of stay point cluster c'_i into parts $\{p_1, p_2, \dots, p_j, \dots, p_l\}$ that match to the place preference vectors in M
 5: **for** all parts p_j in stay point cluster c'_i **do**
 6: initialize a place preference vector $Y = [0, 0, \dots, 0, 0]$
 7: $Y = [0, \dots, 1, \dots, 0, t]$, where t is the length of the time of p_j , and 1 is corresponding to the semantic mark of c'_i
 8: Use semantic pattern fusion to fuse Y and its corresponding place preference vector in M
 9: **end for**
 10: **end for**

4. AUTOMATIC SEMANTIC MAP GENERATION

4.1 Semantic identification of unknown regions

If people with known semantic patterns are active in an unknown region, we can use the known semantic patterns to identify the semantic information of spatial objects and generate semantic map based on their trajectory in the unknown region. Firstly, we divide the duration of any staying point cluster with semantics c'_i into several time slices in the same time interval as the semantic pattern $M = [X_1, X_2, \dots, X_j, \dots, X_m]$. For each slice, we extract its semantic information from the corresponding place preference vectors of semantic pattern. Then we set a location preference vector $Y_j = [y_{j1}, y_{j2}, y_{j3}, \dots, y_{jn}, t_j]$ for each time slice p_j , where t_j is the duration of time slice, y_{ji} is the value corresponding to the i -th object in the location preference vector X in the semantic pattern M that is the same as the time slice p_j . By fusing a stay point cluster c'_i for all time-sliced location preference vectors, we can get a probability distribution vector, which represents the probability of a stay point cluster according to a certain object. Finally, the final probability of the semantic type to the object can be obtained calculated by fusing all probability distribution vectors belonging to the same object in the stay point set C' . See Algorithm 4 for the pseudo code.

Algorithm 4 Automatic Semantic Map Generation

INPUT: Semantic pattern $M = [X_1, X_2, \dots, X_j, X_m]$, Stay point cluster set with corresponding objects $C' = \{c'_1, c'_2, \dots, c'_i, \dots, c'_n\}$

OUTPUT: Objects with semantic mark in an unknown region

- 1: **for** all stay point cluster c'_i in stay point cluster set C' **do**
 - 2: Divide the time of stay point cluster c'_i into parts $\{p_1, p_2, \dots, p_j, \dots, p_l\}$ that match to the place preference vectors in M
 - 3: **for** all parts p_j in stay point cluster c'_i **do**
 - 4: initialize a place preference vector $Y_j = [0, 0, \dots, 0, 0]$ for p_j
 - 5: $Y_j = [y_{j1}, y_{j2}, y_{j3}, \dots, y_{jn}, t_j]$, where t_j is the length of the time of p_j , and y_{ji} is the value of i -th object type in the corresponding place preference vector of p_j
 - 6: **end for**
 - 7: Fuse place preference vectors for all parts of c'_i through semantic pattern fusion to get a probability vector $Y_i = [y_{i1}, y_{i2}, y_{i3}, \dots, y_{in}, t_i]$ for c'_i , where t_i is the time of stay in c'_i and y_{ik} is the probability that the corresponding object of c'_i belongs to the k -th object type
 - 8: **end for**
 - 9: Fuse probability vectors of the same object and get the final probability vectors of the corresponding object set of the stay point cluster set
-

4.2 Activity balance weight

We find that defining the semantic type with the longest stay time may lead to bias while the activity corresponding to the semantic type takes much more time than another. For instance, dining usually takes only 15 minutes in canteen while studying always takes a couple of hours in teaching object. In this situation, the spatial object of canteen may be wrongly recognized as teaching objects. Therefore, we adopt activity balance weight to eliminate the bias. An activity balance weight is the reciprocal of the average duration of the activity corresponding to the semantic mark. When detect the semantic mark of a spatial object, the

probability of each semantic mark should be multiplied by the corresponding activity balance weight.

5. EXPERIMENTAL EVALUATION

In this section, we conduct experiments on real mobility data sets to verify the feasibility and performance of our method.

5.1 Dataset

Our data set is collected by the volunteers from Wuhan University who shared location data using our mobile app. The data source includes the trajectory data of 6 volunteers for 4 months and the map of Wuhan University. To avoid sudden changes in semantic patterns, we mainly collected non-holiday trajectory data of student volunteers.

5.2 Experimental results and discussion

To verify the validity of our method, we designed and implemented two experiments.

The first experiment is to apply one person's trajectory data to extract semantic patterns, and then use this semantic pattern to match another person's trajectory data to extract the related semantics in the behavior trajectories. The experimental result is shown as Table 1.

Parameters Setting	Number of Objects	Accuracy Rate
15 minutes with activity balance weight	108044	58.54%
One hour with activity balance weight	108044	57.81%
15 minutes without activity balance weight	108044	50.53%
One hour without activity balance weight	108044	50.51%

Table 1. The result of Experiment 1

The second experiment used the same person's semantic pattern to predict the semantics of the region where he was active in another period of time. The experimental result is shown as Table 2. Obviously, the recognition accuracy of Experiment 2 is higher than that of Experiment 1, this is because the semantic pattern of the same person is more stable than that of the same class of person.

Parameters Setting	Number of Objects	Accuracy Rate
15 minutes with activity balance weight	2710	70.0%
One hour with activity balance weight	2710	64.7%
15 minutes without activity balance weight	2710	58.6%
One hour without activity balance weight	2710	57.5%

Table 2. The result of Experiment 2

Experimental results show that activity balance weights can significantly improve accuracy. In addition, the time period of 15 minutes can achieve better accuracy than the time period of 1 hour. The reason is that using a short time interval, like 15 minutes, has more accurate description of each time interval. In contrast, a long-time interval, like one hour, has high

generalization performance to contain an activity which can happen in a large time range. For example, if a semantic pattern with time interval of 15 minutes shows that a person only eats lunch from 12:00 to 12:15, we can't get right semantic mark with that semantic pattern when the trajectory data shows that another person eats lunch out of that time range.

The existing experiment results have limited recognition accuracy, which is mainly due to the limited trajectory data used in the experiment and the limited accuracy of the extracted behavior patterns. The identification effort is to increase the accuracy while increasing the number of trajectory data to obtain more precise semantic pattern. In addition, the recognition accuracy is also related to the density of features in the area to be recognized. The accuracy of feature recognition in low-density areas is higher than that in high-density areas.

6. CONCLUSION

In this paper, we study the method of identifying the semantics of unknown spatial objects based on personal trajectory data. To achieve this goal, we firstly present a hierarchical matching strategy to match the semantics of spatial objects with stay point clusters. Then we formulate the semantic pattern for personal trajectory data and adopt a fusion method for different patterns for the same person or group. Finally, we combine the known semantic pattern and the trajectory information of unknown region to realizing the recognition of semantics. The experimental results show the feasibility and practicality of the new method. Undoubtedly, this study provides an innovative and effective way of semantics' recognition and automatic semantic map generation using trajectory data.

ACKNOWLEDGEMENTS

This work was supported by National Science Foundation, China (41701518).

REFERENCES

- González M. C., Hidalgo C. A., and Barabási A. L., 2008. Understanding individual human mobility patterns. *Nature*, 453(7196): 779-782.
- Griffith, D. A., Y. Chun, M. E. O'Kelly, B. J. L. Berry, R. P. Haining, and M. P. Kwan., 2013. *Geographical Analysis: Its First Forty Years*. *Geographical Analysis* 45(1), 1–27.
- Hongmin Liu, Xincheng Tang, and Shuhan Shen, 2019. Depth-map Completion for Large Indoor Scene Reconstruction. *Pattern Recognition(X)*, X.
- Krogh B., Andersen O., Lewis-Kelham E., Pelekis N., Theodoridis Y., and Torp K., 2013. Trajectory based traffic analysis. In *Proceedings of the 21st ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems*. ACM, 536–539.
- Kwan M.P., 2004. GIS Methods in Time-Geographic Research: Geocomputation and Geovisualization of Human Activity Patterns. *Geografiska Annaler B*, 86, 267–80.
- Kwan, M. P., 2013. *Beyond Space (As We Knew It): Toward Temporally Integrated Geographies of Segregation, Health, and Accessibility*. *Annals of the Association of American Geographers* 103(5), 1078–86.
- Lv M., Chen L., Xu Z. et al., 2016. The discovery of personally semantic places based on trajectory data mining. *Neurocomputing*, 173(JAN.15PT.3): 1142-1153.
- Raubal M., H. J. Miller, and S. Bridwell., 2004. User-Centred Time Geography for Location-Based Services. *Geografiska Annaler B*. 86, 245–65.
- Richardson, D. B., N. D. Volkow, M.P. Kwan, R. M. Kaplan, M. F. Goodchild, and R. T. Croyle., 2013. Spatial Turn in Health Research. *Science*, 339(6126), 1390–92.
- Rossi L., Walker J., and Musolesi M., 2015. Spatio-temporal techniques for user identification by means of GPS mobility data. *Epj Data Science*, 4(1):11.
- Shen, Y., M.P. Kwan, and Y. Chai., 2013. Investigating Commuting Flexibility with GPS Data and 3D Geovisualizations: A Case Study of Beijing, China. *Journal of Transport Geography*, 32, 1–11.
- Shoval, N., H.-W. Wahl, G. Auslander, M. Isaacson, F. Oswald, Landau, and J. Heinik., 2011. Use of the Global Positioning System to Measure the Out-of-Home Mobility of Older Adults with Differing Cognitive Functioning. *Ageing & Society*, 31, 849–69.
- Wesolowski, A., N. Eagle, A. J. Tatem, D. L. Smith, A. M. Noor, and R. W. Snow., 2012. Quantifying the Impact of Human Mobility on Malaria. *Science*, 338, 267–70.
- Zhang D. , Lee K. , and Lee I., 2019 . Semantic Periodic Pattern Mining from Spatio-temporal Trajectories. *Information ences*, 502, 164-189.
- Zhang W , Zhou S , Yang L , et al., 2019. WiFiMap+: High-Level Indoor Semantic Inference with WiFi Human Activity and Environment. *IEEE Transactions on Vehicular Technology*, PP(99):1-1.
- Zhou, B., Elbadry, M., Gao, R., and Ye, F., 2019. Towards scalable indoor map construction and refinement using acoustics on smartphones. *IEEE Transactions on Mobile Computing*, 1-1.