



UNIVERSITÀ DEGLI STUDI DI CATANIA
DOTTORATO IN INFORMATICA (INTERNAZIONALE)
XXXIII CICLO

Claudia Cavallaro

Analysis of large GPS trajectories datasets via multi-agent techniques

TESI DI DOTTORATO DI RICERCA

Relatore:
Chiar.mo Prof O. Muscato
Chiar.mo Prof E. Tramontana

Coordinatore: Chiar.mo Prof. S. Battiato
ANNO ACCADEMICO 2019/2020

Contents

Abstract	4
Acknowledgements	4
1 Introduction	8
1.1 Contributions	9
1.2 Aims and Approach	10
2 Related work	13
2.1 Trajectory monitoring	13
2.2 Extraction and recommendation of places of interest	16
3 Preliminaries and Basic Definitions	20
3.1 Geographic definitions	20
3.2 Other different distances used to compare trajectories	22
3.2.1 Fréchet Distance e DDF	25
3.3 Recommendation Systems and collaborative filtering	27
3.4 Multi-Agent System	28
3.5 Clustering techniques	28
3.5.1 K-Means	29
3.5.2 DBSCAN	30
3.6 Apriori algorithm	31
3.7 Sentiment Analysis	34
4 Problem Analysis	36
4.1 Proposed Datasets	36
4.1.1 GeoLife	36
4.1.2 Taxi	37
4.1.3 Truck	37
4.1.4 Cabspottingdata	37
4.2 Proposed Approach	38
4.2.1 StayPoints extraction	45
4.2.2 Solution for Points Of Interest detection	45

4.2.3	POIs Recommendation (proposed Multi-Agent system)	48
4.2.4	Privacy Preserving	51
4.2.5	People movements identification	51
4.2.6	Mining corridors from GPS trajectories, through Apriori	52
4.2.7	Probability of movement for GPS points	58
5	Evaluation	62
5.1	SPs and POIs obtained	62
5.2	Movement of people detected	69
5.3	Corrispondence between flows and POIs	72
5.4	Relations with the proposed multi-agents system	79
5.5	Test performed to investigate probability of users' movement	81
5.6	Different methods by varying the corridor parameters	90
5.6.1	Filtering process	93
6	Final discussion, remarks and future work	96
6.1	Epilogue	96
6.2	Conclusions	98
6.3	Open Issues	99
	Bibliography	100

*Un passo alla volta,
senza fretta ma,
senza sosta.*

Johann Wolfgang Goethe

Abstract

1 The widespread use of mobile devices and Internet services on users' locations offer the oppor-
2 tunity to acquire information related to users' trips in real time. This availability has given
3 rise to several studies based on geospatial trajectories but, because of the large volume of the
4 collected information, processing them together is usually difficult. In this thesis we illustrate
5 an approach that, using a multi-agent system, provides personalized recommendations of Points
6 Of Interest (POIs), such as libraries, museums, restaurants, etc. to users. In our context, an
7 agent is an application that improves user navigation in a city. It collects opinions, in terms
8 of scores, that quantify the level of satisfaction in visiting a certain place in a certain period
9 of time. In this approach, interesting positions emerge from the analysis of the collected data,
10 hence scores and suggestions may be available for any large city in any place, when a sufficient
11 number of people provide data. In addition, the next places to visit are suggested according to
12 people's behavior and preferences.

13 Other directions explored are the identification of the flows of multiple users and the intention
14 to predict the paths that will be taken by a user on the basis of the common paths, already known,
15 of other individuals. Given a large dataset of geographic trajectories in an urban metropolitan
16 area, an efficient strategy for detecting corridors is also proposed. These can be defined as
17 geographical paths, of a minimum length, commonly crossed by a minimum number of different
18 users. This approach is important for transportation analytics because it allows to detect missing
19 lines in public transport systems and also to advise private users which route to take to move
20 from one part of the city to another based on the behavior of users who have provided their
21 GPS logs.

22 Although people like to visit popular places, due to health problems and due to the recent
23 restrictions currently in place around the world for covid-19 influenza pandemic, meetings should
24 be avoided. When planning a trip, one must consider both the attractiveness in terms of general
25 interest for the destinations and the density of people who gather there. In the final part of the
26 thesis, we propose a recommendation system that aims to offer to users some suggestions on
27 useful routes and destinations which balance liveliness and overcrowding. First, we use datasets
28 that store GPS locations as the basis for route and destination statistics. Then, we use an
29 accurate probability algorithm that estimates the number of people moving from one place of
30 the city to another and consequently show a list of destinations to users. Destination points are
31 filtered according to the user's preferences on the density of people. A multi-agent system is
32 used to manage user's requests to find a route for a journey, statistics on possible destinations

33 and suggestions for users. Thanks to our solution, we can inform users about suitable routes
34 and destinations, as well as alert them when a favorite destination is overcrowded.

Acknowledgements

35 I would like to thank my Ph.D. supervisors Professors Emiliano Tramontana and Orazio Muscato
36 for their work and their presence. Professor Tramontana has guided and supported me over the
37 years, encouraging my choices. Professor Muscato has given me precious scientific advices, and
38 he has passed on to me the passion for research since my master.

39 In these years I have had the privilege of carrying out a period of study and research at
40 the “Universitat de Barcelona (UB)”. For this I thank Professor Jordi Vitrià who accepted
41 me for visit in his department and guided me in research for 6 months and the “Departament
42 de Matemàtiques i Informàtica” of UB which provided me with a peaceful and stimulating
43 environment to work. This great experience has allowed me to learn many things, grasp new
44 aspects of scientific research and get to know new realities.

45 A sincere thanks to Professor Sebastiano Battiato for his important guide as Coordinator of
46 the PhD course, for his availability and for his humanity.

47 I would like to thank Professor Paolo Bellavista for his hospitality at the Department of
48 Computer Science and Engineering (DISI) of University of Bologna in the final period of my
49 PhD.

50 Thanks to my father, a fundamental presence in my life. Thanks for his patience, for his
51 encouragement and for understanding the importance that the study and the research has to
52 me.

53 A hug and a final thank goes to all those people I met at the department of Mathematics
54 and Computer Sciences of University of Catania, especially my colleagues Francesco S., Gabriella
55 V., Antonio S. and Fabio D. for their advices.

56 I acknowledge the support provided by my PhD grant PO FSE 2014 – 2020 founded by
57 Regione Siciliana.

Chapter 1

Introduction

58 The research field explored, during the PhD, is Geo-localized data analysis. This analysis has
59 become an important field of study due to the increase in the volume of geographic information,
60 obtained from cell phones and GPS devices. An example of application for the analysis of
61 geographic data is to help understand the behavior, movement of people and the models in the
62 cities.

63 Given the widespread use of mobile devices that track their geographical location, it has
64 become increasingly easy to acquire information related to users' trips in real time. This avail-
65 ability has triggered several studies based on user's position, such as the analysis of flows of
66 people in cities, and also new applications, such as route recommendation systems.

67 The huge amount of data to analyze has led researchers to develop computational tools and
68 data mining techniques combined with machine learning algorithms to allow better management
69 and understanding of geographic information. Collecting the urban trajectories, for example, it
70 might be helpful for people to make better decisions when they do visit an unknown location
71 and there are interesting places in the city or for knowing how users interact with each other.
72 It is possible to collect geographic information (e.g. from travel diaries, GPS data, image tags)
73 from different sources such as smartphones, GPS devices, social networks, etc.

74 For Microsoft Researchers, the GeoLife project [104], [105] was started to predict the mode
75 of transport between locations in a geographical region and support users to know how such
76 places can be reached.

77 Geographical locations are widely used in applications such as recommendation systems.
78 The widespread use of location-based Internet services (eg Google Maps) offers the opportunity
79 to collect user locations.

80 In the first part of this research we propose a multi-agent system for creating recommenda-
81 tions of Points of Interest (POIs). Using collaborative agents and a centralized server, POIs are
82 created from a dataset provided by Geolife. The server dynamically acquires information from
83 agents, which are stored on the user's mobile device, creating new suggestions on the next place
84 to visit. Another theme developed is the search for flows of multiple users. Discovered flows
85 of people are useful, e.g., to find common itineraries, to suggest reachable spots to users, or to
86 submit improvements on the public infrastructure.

87 In fact, in the second part of this work we propose a new algorithm to detect corridors in
88 the same geographical dataset. Corridors can be defined as geographical paths, with a minimum
89 length, that are commonly traversed by a minimum number of different users. We propose an
90 efficient strategy based on the Apriori algorithm to extract frequent trajectory patterns from the
91 geo-spatial dataset. After that, we refine the results obtained with a subsequent filtering step, by
92 using a Radius Neighbors Graph. This algorithm is relevant for transportation analytics because
93 it is the base to detect lacking lines in public transportation systems and also to recommend to
94 private users which route to take when moving from one part of the city to another on the basis
95 of behavior of the users who provided their logs.

96 Currently, organising a trip should take into account the number of people that will gather
97 at their chosen destination points, as it is necessary to avoid visiting a place that will become
98 overcrowded to comply with the restrictions due to the covid-19 influenza pandemic. Hence, an
99 estimate of the number of people that will be in some place in a future time can be valuable for
100 people moving and in situations where they could choose visiting some other place. Different
101 techniques have been proposed to mine the knowledge of different users from GPS logs such as
102 a probabilistic model to predict movements of people.

103 The proposed approaches have been validated by means of several experiments on data
104 concerning locations, for eliciting the meaningful places, starting from data available on taxis,
105 trucks, or people movements and also for predicting gatherings.

1.1 Contributions

106 Some results described in this thesis were published in journals:

- 107 • Cavallaro, C.; Verga, G.; Tramontana, E.; Muscato, O. *Eliciting Cities Points of Interest*
108 *from People Movements and Suggesting Effective Itineraries*. *Intelligenza Artificiale*, 14
109 (1), 49-61, 2020.
- 110 • Cavallaro, C.; Vitrià, J. *Corridor Detection from Large GPS Trajectories Datasets*. *Ap-*
111 *plied Sciences*. 2020; 10 (14):5003.

112 and conferences:

- 113 • Cavallaro, C.; Verga, G.; Tramontana, E.; Muscato, O. *Suggesting Just Enough (Un)Crowded*
114 *Routes and Destinations* Proceedings of the 21th Workshop From Objects to Agents,
115 Bologna, Italy. CEUR Workshop Proceedings, 2020. Vol. 2706, pp. 237-251. ISSN=1613-
116 0073.
- 117 • Cavallaro, C.; Verga, G.; Tramontana, E.; Muscato, O. *Multi-Agent Architecture for Point*
118 *of Interest Detection and Recommendation*. Proceedings of the 20th Workshop From Ob-
119 jects to Agents, Parma, Italy. CEUR Workshop Proceedings, 2019, Vol. 2404, pp. 98-104.
120 ISSN=1613-0073.

1.2 Aims and Approach

121 Nowadays, thousands of people are using their mobile devices to access new information in
122 relation to their geographic location. This innovation has given rise to new services, which for
123 example read GPS coordinates to receive information on nearby Points of Interest (POIs).

124 The first purpose of my research is to recover POIs and find the similarities between users
125 based on their StayPoints (SPs). Taking advantage of a multi-agent system, we propose an ap-
126 proach providing personalized recommendations of places of interests, such as libraries, museums
127 or restaurants, to users.

128 Current recommendation systems for tourist attractions or services available for a city are
129 generally managed manually, e.g. a city council publishes a list of public services or a publisher
130 indicates the top places, or POIs, for a place or region [4]. A static list of POIs may be insufficient
131 for the support that people expect, since for a large city or region there may be hundreds of
132 points, each with its own specific offer, where each point can vary over time and sometimes
133 with a high frequency. The manual care of such a list could be cumbersome, and in some cases
134 still a substitute for an updated list. In addition, each POI identified manually is not usually
135 associated with an indication of the most advantageous time period for visitors and with the
136 real-time conditions of the spot [26], [48].

137 Thanks to the availability of mobile devices, users can provide their comments on the places
138 they visit in real time. Such comments can be collected and be useful for other people looking
139 for advice. The proposed approach uses an agent-based solution to collect location data and
140 user satisfaction in order to offer to users the next place they could visit. This is based on
141 their current location, and on the time to reach the destination. The solution proposed in this
142 thesis may be able to compose a personalized list of POIs for each individual user based on
143 their previous position and trajectory, selecting real-time data, collected on a server, on user
144 experiences in some positions. While users provide their comments and rate a place, their precise
145 location and identity are preserved, and they are only kept on the user smartphone, whereas
146 peer agents and the server do not handle such data.

147 We identify significant places for groups of users from the behaviors extracted from the GPS
148 trajectories of different people and we want to predict the future behavior of a specific user
149 directed to these locations. We also want to understand the correlation among user behaviors.

150 The multi-agent system proposed, which exchanges information with the centralized server,
151 has been used for the creation of a recommendation tool for POIs. It provides users with: (i)
152 a list of POIs, and for each point, (ii) additional information based on real-time data collected
153 by other users, which help them choose the next destination with greater awareness. A key
154 objective is to have this information as close as possible to real-time data, and ranks places
155 based on feedback from other users, the most frequent time slots and the time spent visiting a
156 place by other users.

157 The approach offers a better experience by giving additional dynamic data (such as popular-
158 ity, number of users in real-time) to the list of POIs and by exploring their temporal relationships.
159 In fact, for POIs, which we determine using the DBSCAN algorithm, we take into consideration
160 the time intervals in which users have visited them, to offer a more advanced service. Finally,

161 the approach was designed to preserve user privacy, i.e. it does not reveal the exact location of
162 users.

163 The analysis of large human mobility datasets has the potential to provide many useful
164 suggestions to public operators as well as to individual users, but in order to provide this
165 information in a timely manner, efficient algorithms are needed.

166 In this scenario, corridor detection has emerged as one of the key elements to make informed
167 decisions about public transportation systems as well to recommend optimal routes to individual
168 users. In this case, the use of a brute-force algorithm is not an alternative due to its compu-
169 tational cost and more efficient solutions are necessary. To this end, we are proposing in this
170 thesis a new approach, based on the Apriori algorithm, which is suitable to the use in very large
171 datasets. In order to show its performance, we have analyzed the dataset of the GeoLife project
172 of Microsoft Research Asia, which includes geographical trajectories of 181 users that uploaded
173 GPS recordings corresponding to their routes in Beijing and its surroundings.

174 The algorithm Apriori, introduced by Agrawal et al. in [1], is a popular data mining algorithm
175 for market basket analysis to detect which items are frequently bought together. This algorithm
176 can be used in our scenario to reduce computational complexity. The key point of its application
177 is to consider different assignments for its *items* and *baskets*. Specifically, we consider two
178 alternating roles for *items* and *baskets*: in the first case, GPS points are represented as a set
179 of *items* and trajectories as a set of *baskets*, but later we reverse this assignment and consider
180 GPS points as a set of *baskets* and trajectories as a set of *items*. This strategy is particularly
181 suitable in our case, where we have to identify corridors for a large volume of data and the naive
182 application of Apriori would exceed any memory setting.

183 Both versions of the Apriori algorithm can be seen as filters that act on the set of points
184 and trajectories. In the first case, we are able to detect sets of points that are shared by at
185 least a number of trajectories. These points are not necessarily aligned to form trajectories.
186 In the second case, we detect sets of trajectories that contain at least a number of points. By
187 combining both approaches, the set of candidates that must be checked to determine the final set
188 of trajectories is reduced to a very low cardinality. The keypoint for this reasoning is to be
189 aware that the complexity of Apriori depends not only on the number of items and bags, but
190 on the number of different itemsets that are present in at least k bags. In the original case of
191 supermarket application, itemsets of 2 elements represent the bottleneck of the methods. In
192 our case, the bottleneck is created by itemsets with higher cardinality. By reversing the roles of
193 items and bags and doing a final check of the remaining candidates, we bypass this bottleneck.

194 The Apriori algorithm works in a discretized space and its results must be refined. To this
195 end, we use the Radius Neighbors Graph, which uses the mapping graph together with the
196 adjacency structure of the GPS points in the neighborhood of the points of a fixed trajectory.

197 To communicate to users how many people are in some places, the statistics accumulated
198 over time of recorded rides are usually used to estimate a measure of traffic or gatherings [20],
199 [94], [11]. Moreover, both popular online services and other apps just count the number of
200 people currently present in some place [28], [9], [58]. However, statistics gathered in the past
201 cannot be a reliable indication for the current situation that has to cope with e.g. restrictions
202 on gatherings, lower capacity of public transport means, etc. due to the influenza pandemic.

203 Additionally, a kind of real time measures of gatherings does not let other people plan their trip,
204 hence understanding whether, e.g., one hour later when arriving at the destination, the place
205 will still be (un)crowded.

206 A better estimate is therefore needed which takes into account: (i) the current amount of
207 people in some place, and (ii) the statistics on the number of people that being in some origin
208 place typically flow to another place to visit later on. In addition, an app behaving as an
209 assistant agent is needed to timely inform interested people.

210 Our latest work, presented in the final part of this thesis, proposes an approach to determine
211 the probability that users move along certain paths. Given the recordings of different user
212 locations, we compute the probability for a user that being in place A will move to another place
213 B (i.e. a possible destination), so when he arrives at place B he will contribute to the number
214 of people who gather there. By computing in advance the probability that it will go to place B
215 in a future time, we can guess if a place will become overcrowded. Our proposal for estimating
216 people's destinations is based on the analysis of the co-occurrence of places statistically visited
217 by a number of people above a threshold.

218 In addition, the app provides means to collect data on the current amount of people in a
219 place and thus on their trajectories. When we collect user data via the app, we ensure that user
220 privacy is preserved by providing only an approximate location to a central server.

221 Our approach can be useful in many contexts where estimating the number of people in
222 advance can be a crucial factor for a better service, such as when organizing public transport, or
223 for retailers, etc. In addition, it could be enriched with data, coming from competent authorities,
224 which reveal some places where a covid-19-positive was found. Then, using our paths found, we
225 could provide probabilities about other places where the infection may have spread.

226
227 The rest of the thesis is organised as follows. In Chapter 2 we survey related works, in Chapter
228 3 we clarify some definitions and techniques known in the literature. Chapter 4 describes the
229 datasets used and the approach proposed. Chapter 5 shows the test performed and the results
230 obtained. Finally in Chapter 6 we discuss the conclusions and draw future works.

Chapter 2

Related work

231 This chapter aims to present an overview of the state of the art on geospatial data analysis.
232 In Section 2.1 we present different and general techniques on the study of trajectories, while in
233 Section 2.2 we go into detail on works similar to those covered in this thesis, such as the extraction
234 and recommendation of places of interest. We show below the differences between our proposed
235 work and other papers in the literature, highlighting the innovative points presented in this
236 thesis.

2.1 Trajectory monitoring

237 There are several approaches in the literature regarding the analysis of geographic trajectories.
238 In fact, this is today an important research topic due to the increase in the volume of geographic
239 information.

240 Zygoras and Gunopulos propose in [108] the analysis of real GPS trajectories collected
241 from taxis operating in the city of Porto, moving buses in Dublin and Atlantic hurricanes from
242 1950 to 2004. In order to detect corridors, they discretize the trajectories using a grid and
243 decompose the space in different sets of frequently observed locations. The Latent Dirichlet
244 Allocation (LDA) model is applied to the trajectory dataset to extract frequent traffic patterns.
245 Their LDA formulation is similar to the one commonly used for natural language processing,
246 but replacing documents by grid cells and words by trajectories, respectively. LDA learns the
247 probabilistic distributions of hidden variables that are introduced in order to discover patterns
248 in the dataset. They apply a hierarchical clustering algorithm to the subtrajectories of each
249 frequent set, following a bottom-up approach, using Dynamic Time Warping (DTW) in order
250 to measure the distance between two subtrajectories. Finally, another algorithm selects the set
251 of corridors from the set of candidate corridors minimizing the MDL principle.

252 Bicocchi et al. [12] propose a system that suggests daily and local transport sharing opportu-
253 nities for short-term trips, by analyzing the traces of urban mobility in Milan and Turin. Based
254 on “Call Detail Record” (CDR) data, they propose algorithms to recognize similar paths that
255 can be used to recommend shared rides. The data used for the experiments was provided by a
256 telecom mobile operator and support a travel recommendation system for multiple users. Gen-

257 eral mobility routines are identified through an extension of the probabilistic generative LDA
258 model, performed on the set of available trips, including among others, leisure mobility rou-
259 tines and commuting trips. Both referred works seek paths that are common to multiple users,
260 organizing the movements made in time slots, but with different approaches. In both cases
261 geographical distance is measured through the Haversine formula. Their data are based on a
262 sequence of intermittent positions. In fact, they identify the location of the user with respect to
263 the cell when the user receives or makes calls or messages. Their goal is to extract 10 different
264 positions shared by multiple users; our approach instead is to identify people movements and to
265 detect the foreseeable routes by computing the number of gathering in some place.

266 Buchin et al. [16] consider the problem of detecting commuting patterns and propose algo-
267 rithms that cluster subpaths of given trajectories. The idea is to select a reference trajectory
268 and then find all subtrajectories that are close to this one by using the Fréchet distance. They
269 consider the optimisation problem of finding the longest cluster of a fixed size and the largest
270 cluster of a fixed length. The Fréchet distance is a distance measure for continuous shapes,
271 such as curves and surfaces, and it is defined using reparameterisations of the shapes. Since it
272 takes the continuity of the shapes under consideration, it is generally considered to be a more
273 appropriate distance measure for curves than the Hausdorff distance. For polygonal curves, the
274 discrete Fréchet distance (DDF) is a natural variant of the Fréchet distance. In this paper both
275 are used.

276 Rolim et al. propose in [78] a method to identify movement patterns of sets of trajectories
277 and analyze simulated trajectories in a region of Itaim Quarter in Santa Catarina - Brazil.
278 They consider the frequency distribution of points for identifying a set of frequent regions. The
279 framework consists of two phases: a partitioning phase, in which the trajectories are segmented
280 using the Minimum Description Length principle (MDL) and a clustering phase, that considers
281 the density to group similar segments in a cluster using the Fréchet distance as a measure of
282 similarity between curves.

283 Devogele et al. [33] describe a new algorithm to compute DDF which aims to lower com-
284 putation time and improved precision. It includes three different improvements: the first one is
285 the Douglas & Peucker filtering process. Indeed, for GPS trajectories, the number of positions
286 can be dramatically reduced. The second one is a pruning process. Only a small part of the
287 two matrices required to compute the discrete Fréchet distance are computed. The last one is
288 an improvement of the accuracy of the DDF. The proposed method is more complex than the
289 regular discrete Fréchet distance, but CPU time is reduced. In terms of accuracy, a balance
290 between CPU time and precision is required.

291 Zheng [103] conducts a survey about the major research in trajectory data mining. He
292 evaluates trajectory data preprocessing and management and some other data mining tasks
293 (trajectory pattern mining, outlier detection and trajectory classification). The paper also in-
294 troduces some methods that transform trajectories into matrices and graphic tensors, to which
295 different machine learning and data mining techniques can be applied.

296 Bian et al. [10] study the problem of trajectory grouping. Trajectory clustering has been
297 applied in pattern recognition, data analysis and machine learning, and it is prevalent in some
298 application fields such as the prediction of the movement of objects, traffic monitoring, under-

299 standing of activity and anomalous detection.

300 Crociani et al. in [29] offer an unsupervised learning approach for an automatic lane detection
301 in multidirectional pedestrian flows. Unlike our work, they focus on pedestrian dynamics in a
302 video using a distance function that takes into account the angular distance between the vectors.
303 After that they aggregate instantaneous information on both position and speed of pedestrians
304 to form clusters on short time windows. We use a fixed grid to group people positions into cells,
305 and then compute the probabilities using the data arriving from the agents in real time.

306 Zheng et al. [106] aim to mine interesting locations and classical travel sequences in a given
307 geospatial region. They first model multiple individuals' location histories with TBHG (tree-
308 based hierarchical graph). Next, supported by the TBHG, they propose a HITS (Hypertext
309 Induced Topic Search) inference model, which consider user's location as a directed link from
310 the user to the location.

311 In [80] Sakamane et al. describe multiple methods for inferring commuter route choice from
312 cell phone network data and suggest users the best routes to follow. Based on a calendar year
313 of CDR data collected by mobile users in Portugal, they interpolate waypoints of the route,
314 calculate the shortest distance between a choice of possible routes and the positions of mobile
315 use and also exploiting the Voronoi cells, assigning a choice route in coverage areas. Since the
316 real positions of the users are not as frequent as our fine-grained GPS recordings nor available
317 in the time frames in the absence of telephone activity, the tracks are deduced by choosing the
318 most plausible ones among the paths suggested by Google Maps. Their work proposes to draw
319 an inference for travel routes based on the statistics obtained, based on frequency in certain
320 time intervals and origins-destinations of fixed routes. The noise of the processed data is filtered
321 using DBSCAN as the clustering algorithm and the commuting radius as an admissible spatial
322 interval.

323 Zou et al. [107] present an algorithm for modelling various movement phenomena, such as
324 that of the movement of aircrafts in the airports of Hong Kong and Macau over a period of 30
325 days. For that purpose, a 4D time density is calculated, representing 3D spatial coordinates as
326 well as speed. Using this representation they can detect hotspots and trajectory convergences.

327 In [27] Chessa et al. deal with mobility through the PartecipAct campaign of the University
328 of Bologna, part of the projects called Mobile Crowd Sensing (MCS). This research aims to
329 draw some general inferences on the usability and level of realism of these datasets and evaluate
330 research solutions for MCS. The similarity between users, rather than being measured based
331 directly on interests and therefore places of visit in common, is evaluated through social graphs:
332 the more connections you have in common, the greater the probability for the same group of
333 people to frequent the same places on long periods (e.g. university colleagues) and therefore to
334 share the same habits (gyms, recreational activities in the same city). The participants who
335 voluntarily collect data in PartecipAct are all students (170 in total), the datasets we explored
336 are heterogeneous users. Furthermore, the client they use asks users if they want to perform
337 certain activities, and only in case of confirmation does it collect the data (active detection). The
338 back-end takes care of processing and archiving the detected data. Geonotification associates
339 activities with one or more geographical areas and automatically notifies users as soon as they
340 enter them. To verify the popularity of certain areas, the co-location of the users is extracted

341 from the device’s position via bluetooth: it is assumed that two devices are co-located and able
342 to communicate if they are positioned within 10 meters of each other for at least 150 seconds.
343 Frequent zones are then aggregated via DBSCAN.

344 Hosseinpoor et al. [47] introduce an approach for identifying critical points on a large volume
345 of real data. Three trajectory descriptors, namely the curvature areas, the turning points and
346 the self-intersecting points, are determined by detecting the critical point of trajectories using a
347 convex hull (TCP-CH) algorithm. These characteristics allow the detection of the directions of
348 change, represent the changes of the path and the presence of double-level intersections. They
349 are useful in order to show, in addition to anomalous values, the geometric properties of the
350 trajectories: shape, complexity, direction and distance. The approach allows, by selecting some
351 trajectory parameters, to identify the number and positions of the points of interest and identify
352 from them the similarities of the trajectories.

353 An accurate monitoring of user mobility provides support for efficient resource usage. E.g.,
354 it could help avoid traffic congestion [20],[59],[96].

2.2 Extraction and recommendation of places of interest

355 This work is an intersection of multiple disciplines, including POI recommendation systems,
356 collaborative filtering and privacy preserving systems.

- 357 1. Over the years, these analyzes have fueled different studies on trajectories, such as com-
358 puting the probability of moving from one POI to another, using, for example, the Markov
359 chains [40], [66] and therefore creating methods which predict the subsequent movements
360 of the individual user from the analysis of their POIs.
- 361 2. Collaborative Filtering (CF) approach is one of the approaches for creating recommenda-
362 tion systems. It creates suggestions using a similarity metric among users. The assumption
363 is that similar users probably have similar tastes. The concept of CF was introduced in
364 1992 by Xerox research staff within the Tapestry project, a system that allowed users to
365 track documents based on comments left by other users [44]. Later, several ratings-based
366 automated recommendation systems were developed, e.g. the GroupLens research system
367 [82] provides a pseudonymous CF solution for Usenet news and movies. Other technologies
368 have also been applied to recommending systems such as Bayesian networks [15], [68], [95]
369 and clustering [25], [95].
- 370 3. Privacy preserving: collaborative filtering techniques have been very successful in e-commerce
371 and in direct application recommendation. They are widely used and very useful but they
372 often fail to protect user privacy (for example in the case of GPS coordinates transmitted
373 with their timestamp), so they have some disadvantages. In [18], [19] privacy violations
374 are addressed with cryptographic systems, which can reduce the risk to the user. In other
375 research works, e.g. in [75], each user first disguises his private data, and then sends it
376 to the data collector. Therefore, a Randomised Perturbation (RP) technique is used to
377 camouflage private data [2]. Moreover, anonymisation techniques can be used, however
378 these introduce some attack problem, making datasets not very useful [77], [90].

379 Unlike other approaches, our proposal includes a solution to identify POIs through the use of
380 the DBSCAN (Density-Based Spatial Clustering of Applications with Noise) algorithm. Then,
381 collaboration filtering is used with the dynamic calculation of ratings based on user experiences.
382 For such a rating we use spatio-temporal variables offering a dynamic and realistic outcomes.

383 This is done by safeguarding the privacy of users because the centralised server only tracks
384 the movements near the POIs. Furthermore, it is important to offer a service that makes the
385 user and her privacy more secure. To do this we have users sharing their position only if they
386 are close enough to a POI and this information is manipulated to ensure user safety. E.g. the
387 position of a user within the radius of a POI will be saved in our central server with an error rate
388 of about 300 meters in order to preserve the user's privacy. This does not corrupt our system
389 data and better protects users.

390 In the literature, several studies use collaborative filtering to suggest itineraries or POIs. In
391 [48], the authors propose time-sensitive trip routes, consisting of a sequence of locations with
392 associated timestamps. In [101], Yoon et al. present a recommendation for itineraries based
393 on multiple user-generated GPS trajectories. In [49], Hsueh and Huang suggest a user-based
394 collaborative filtering with time preference to explore user preferences on places visited and offer
395 a recommended itinerary. In [28], data on Foursquare were used to find clusters considering both
396 spatial and social proximity, and results can be useful to characterise the amount of people and
397 the their flows on portions of a city. In [76], photos on Flickr were analysed to suggest routes that
398 can be pleasant, beautiful or quiet, according to the geo-location of photos and user comments.

399 Moreover, in [9], data gathered from both Flickr and Foursquare were used to identify POIs
400 in Milan (Italy).

401 Compared to the approach we propose, these systems have the following shortcomings:

- 402 1. the nearest points for the user are not identified;
- 403 2. the data is not updated in real time;
- 404 3. POIs are not automatically detected.

405 This work includes automatic calculation of POIs, starting from a series of trajectories. In
406 addition, while other works offer a complete itinerary, urging the user choose an origin and a
407 destination, in this proposal a user can immediately see the closest recommended points and
408 can dynamically change the next point of the itinerary according to current needs. That allows
409 users to search more effectively through travel information and organize the trip.

410 In [3], [54] and [65], the identification of the relevant places through various clustering al-
411 gorithms is based on the data provided by a telecommunications operator who has recorded
412 events such as calls or text messages. This is a statistic on the amount of people close to some
413 place of interest (e.g. a supermarket). For the above approaches, the data is only available to
414 very popular telecommunications operators or Internet service providers. Instead, the analysis
415 proposed in this document can be performed on data collected by users' devices, therefore with-
416 out the external support provided by large telecommunications operators. Furthermore, in our
417 approach we suggest locations that represent places that are meaningful to people, rather than
418 just popular places.

419 The algorithm for detecting personal anchor points like home or places of work or regularly
420 visited sites, unlike our approach is based only on frequency, or if these areas are regularly visited
421 in certain time slots, but not on the time spent in the surroundings. We extract the StayPoints
422 from each log route, and from these we determine the POIs, since if a place has meaning and
423 has passed through them, the user will presumably have spent some time nearby.

424 Furthermore passive mobile positioning is related to the precision of the cell, as each has
425 a certain coverage radius (which has a maximum value of about 35 *km*), therefore the exact
426 position cannot be determined but with an uncertainty given by the coverage of the cell to which
427 it hooks. Further uncertainties are given by the fact that the size of the cells are not fixed and
428 that a GSM mobile phone can hook to one of the two neighboring cells depending on whether
429 one is crowded or with a disturbed signal.

430 In such cases, since a visit is recorded through the cell phone only if the user makes a call
431 or receives a text message, it is not possible to determine how long he remains in that area. In
432 our case the entire path associated with time is marked, this analysis based on trajectories with
433 regularly recorded GPS positions allows us to investigate urban mobility in a more accurate way.

434 In [50], places of interest are collected according to an analysis of geotagged photos and
435 grouping data according to DBSCAN. This work differs from the proposed approach in that it
436 does not consider: time slots for places of interest, comments from people on such places and
437 data updated in real time.

438 An Area Of Interest (AOI) could be visited only during the night and therefore could not
439 be extracted because it is not popular according to this method. An analysis of these areas
440 dependent on the seasons is missing, as a place could be interesting to visit in summer and not
441 in winter or vice versa.

442 In this approach it is not easy to find flows of people, because the whole path is not recorded.

443 Assigning importance to an AOI, through the TF-IDF method (term frequency–inverse docu-
444 ment frequency) [81], based on the tags inserted is influenced by subjectivity: in social networks
445 some users would prefer to use popular tags to attract likes, conversely an AOI that has few
446 tags or little less significant tags would take away importance to the area.

447 In [58], the places of interest are determined using clustering techniques on the geolocation
448 of people according to the distance of Haversine, also used by us. The time parameter for
449 StayPoints is not changed on the basis of the dataset, because this problem is not analyzed by
450 them. Our work offers further contributions in terms of privacy for users, also providing users
451 with real-time comments and the classification of useful time slots for the selected places.

452 In recent years, CF recommendation systems have been supported by multi-agents systems.
453 A multi-agent system is a system with a significant number of independent agents interacting
454 with each other [87], [34]. In recent years, multi-agent systems have been widely used as they
455 are considered suitable for systems with modular architecture, thanks to their independence
456 [87]. In general, agents interact in three ways [85]: (i) each agent can communicate directly
457 with any other agent (called “autonomous agents”); (ii) agents communicate with each other
458 indirectly through an intermediary (called “facilitator”); (iii) all agents communicate with each
459 other through an intermediary, however agents may communicate with each other after the
460 communication has been set up by the intermediary (called “Mediator”). In the second case,

Table 2.1: Differences between this approach and other existing works

Feature	[101]	[48, 49]	[7]	[26]	[84]	[50, 58]	[3, 9, 54, 65]	this approach
visiting time		v						v
visiting order	v	v		v				v
visiting duration	v	v		v				v
location based service				v	v	v	v	v
uses data that is always up to date				v	v	v		v
use multi agent system			v	v	v			v
calculates points of interest						v	v	v
preserve privacy								v

461 the robustness can be low and the overhead is relatively high but the intermediary acts as a
462 protective wall for the privacy of the users because the agents do not communicate directly and
463 process the information received from the users, reducing their work [63] On the other hand, the
464 use of an intermediary has several advantages in terms of synchronization, reusability, scalability
465 and modularity [46].

466 In [56], the authors proposed a multi-agent system that allows users to optimize the energy
467 consumption of their smart homes. Each electrical device is configured as a virtual agent. These
468 agents work simultaneously and together to reduce consumption while ensuring user comfort,
469 energy costs and maximum energy savings.

470 Agent-based recommender systems have been proposed in the last years in different scenarios
471 including the tourism. For example, in [7], Batet et. al. present Turist@, i.e. a system
472 based on multi-agent technology to give personalised tour attraction recommendations more
473 effectively, highlighting the usefulness of finding points near the user. Similarly, [26] illustrates
474 an application to better plan travel decisions based on a multi-agent system. The authors of
475 [84] propose a system that produces recommendations for both individuals and groups.

476 Finally, this approach is the first approach in the tourism sector that deals with preserving
477 user privacy in two ways: by using a centralized server the agents do not directly know each
478 other to exchange information and by extracting only information on POIs, and an approximate
479 position of the user, only when he expresses an opinion on a POI.

480 Table 2.1 summarises the differences between our approach and other existing works in
481 location recommendation. Note that only the proposed approach makes use of gathered data
482 to compute points of interests, which therefore emerge and change while users interact with the
483 provided system; additionally, several precautions have been taken to ensure anonymity of the
484 user identity and their location.

Chapter 3

Preliminaries and Basic Definitions

485 This chapter recalls the basic concepts and definitions used in the thesis and describes some
486 algorithms applied in literature in some related works.

3.1 Geographic definitions

487 Parallels and meridians form, on the Earth's surface, a geographical grid which allows to iden-
488 tify the absolute position of each single point. Latitude and longitude define the geographical
489 coordinates of a place. A basic meridian, which passes for an astronomical observatory located
490 in Greenwich, it has been conventionally set. This meridian is also called zero or origin meridian
491 (has zero longitude).

492 **Longitude** is the angular distance of a point from the fundamental meridian. This corre-
493 sponds to the angle between the meridians plane of the point that we need to detect and the
494 plane of the fundamental meridian. **Latitude** is instead the angular distance of a point from
495 the equator, being measured along the meridian which passes through that point.

496 According to the rules of spherical trigonometry we can determine the shortest curve that
497 connects two points of the Earth (geodesic). Remember that thinking about the Earth as a
498 sphere is an approximation, since in reality it has the shape of an ellipsoid and the terrestrial
499 ray is longer at the equator (6,378 km) and shorter at the north and south poles (6,357 km).

500 **Definition 3.1.1.** *For the Spherical Law of Cosines the distance d between two points, $A = (lat_1,$
501 $lon_1)$ and $B = (lat_2, lon_2)$, of given latitude and longitude is defined by:*

$$d(A, B) = R * \arccos(\sin(lat_1) * \sin(lat_2) + \cos(lat_1) * \cos(lat_2) * \cos(lon_2 - lon_1)). \quad (3.1)$$

502 *The angles used are expressed in radians. This formula approximates the Geoid to a sphere*
503 *of medium quadratic radius, $R = 6,372.795477598$ km, so the distance calculation may have*
504 *an error of 0.3 %, particularly in the polar extremities, and for long distances crossing several*
505 *parallels.*

506 **Definition 3.1.2.** *The Haversine formula is more accurate than the cosine formula for mea-*
507 *suring distanced, due to problems associated with small distances. Important in navigation, the*

508 first table of Haversines was published by Andrew in 1805, but the term Haversine was coined in
 509 1835 by Inman [53]. The Haversine distance [86] $d(A, B)$ between two points A and B , which
 510 gives the great-circle distance in kilometers between two points on a sphere, is defined as follows:

$$d(A, B) =$$

$$2R \arcsin \sqrt{\sin^2 \frac{lat_1 - lat_2}{2} + \cos lat_1 \cos lat_2 \sin^2 \frac{long_1 - long_2}{2}}. \quad (3.2)$$

511
 512 The angles used are expressed in radians, R is the terrestrial radius (average radius 6,371.005076123
 513 km).

514 Since the Earth is divided into equal parts by the meridians, we know that 1 degree of
 515 latitude equals 111.121 km while 1 degree of longitude varies according to where we are because
 516 all parallels have different diameters moving from the equator to the north.

517 A trajectory log is a set of geographic locations for a user in motion, ordered by time, from
 518 the moment he starts recording his journey until the moment he stops.

519 **Definition 3.1.3.** A trajectory T is defined as $T = ((lat_1, long_1, date_time_1); \dots; (lat_n, long_n,$
 520 $date_time_n))$, where lat_i and $long_i$ for $i = 1 \dots n$ are the latitude and longitude in decimal degrees
 521 of the GPS points of the trajectory, $date_time_i$ is the time and the date in which every point was
 522 registered and n is the length of the trajectory T (the total number of its points).

523 Each trajectory in the GeoLife database, that we used for experiments and described in the
 524 next chapter, is identified by a trajectory ID and by the user ID who has traveled it.

525 **Definition 3.1.4.** A sub-trajectory is a sub-sequence of a T trajectory, that is, a part of points
 526 consecutive within the trajectory.

527 $SubT = ((lat_{m_1}, long_{m_1}, date_time_{m_1}); (lat_{m_2}, long_{m_2}, date_time_{m_2}); \dots; (lat_{m_k}, long_{m_k}, date_time_{m_k}))$,
 528 where $1 \leq m_1 \leq m_2 \leq m_k \leq n$.

529 Since the trajectory data can be collected by different devices, the log routes will be of
 530 different sizes and related to different objects in movement and which can be of different sources
 531 and categories:

532 1. People

533 The trajectories in this case refer to people who have registered their movements in the
 534 real world in the form of GPS tracks, and this it can take place both in active and passive
 535 mode for an arbitrary period of time.

536 • Active recording

537 The data of this type of source concern people who spontaneously enable the location
 538 service on the devices in order to obtain GPS tracks for trajectories, for this reason
 539 it is called active recording. This type of data generally affects people in a trip,
 540 who record their track for storage purposes travel itself, or cyclists or runners who

541 record their activities for subsequent sport analysis and comparisons with previous
542 performances. In some social networks it is also possible to recreate a route made
543 from geo-tagged media in different places.

544 • **Passive recording**

545 It is possible to reconstruct a trace relating to a person's movement without the
546 individual activating specifically the GPS detector on their device, in this case we
547 are talking about passive recording. For example by the sequence of telephone cells
548 connected by the mobile devices, with the corresponding time in which you have
549 passed through the associated area, or when using a credit card it is possible to
550 connect the different positions in which you were registered.

551 2. **Vehicles**

552 A multitude of transport vehicles are nowadays equipped with GPS systems that track their
553 journeys, including for example taxis, buses, ships and planes. The related data report the
554 temporal positions recorded by the devices, which can be used for various purposes such
555 as better allocation of resources, traffic analysis or improvements of transport.

556 3. **Movements of animals**

557 Recently the behavior of animals has been widely studied by biologists and zoologists
558 to understand behavior and displacements carried out by groups of animals in different
559 areas, for this purpose different GPS devices have been used for reconstruct trajectories
560 and subsequently analyze their movements.

561 4. **Mobility of natural phenomena**

562 The same principle is applied to natural phenomena, such as hurricanes, tornadoes and
563 ocean currents. The trajectories related to these phenomena capture their movement
564 and allow analysis to be performed deeper on the usual changes and behaviors of these
565 phenomena.

3.2 Other different distances used to compare trajectories

566 Let S be the symbolic representation of a measurement space and let x, y, z be three points
567 in S . A **dissimilarity** or **semimetric measure** is defined as a function $d(x, y) : S \times S \rightarrow R$
568 which satisfies the following conditions:

- 569 1. $d(x, y) = 0 \Leftrightarrow x = y$;
- 570 2. $d(x, y) \geq 0 \forall x, y \in S$;
- 571 3. $d(x, y) = d(y, x) \forall x, y \in S$.

572 The first condition indicates the reflexivity of the relationship, the second requires that the
573 distance, however is not negative, the third finally indicates symmetry.

574 Under the conditions listed above, if the function also satisfies the following "Triangle in-
575 equality":

576 $d(x, y) \geq d(x, z) + d(y, z) \forall x, y, z \in S$
577 is said to be a **metric**.

578 Likewise, a **similarity function** is defined to satisfy the conditions:
579

580 1. Symmetry: $s(x, y) = s(y, x) \forall x, y \in S$;

581 2. Positivity: $0 \leq s(x, y) \leq 1 \forall x, y \in S$.

582 If it also satisfies these conditions

583 3. $s(x, y)s(y, z) \leq [s(x, y) + s(y, z)]s(x, z) \forall x, y, z \in S$ and

584 4. $s(x, y) = 1 \Leftrightarrow x = y$,

585 it is called a **similarity metric**.

586 For a data set with input patterns, we can define an symmetric matrix, called proximity
587 matrix, whose (i, j) th element represents the similarity or dissimilarity measure for the i th and
588 j th patterns $(i, j = 1, \dots, N)$.

589 A fundamental part of each clustering algorithm is an appropriate choice of a distance or
590 dissimilarity measure that allows to translate numerically the concepts of association between
591 similar elements and distinction between elements belonging to different clusters.

592 Given two sets of data X and Y of length m :

$$\begin{aligned} X &= x_1, x_2, \dots, x_i, \dots, x_m \\ Y &= y_1, y_2, \dots, y_j, \dots, y_m \end{aligned}$$

595 it must be established how close the two series are.

596 **Definition 3.2.1.** *The **Euclidean distance** between X and Y is defined as follows:*

$$D(X, Y) = \sum_{i=1}^m \sqrt{(x_i - y_i)^2}.$$

599 **Definition 3.2.2.** ***Minkowski's distance** between X and Y is defined as follows:*

$$D(X, Y) = (\sum_{i=1}^m (x_i - y_i)^r)^{1/r}.$$

600 where r is an input parameter that can take:

- 601 • $r = 1$: *Manhattan distance*;
- 602 • $r = 2$: *Euclidean distance*;
- 603 • $r \rightarrow \infty$ *is the maximum distance between the components of a vector.*

607 **Definition 3.2.3.** *The Lagrange-Tchebychev distance is defined as follows:*

$$608 \quad D(X, Y) = \max_{1 \leq i \leq m} |X_i, Y_i|. \\ 609 \\ 610$$

611 **Definition 3.2.4.** *The Mahalanobis distance is defined as follows:*

$$612 \quad \text{Mahal}(X, Y) = (X - Y) \sigma(X, Y)^{-1} (X - Y)^T \text{ where } \sigma(X, Y) \text{ is the covariance matrix.} \\ 613$$

614 Many applications, such as those that use financial and marketing trends, require certain
615 temporal similarity criteria on the data. Look for patterns similar to some known behaviors can
616 help predict or test hypotheses.

617 We use the following conventions. If S and Q are two sequences:

- 618 • $\text{len}(S)$ is the length of S ;
- 619 • $S[i : j]$ is the sub-sequence in S identified from position i to j ;
- 620 • $d(S, Q)$ is the distance between the two sequences.

621 We can classify, given a sequence, two types of similarity queries.

- 622 1. **Whole matching:** given N data sequences S_1, S_2, \dots, S_N and a query sequence Q , all of
623 the same length, we want to find all the sequences that they are distant at most ϵ from Q :
624 $d(S_i, Q) \leq \epsilon$.
- 625 2. **Subsequence matching:** given N data sequences S_1, S_2, \dots, S_N in length arbitrary, a Q
626 query sequence and a tolerance ϵ , we want to identify the S_i sequences ($1 \leq i \leq N$) which
627 contain subsequences whose distance from Q is smaller than or equal to ϵ .

628 If sequences X and Y have respectively n and m element, we can align these two sequences
629 using the **Dynamic time warping (DTW) algorithm** [8].

630 In fact, DTW can lead to a measure of distance between the two aligned sequences. This
631 algorithm is particularly useful for treating the sequences in which individual components have
632 characteristics that require over time, and for which the simple linear expansion or compression
633 of the two sequences does not bring satisfactory results.

634 For example, similarities in walking could be detected using the DTW, even if one person
635 walked faster than the other, or if there were accelerations and decelerations during an obser-
636 vation. It has been used in various fields of application, from voice recognition to recognition of
637 motor activities.

638 Example

639 Consider an $n \times m$ matrix called D , where $d(i, j)$ is the distance between x_i and y_j .

640 There are many metrics used for computing distances. Typically Euclidean distance is used
641 and therefore, at this point it is possible to build an alignment path between the two sets of

642 data. The optimal alignment between X and Y data series is obtained by minimizing the local
 643 distance $d(i, j)$ between the series of points of X and of Y .

644 A matrix D $n \times m$ is therefore constructed with the following procedure:

645 the following values are assigned to positions $(0,0)$, $(1,1)$, $(i,0)$, $(0,j)$:

646 $D(0,0)=0$,

647 $D(1,1)=d(1,1)$,

648 $D(i,0)=\infty$ for $1 \leq i \leq n$,

649 $D(0,j)=\infty$ for $1 \leq j \leq m$.

650 The remaining values of the matrix are assigned using a recursive procedure:

651 $D(i,j)=d(i,j)+\min\{D(i-1,j-1), D(i-1,j), D(i,j-1)\}$.

652 A further step is required to obtain the optimal alignment between X and Y , that is the
 653 optimal warping path W . W is defined as a continuous path of w_k elements in the matrix D ,
 654 where each w_k corresponds to an element $D(i, j)$ of the matrix D , chosen so as to optimize the
 655 alignment between X and Y .

656 The complete warping path is described as: $W = w_1, w_2, \dots, w_k, \dots, w_K$

$$657 \max(m, n) \leq K < m + n - 1.$$

658 The search for the path is bound by the following conditions:

- 659 1. **Boundary conditions.** $w_1 = (1, 1)$ and $w_k = (n, m)$. This means that warping path
 660 must start and end in the elements $D_{1,1}$ and $D_{n,m}$ of the matrix.
- 661 2. **Continuity.** Given $w_k = (a, b)$, then $w_{k+1} = (a', b')$, where $a' - a \leq 1$ and $b' - b \leq 1$.
 662 This forces warping to move on adjacent cells horizontally diagonally or vertical.
- 663 3. **Monotonicity.** Given $w_k = (a, b)$, then $w_{k+1} = (a', b')$, where $a' - a \geq 0$ and $b' - b \geq 0$.
 664 This ensures that the path continues at each step towards $D_{n,m}$ and does not go back to
 665 $D_{1,1}$.

666 The optimal path is obtained by minimizing the: $DTW(X, Y) = \min\left(\frac{\sqrt{\sum_{k=1}^K w_k}}{K}\right)$.

667 The constant normalization for the K value is used to compensate for the fact of have
 668 warping path with different lengths in order to optimize the alignment local between datasets.
 669 To reconstruct the path, you need to memorize another structure in addition to D , used to store
 670 the direction of the path. The path is obtained a posteriori using this pointer.

3.2.1 Fréchet Distance e DDF

671 For an intuitive definition we can imagine a man crossing a finished curved path while walking
 672 his dog on a leash, with the dog crossing a separate one. We suppose that the dog varies his
 673 speed to keep as loose in his leash as possible and the dog and this owner have to walk without
 674 backtracking from an endpoint each other: the Fréchet distance [39] between the two curves is
 675 the length of the short leash sufficient both to cross their separate paths.

676 The Fréchet distance is a measure of similarity between curves that takes into account the
 677 location and ordering of the points along the curves. Therefore it is often better than the
 678 well-known Hausdorff distance.

679 A parameterized curve in \mathbb{R}^d can be represented as a continuous function $f : [0, 1] \rightarrow \mathbb{R}^d$. A
 680 monotone reparametrization α is a continuous non decreasing function : $\alpha : [0, 1] \rightarrow [0, 1]$ with
 681 $\alpha(0) = 0$ and $\alpha(1) = 1$. Given two curves $f, g : [0, 1] \rightarrow \mathbb{R}^d$, their Fréchet distance , $\delta F(f, g)$, is
 682 defined as:

$$\delta F(f, g) = \inf_{\alpha, \beta} \max_{t \in [0, 1]} d(f(\alpha(t)), g(\beta(t)))$$

683 where $d(x, y)$ denotes the Euclidean distance between points x and y , and α and β range
 684 over all monotone reparametrizations.

685 Alt and Godau described the first known algorithm [5] for the continuous Fréchet distance
 686 between two polygonal curves in Euclidean space, based on the principle of parametric search.
 687 The execution time of their algorithm for two polygonal curves with m and n segments is
 688 $O(m * n * \log(m * n))$.

689 The **discrete Fréchet distance** (DFD), also called coupling distance, is an approximation
 690 of the Fréchet metric for polygonal curves, defined by Eiter and Mannila [35]. The discrete
 691 Fréchet distance considers only leash positions where the endpoints are located at the vertices
 692 of the two polygonal curves and never within an edge. This particular structure allows the
 693 discrete Fréchet distance to be calculated in polynomial time by an easy dynamic programming
 694 algorithm. DFD is computable in $O(m * n)$ time.

695 **Definition 3.2.5.** *The discrete Fréchet distance between two polygonal curves $a : [0, m] \rightarrow \mathbb{R}^k$*
 696 *and $b : [0, n] \rightarrow \mathbb{R}^k$ is defined as:*

$$DFD(a, b) = \min_{\substack{\sigma: [1:m+n] \rightarrow [0:m], \\ \beta: [1:m+n] \rightarrow [0:n]}} \max_{s \in [1:m+n]} \{d(a(\sigma(s)), b(\beta(s)))\}$$

698 where σ and β range over all discrete non-decreasing onto mappings of the form $\sigma : [1 : m + n] \rightarrow [0 : m]$, $\beta : [1 : m + n] \rightarrow [0 : n]$.

700 This similarity measure has been applied to different problems, from morphing and hand-
 701 writing recognition, to the alignment of the protein structure, to computer graphics, as well
 702 as geographic applications, because it captures perceptual and geographical similarity between
 703 discrete trajectories.

704 DFD can tolerate non-uniform sampling rate. Other distance measures require that points
 705 along the trajectories are uniformly and densely sampled, but for real datasets it is difficult
 706 to have. Another characteristic of Discrete Fréchet distance is the local time shifting, that is
 707 the ability of tolerating short term discrepancies (e.g., missing samples, measurement errors) in

708 aligning two trajectories.

709
710 Such distances can be used in our case to compare which GPS subtrajectories are closest to
711 each other. We initially calculated the distance between pairs of points for different trajectories
712 using Haversine’s formula to find flows (see Sections 4.2.5 and 5.2) and then to determine the
713 frequently visited areas more quickly we used the Apriori algorithm (definitions in Section 3.6)
714 in Sections 4.2.6 and 5.6.

715 3.3 Recommendation Systems and collaborative filtering

716 It is difficult to extrapolate relevant information from a large amount of data such as the corre-
717 lation between individuals. The similarity between users is very important for example in the
718 customers and in the commercial companies. Many searches have been based on the transactions
719 of supermarkets and bookstores or on the behavior of the users in the web communities, to allow
720 an individual to discover potential friends based on interests similar in the areas of books, films
721 and music. In addition, this would allow sellers to improve their sales and marketing strategies
722 by recommending the right products to consumers.

723 Similarity is also important in geographic information systems, also exploring the correlation
724 between geographic location. We think that if there are users who share the same location history
725 they could share similar interests and preferences. The more places they have in common, the
726 more they could be connected. And so we want to mine the similarity based on GPS trajectories
727 of different users generated in the real world.

728 Recommendation systems use community opinion to identify more effectively their interest
729 among a range of content. A well known technique used in such systems is called Collaborative
730 Filtering [83] when trying to predict the evaluation of a product for a given user. The general
731 idea behind it is that similar users vote similarly on similar items. Therefore, if the similarity
732 between users and items is established, you can make a potential prediction for a user’s rating
733 for some items.

734 It has also been explored in social networks to facilitate people to identify potential friends
735 and content that is interesting to them on the Web. One of the most commonly used algorithms
736 is the closest neighborhood approach. In our work we extend the direction of similarity from
737 people’s online behaviors to the chronology of places visited in the real world.

738 Recommendation systems are divided into three main categories [13]: collaborative filtering,
739 (CF), content-based filtering and (CB) and hybrid filtering (HF). Content-based filtering [64]
740 makes recommendations based on user choices made in the past (e.g., for a person who likes
741 a carbonated drink such as cola, the system offers recommendations for similar soda drinks).
742 Collaborative filtering [83] allows users to give ratings about a set of elements, so that when
743 enough information is stored on the system, it is possible to make recommendations to each
744 user based on information provided by those users that have the most in common with them
745 (for example, if Bob and Alan have seen the same horror films, one of the films of the same
746 category seen by Bob is suggested to Alan and vice versa). The hybrid technique is a mixture
of the first two. In our approach, suggestions are chosen based on the experiences and choices

747 of other users, falling back into the CF systems. These experiences are documented by agents
748 who observe user behaviour in real time.

3.4 Multi-Agent System

749 Several studies have proposed the use of multi-agent systems (MAS) [30] in a wide range of
750 different domains. In 1998, a study described a supporting system for suggesting possible pur-
751 chases during shopping based on the GPS position with the use of agents [37]. In general, there
752 are two main approaches to MAS developments: centralised policies (CMAS) and decentralised
753 policies (DMAS) [99].

- 754 • A centralised approach consists of taking all of the decisions in one place. In a typical
755 CMAS, a central server collects all the relevant data that come from the different actors
756 (that is, agents) and identifies the decisions for each agent according to the global system
757 state. The centralised view of the system can be described by a multi-agent Markov
758 decision process model, an example is presented in [14].
- 759 • A decentralised approach consists of making each entity responsible for its own decision.
760 In a typical DMAS, an agent cannot see other agents local states and local actions, and
761 has to decide the next local action on its own. Thus, each agent has only a partial view of
762 the systems global state, and different agents have different partial views. A good example
763 is in [100] whose authors propose a decentralised multi-agent decision process framework
764 that provides the basis for a decision-theoretic study of decentralised policies.

765 The decentralized architecture has advantages in synchronisation, reusability, scalability, and
766 modularity [43], [46], [92]. However, the complexity of decentralised systems is greater than
767 that of centralised ones. Although decentralisation shows obvious advantages, decentralisation
768 also has its own drawbacks, including that agents cannot predict the group behavior based
769 only on the available local information, possible instability, and sub-optimal decisions. Due to
770 the importance of total knowledge, the choice for the approach we proposed fell into the first
771 category. In this work, apps on smartphones are agents that indirectly communicate among
772 themselves and take some decisions on behalf of the user. In addition, the centralized server is
773 able to filter information by offering advice to users without sending their sensitive data; this
774 preserves the user’s privacy.

3.5 Clustering techniques

775 Cluster analysis is the search for groups of objects such that objects appear members of a group
776 are “similar” to each other and different from the objects in the others groups. The intercluster
777 distances (between different cluster data) are maximized, intracluster distances (between data
778 belonging to the same cluster) are minimized. Clustering is referred to as unsupervised classi-
779 fication (unsupervised learning): as for classification, the purpose is to segment the data, but

780 without assigning class labels. In fact, we do not have pre-defined classes, but each cluster can
781 be interpreted as a class of similar objects.

782 Clustering is a collection of clusters. An important distinction is between:

783 (i) Partitioning clustering: a division of objects into subsets (clusters) not overlapping, where
784 each object belongs exactly to a cluster.

785 (ii) Hierarchical clustering: a set of nested clusters organized as a tree hierarchical.

786 Among the types of partitional clusters we find the center-based one, where a cluster is a
787 whole of points such that a point in the cluster is closer (or more like) to the “center” of the
788 clusters rather than in the center of each other. The center of a cluster is called centroid, the
789 average of all points belonging to the cluster, or mediod, the most “representative” point of the
790 cluster.

3.5.1 K-Means

791 K -means [38] is a partitioning clustering technique distance-based, in which each cluster is
792 associated with a centroid. Each point of our data is assigned with the cluster whose centroid
793 is closest. The number of clusters, K , must be specified as input, as well as the matrix of the
794 data to be clustered (in our case the pairs of latitude, longitude of the GPS points).

795 The choice of the number of K clusters must be well thought out. Generally it is set between
796 a minimum limit of 3-5 (to avoid too large clusters) and a maximum (clusters too small would
797 undermine the simplification process that underlies the grouping).

798 Given K , the K -means algorithm is implemented in iterative steps, and detects center-based
799 cluster:

- 800 1. Choose K points that represent the initial centroids (means) of the clusters. You can
801 randomly choose K observations, for initializing the cluster center.
- 802 2. Calculate distances from all observations to each centroid.
- 803 3. Assign each object to the center of the nearest cluster.
- 804 4. Calculate the average of the observations in each cluster to get new K 's centroid positions
805 (distance recalculation).
- 806 5. Repeat steps 2 through 4 until the cluster assignments change or the maximum number
807 of iterations is reached.

808 It is possible to use different metrics for the K -means such as Euclidean distance, Manhattan
809 distance, cosine, correlation, hamming.

810 There are various evaluation measures to estimate the quality of the clustering obtained for
811 the chosen parameters, such as Calinski-Harabasz Index [17], Davies-Bouldin Index [31], Sum of
812 Squared Error and Silhouette [79]. The Calinski-Harabasz index is computed as a ratio of inter-
813 cluster variance to within-cluster variance, whereby good clustering solutions tend to have large
814 inter-cluster variation and smaller within-cluster variation, so that the optimal cluster number
815 corresponds to the largest ratio value obtained.

3.5.2 DBSCAN

Other partitioning clustering methods are based on the notion of density. Their idea general is to grow a given cluster until the density (number of objects or points of data) in a neighborhood of a certain radius does not exceed a certain threshold. This method can be used to filter noise and discover arbitrarily shaped clusters.

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) [36] and its OPTICS [6] extension build clusters based on a density-based connectivity analysis.

Some new definitions are required to understand how density-based clustering algorithms work. DBSCAN requires only two parameters (no prior knowledge of the number of clusters): the maximum radius of the neighbourhood, ϵ , and the minimum number of points in an ϵ neighbourhood of a point, $MinPts$.

The neighborhood of radius ϵ of a given object p of the dataset D is called “eps-neighborhood” object: $N_\epsilon(p) = \{q \in D \text{ s.t. } dist(p, q) \leq \epsilon \}$.

The shape of a neighbourhood depends on the distance function.

Given ϵ and $MinPts$, this algorithm categorizes the objects into exclusive groups.

- A point is a **core point** if it has more than a specified number of points ($MinPts$) within ϵ . These are points that are the interior of a cluster.
- A **border point** has fewer than $MinPts$ within ϵ , but is the neighborhood of a core point.
- A **noise point** is any point that is not a core point nor a border point.

The algorithm follows these steps:

1. randomly select a point p ;
2. retrieve all points density-reachable from p wrt. ϵ and $MinPts$;
3. if p is a core point, a cluster is formed;
4. if p is a border point, no points are density-reachable from p , then visit the next data point;
5. continue the process until all points have been processed.

The main principle is therefore that a **cluster** is defined as a maximum set of points related to density. It is separated by low density regions (which represent noise).

DBSCAN discovers groups of arbitrary forms (spherical, elongated, linear). Its pseudocode is

in Algorithm 1.

Algorithm 1: DBSCAN

Procedure: *DBscan(dataset D, ϵ , MinPts)*

$C = [-]$

foreach *unvisited point P in dataset D do*

 mark P as visited

$N_\epsilon(p) = \text{getNeighbors}(P, \epsilon)$

if *size($N_\epsilon(p)$) < MinPts: then*

 | mark P as NOISE

else

 | $C = \text{next cluster}$

 | $\text{Expand_cluster}(P, N_\epsilon(p), C, \epsilon, \text{MinPts})$

end

end

Procedure: *Expand_cluster($P, N_\epsilon(p), C, \epsilon, \text{MinPts}$)*

add P to cluster C **foreach** *point P' in $N_\epsilon(p)$ do*

if *P' is not visited then*

 | mark P' as visited

 | $N'_\epsilon(p) = \text{getNeighbors}(P, \epsilon)$

 | **if** *size_of($N'_\epsilon(p)$) \geq MinPts then*

 | $N_\epsilon(p) = N_\epsilon(p)$ joined with $N'_\epsilon(p)$

 | **end**

end

if *P' is not yet member of any cluster then*

 | add P' to cluster C

end

end

Procedure: *regionQuery(P, ϵ)*

Return all points within P 's ϵ -neighborhood (including P)

The two parameters can be determined by a heuristic and must be chosen according to the application environment.

Among the advantages we find that it handles noise well and is particularly suitable for spatial data sets: geomarketing, tomography, satellite images.

In the context of GPS trajectory analysis we have used clustering techniques in order to group points of different trajectories spatially close together quickly.

3.6 Apriori algorithm

The algorithm Apriori, introduced by Agrawal et al. in [1], is a popular Data Mining algorithm for Market Basket Analysis to detect which items are frequently bought together.

854 **Market Basket Analysis** is indeed a technique which aims to use information on what,
855 when and how much you buy for building models of purchasing behavior. These are then used
856 to establish the arrangement of the shelves, the composition of promotions, the issue of discount
857 coupons.

858 The data mining technique that is correlated with the market basket analysis is the automatic
859 generation of **association rules**. This is based on the search for recurring configurations within
860 the data relating to *sales transactions* or *basket*. By **sales transaction** is meant the set of
861 products (the contents of a shopping cart for a customer), accompanied by the information
862 relating to the time and place of purchase. An **association rule** says that there is a strong
863 correlation between the purchase of two or more products.

864 Association rules can also be applied outside of this context, when it is possible to identify
865 a transaction or “basket”: it represents the choice of different elements by a “user”; in our case
866 a user saves his trajectories, which contain his positions in the geographic network, recorded as
867 time changes.

868 Regarding the Apriori algorithm, we have to define the concept of *itemset* and its *support*.
869 A set of unique items is usually referred as **itemset**, and an itemset containing K items is called
870 **K -itemset** or **K -tuple**. An itemset satisfies **minimum support** (min_sup), if the occurrence
871 frequency of the itemset is greater than or equal to min_sup . If an itemset has minimum support,
872 then it’s called **frequent itemset**.

873 Apriori consists of a “Join” function, which generates a new combination of elements (the
874 candidates) and a “Prune” function in which combinations of elements that do not satisfy the
875 minimum support are eliminated. Its pseudocode is presented in Algorithm 2.

876 In this iterative search algorithm the set of elements of cardinality K is used to analyze the
877 set of elements of cardinality at $K + 1$. The algorithm is based on a simple principle: if a set of
878 items is frequent, all subsets of this set are also frequent.

879 The advantage is that at each Apriori level, the algorithm checks whether a tuple is frequent
880 only if it is made up of other frequent tuples with smaller support. This search space reduction
881 strategy depends on whether the support of an itemset never exceeds the support of its subsets.
882 This property is known as *antimonotonicity* of the support measure.

883 In fact, a “small” itemset occurs in all transactions where there is a larger “itemset”, plus
884 possibly other transactions. For example if $\{A, B, C\}$ is a frequent itemset, a transaction
885 containing $\{A, B, C\}$ must clearly also contain its subsets $\{A, B\}$, $\{A, C\}$, $\{B, C\}$, $\{A\}$, $\{B\}$
886 and $\{C\}$. Each subset occurs in the transaction database at least as much as the item $\{A, B,$
887 $C\}$ and consequently it is itself a frequent itemset. Conversely, if an itemset like $\{A, B\}$ is not
888 frequent, then none of its own subsets is frequent.

889 The “Join function”, procedure for generating candidates (*Apriori_gen* in Algorithm 2), takes
890 as input the set of $(K - 1)$ -itemset frequent and returns a subset of all frequent K -itemsets. The
891 purpose of the procedure is to return the fewest candidates, so you have to do the calculation
892 of support on as few candidates as possible.

893 By discretizing the data through a grid, adapting the roles of itemsets and baskets of this
894 algorithm to our context, we find with Apriori algorithm the longest and closest sub-trajectories

contained in the cells of the grid and shared by a minimum number of users.

Algorithm 2: Apriori

Input: D' : transaction database; min_sup : the minimum support threshold for every itemset.

Output: frequent itemsets of D'

Procedure: $\text{Apriori}(D')$

$L_1 = \text{find_frequent_1-itemsets}(D')$

for ($k = 2; L_{k-1} \neq \emptyset; k++$) **do**

$C_k = \text{Apriori_gen}(L_{k-1})$

foreach *transaction* $t \in D'$ **do**

$C_t = \text{subset}(C_k, t)$

 // get the subsets of t that are candidates

foreach *candidate* $c \in C_t$ **do**

 | $c.\text{count}++$

end

end

$L_k = \{c \in C_k \mid c.\text{count} \geq \text{min_sup}\}$

end

return $E = \bigcup_k L_k$

Procedure: $\text{Apriori_gen}(L_{k-1}$: frequent($k-1$)-itemsets)

foreach *itemset* $\mathcal{L} \in L_{k-1}$ **do**

foreach *itemset* $\mathcal{N} \in L_{k-1}$ **do**

if ($\mathcal{L}[1] = \mathcal{N}[1] \cap (\mathcal{L}[2] = \mathcal{N}[2]) \cdots \cap (\mathcal{L}[k-2] = \mathcal{N}[k-2])$)

$\cap (\mathcal{L}[k-1] < \mathcal{N}[k-1])$ **then**

$c = \mathcal{L} \times \mathcal{N}$

 // join step: generate candidates

if $\text{has_infrequent_subset}(c, l_{k-1})$ **then**

 | delete c

 // prune step: remove unfruitful candidate

else

 | add c to C_k

end

end

end

end

return C_k

Procedure: $\text{has_infrequent_subset}(c$: candidate k -itemset;

L_{k-1} : frequent($k-1$)-itemsets)

foreach ($k-1$)-subset s of c **do**

if $s \notin L_{k-1}$ **then**

 | **return** TRUE

else

 | **return** FALSE

end

end

3.7 Sentiment Analysis

Sentiment Analysis (SA) also called Opinion Mining is a field of Natural Language Processing (also studied in Data Mining) that analyzes people's opinion, extracting assessments and emotions from their written language referring to entities such as services, products, events, attractions.

The term *Sentiment Analysis* first appeared in Nasukawa and Yi, 2003 [71]. With the growth of social networks and the amount of texts such as reviews, forums, discussions, the expansion of this research area and the great application of the SA in many social, economic, political and business areas has led. For example, consumers want to know other users' reviews about a product before buying it and voters want to know about a political figure before making a voting decision for the election.

This volume of text opinions is not easily decipherable, especially for long texts, and a human reader has difficulty identifying the most relevant information. Furthermore, humans do not consider clear criteria for interpreting and classifying the sentiment of a text: it is estimated that different individuals agree only about 60-65% of the time when evaluating sentiment for a particular text: therefore an automated SA system is necessary.

In addition, sentiment analysis can be used to identify critical information during specific scenarios in real time. Therefore, the advantages of SA can be summarized in Scalability, Coherent Criteria and Real-time Analysis. Moreover, to identify the opinion, the SA systems classify the attributes of the expression as *Polarity*, that is if an opinion is positive or negative and the *Object*, that is what we are talking about, the intensity of that opinion and the relevance of the object of analysis with respect to the context.

There are different types of sentiment analysis and the SA tools are mainly distinguished in systems that focus on polarity (positive, negative, neutral); in systems that detect feelings and emotions (angry, happy, sad, etc.) and those that identify intentions (e.g. not interested) [51],[60]. If it is important to accurately distinguish the level of polarity of the opinion, we consider the analysis of the fine-grained sentiment, which considers the following categories:

- Very positive
- Positive
- Neutral
- Negative
- Very negative.

It could be mapped to a 5-star rating, for example: Very positive = 5 stars and Very negative = 1 star.

Sentiment analysis is studied at three levels:

- Document level: in which the positive or negative opinion of the entire document concerning a single entity is classified, but it is not applicable to documents that express comparisons and assessments on multiple entities.

- 934 • Sentence level: in which the positive, negative or neutral opinion of each sentence is
935 assessed.
- 936 • Entity and aspect: o level of functionality, which is made up of a feeling (negative or
937 positive) and a target (of opinion). For example, in the sentence “The call quality of the
938 mobile X is good, but the battery life is short”, two different aspects of the same mobile
939 phone with different feelings (positive and negative) are emphasized.

940 Pre-processing of the text is necessary as the texts often contain grammatical errors or the
941 same word can be written in different ways or abbreviated by different individuals. If you want
942 to compare opinions relating to the same entity through multiple comments written in different
943 languages, the procedure usually used is to automatically detect the language in the texts, then
944 form a customized model for a chosen language (e.g. English) and finally execute analysis.

945 Sentiment words, also called opinion words, make up the Sentiment Lexicon or Opinion
946 Lexicon. Over the years, some researchers have designed algorithms to compile these lexicons,
947 even if problems are found that are difficult to manage in some cases, such as for sentences that
948 contain sentiment words but do not express any opinion, sentences contain sarcasm, sentences
949 that do not contain sentimental words but may imply opinions; words that in one domain can
950 have a positive meaning but in another negative and finally the presence of opinion spamming
951 used for example to promote or discredit products on behalf of a company by posting false
952 opinions.

953 Short documents are easier to analyze, such as Twitter posts that contain internet jargons
954 and emoticons and that have a limit length of 140 characters, so they are directed to the point.

955 An opinion can therefore be defined as a pair (g, s, t) where g is a target, that is, any entity
956 or aspect of it on which an opinion can be expressed; s is a sentiment on the target, which can
957 be positive, negative or neutral (sentiment orientations or polarity) or a numerical score that
958 expresses the intensity of the sentiment and t is the time in which the opinion is expressed.

959 Since even a single sentence can express more than one sentiment, the algorithms, belonging
960 to the so-called knowledge-based techniques, analyze an opinion of a text add up the sentiment
961 scores of its various terms considering negatives and intensifiers. If the final score is less than 0,
962 the review is negative, otherwise positive. During the classification of sentiment, the opposite
963 words present and the words of negation or adverbs are taken into account as almost and only,
964 which change the orientation of the sentiment.

965 Monitoring the opinion to vary the time of a given entity can help make improvements on it
966 and understand which aspects need to be changed, and also check through the users’feelings if
967 these changes have given effect.

968 In the multi-agent architecture proposed by us in this work, we will use sentiment analysis
969 to analyze the comment that a user enters on each place visited (through his agent) in order to
970 assign an interest score and direct subsequent visit options based on this.

Chapter 4

Problem Analysis

971 In this chapter we describe the large-scale datasets used in our experiments, which show the fea-
972 sibility of our method. They have various characteristics, because they are detected by different
973 types of devices and users (tourists, students, commuter workers). The recorded tracks are real
974 GPS trajectories, which we preferred over CDR raw data because of finer grain (compared to
975 the latter). Then, using this data as a basis for our analysis, we explain the method used. The
976 problem we faced was to find SPs and POIs from big geo-spatial data and, starting from these,
977 to propose a multi-agent system that could satisfy the needs of users to reach places of interest.
978 The real-time data, collected by the architecture in action, is added to those previously saved
979 for a more realistic and accurate result.

980 In the same way we start from the tracks recorded by the taxis, relating to Cabspotting data,
981 to show how the forecast of overcrowded places and the alert to users takes place, applying some
982 tools used by market basket analysis.

4.1 Proposed Datasets

983 The GeoLife dataset presented in Section 4.1.1 is used for the extraction of Points Of Interest and
984 the identification of flows of people, also through a more efficient strategy for finding corridors.

985 Taxis and trucks of the next two sections confirm the usefulness of the approach presented
986 for the SPs and flows and the veracity of the results found, as the POIs are verified by the nature
987 of the dataset.

988 The flow detection algorithm is performed on the Cabspotting data (dataset described in
989 Section 4.1.4) and we used it mainly as a basis for the statistics for the co-occurrence of highly
990 frequented areas. With this latest test, we show how to avoid crowding by predicting the massive
991 movement of people to stationery areas.

4.1.1 GeoLife

992 The GeoLife dataset was collected for a Microsoft Research Asia project for more than four
993 years (from April 2007 to October 2011) and is freely downloadable from their website [42]. It

994 consists of a set of GPS trajectories, ordered sequences of timestamped points with latitude and
995 longitude, recorded by 182 users during their daily movements.

996 The 17,621 total trajectories were recorded by cell phones or other portables GPS devices, for
997 a total distance of 1,292,951 km with a total duration of 50,176 hours. Most of the trajectories
998 were logged in a dense representation, e.g. every 3 seconds or every 5~10 meters per point.

999 73 of the users participating in the GeoLife project had labeled their trajectories with the
1000 mode of transport: bus, taxi, car, foot, bike, subway, train, plane, boat, race and motorbike.
1001 The registered routes cover more than 30 cities in China, but also some cities of Europe and
1002 USA, but most are distributed over Beijing.

1003 This dataset collects users' outdoor movements, such as home-to-work life routines or trips
1004 for entertainment and sports activities (shopping, sightseeing, restaurants, hiking and cycling).

4.1.2 Taxi

1005 Taxi trajectory data is the GPS trajectories collected by taxis equipped with GPS sensors in
1006 Beijing during a period from October 30, 2010 to November 30, 2010. The sampling rate was
1007 one point per minute. Currently, a small part of Taxi (101 trajectories) are available [91]. The
1008 sampling rate was one point per 60 s, and a trajectory has 39,100 data points on average. Each
1009 sample records detailed information of taxis including taxi ID, location, time stamp, occupancy
1010 indicator, orientation and speed. The time at which each GPS point was recorded appears as a
1011 timestamp in Unix time format. Unix time (also known as Epoch time) is a system for describing
1012 a point in time. This is the number of seconds that have passed since the Unix era, which is
1013 00:00:00 UTC on January 1, 1970, minus leap seconds.

1014 The dataset has good temporal and spatial coverage, including weekdays, weekends and pub-
1015 lic holidays of all the main urban areas in Beijing (longitude and latitude in [115.421387; 39.437614]
1016 x [117.321785; 40.609333]).

4.1.3 Truck

1017 Truck trajectory data [91], referred to as Truck, is the GPS trajectories collected by trucks
1018 equipped with GPS sensors in China during a period from August 2015 to October 2015. The
1019 sampling rate varied from 1 s to 60 s. Currently, a small part of Truck (101 trajectories) are
1020 available, and they mostly have around 50 to 90 thousand data points. Every GPS point is
1021 represented by latitude, longitude and timestamp in Unix time format.

1022 The space covered by the registered paths is the large area [86.882817; 0.230753] x
1023 [172.467424; 43.405276].

4.1.4 Cabspottingdata

1024 This dataset [74] includes the trajectories collected in May 2008 by 536 taxis, for a total of
1025 11,219,424 GPS points. Mobility traces are provided by the Exploratorium - the museum of
1026 science, art and human perception through the cabspotting project (<http://cabspotting.org>)

1027 To gather data each vehicle was outfitted with a GPS tracking device that was used to
1028 efficiently reach customers. The data were transmitted from each taxi's cab to a central receiving
1029 station, and then delivered in real-time to dispatch computers via a central server.

1030 For each mobility trace file, associated to a taxi ID, contains in each line: (latitude, longitude,
1031 occupation, timestamp). Where latitude and longitude are in decimal degrees, the occupation
1032 indicates whether a taxi has a passenger (1 = busy, 0 =free) and the time is in the UNIX era
1033 format.

1034 The area covered by these routes corresponds to the county of San Francisco of USA and its
1035 surroundings in California, with minimum and maximum longitude and latitude of the range
1036 $[-127.08143; 32.8697] \times [-115.56218; 50.30546]$. The total size of the trajectories registered with
1037 a customer on the taxi consists of 5,017,659 points.

4.2 Proposed Approach

1038 When an individual travels to an unknown city, he looks for information about the most inter-
1039 esting sites, and the ways to follow to connect them. A simple map could confuse him, therefore
1040 a user through his mobile device searches very often for updated information.

1041 This is the motivation of our work regarding the extraction of GPS points and user cus-
1042 tomizable routes.

1043 In this section, attention is given to finding significant user locations, this analysis is possible
1044 thanks to the growing availability of a large amount of data regarding individual trajectories.

1045 To achieve this, the data to be analyzed mainly comes from cellular or mobile GPS devices,
1046 which allow you to keep a track of the history of positions by GPS paths. The first tests were
1047 carried out on the GeoLife dataset provided by Microsoft Research Asia.

1048 GPS points provides fine-grained information on the spatio-temporal trajectory with respect
1049 to, for example, CDR (Call Detail Record) data, but also they are more difficult to manage. In
1050 fact, given the huge amount of these data to be analyzed, it is not easy to process them through
1051 algorithms with acceptable times.

1052 For this research we have chosen the range of longitude and latitude of $[116.1; 39.7] \times$
1053 $[116.7; 40.13]$ (Beijing metropolitan area). In this area, of 51 per 48 *km*, there are 18,401,631
1054 GPS points.

1055 By plotting the trajectories we can see a second problem: some paths appear broken / not
1056 continuous probably due to the presence of buildings or tunnels that disturb the GPS signal.
1057 Therefore in some areas the GPS device did not register, but we have decided not to tackle this
1058 problem.

1059 Before analyzing the GeoLife paths, we performed a cleaning in order to remove some of
1060 its inconsistent data because these may influence the study results. For each trajectory, after
1061 having computed the instantaneous speed of each point with the consecutive one, we considered
1062 the GPS points with a speed between 0 to 100 *m/s*.

1063 After having computed this, there were some undefined speeds, because the time was zero
1064 probably due to the GPS device stops working properly not recording the elapsed time: with
1065 the speed filter we eliminated the points with these anomalies. The total number of recordings

Table 4.1: Information about different time slots of GeoLife dataset.

Time Slot	Total number of		
	<i>GPS Points</i>	<i>Trajectories</i>	<i>Users</i>
1	3,978,234	5,878	156
2	3,729,429	4,302	150
3	4,976,744	6,613	166
4	3,107,232	4,702	168
5	889,076	1,505	114
6	1,341,196	2,537	129

we have worked on was 17,976,308 points, therefore from the area considered we have deleted 425,323 points.

GeoLife dataset was published in .plt, a vector graphic format, converted by us in .csv files to be processed.

The first step of this work was to build a rectangular grid of dimensions compatible with the maximum and minimum latitudes and longitudes previously chosen: in order to have an accurate statistic we decided to divide our rectangular grid into intervals of 100, 200, 400 and 800 meters.

Since, approximately, for the whole area under examination 100 meters of distance between two points at the same longitude correspond to a difference of latitude of 0.0009 degrees (direction 90 degrees), while for points on the grid with same latitude 100 meters are about 0.0012 degrees of longitude difference (direction 180 degrees), we can give the appropriate divisions of the grids.

Then, we analyzed how many points of different paths fall into our boxes transforming our grid into a matrix M , where every element of M contains the total occurrences of GPS points that fall in that square (100 * 100 meters for example) and the position of the element $m_{i,j}$ corresponds to a certain square of the created grid.

Through an histogram we can evaluate which subsets of area the greatest traffic of users. The contour plot (graph of the isolines) of the matrix M allows you to immediately go back to the main “hotspots”, that are the areas more visited.

From Figures 4.1 and 4.2 we can see some hotspots, around the Tsinghua University, at Yuanmingyuan Park, the Hepingli Residential District and Peking University.

We divided the trajectories into 6 time slots from 4 hours each Slot1 = [00:00:00, 03:59:59], Slot2 = [04:00:00, 07:59:59] and so on, in order to analyze the traffic in Beijing during different time slots (Figures 4.3, 4.4, 4.5, 4.6, 4.7, 4.8). More details are visible in Table 4.1.

An extensive series of experiments in the other two datasets, Taxi and Truck, has been performed in order to study the movements of users in different situations. We mined shared routes in real-word and not only drop-off or pick up taxi. The heatmap of Taxi trajectories is visible in Figure 4.9.

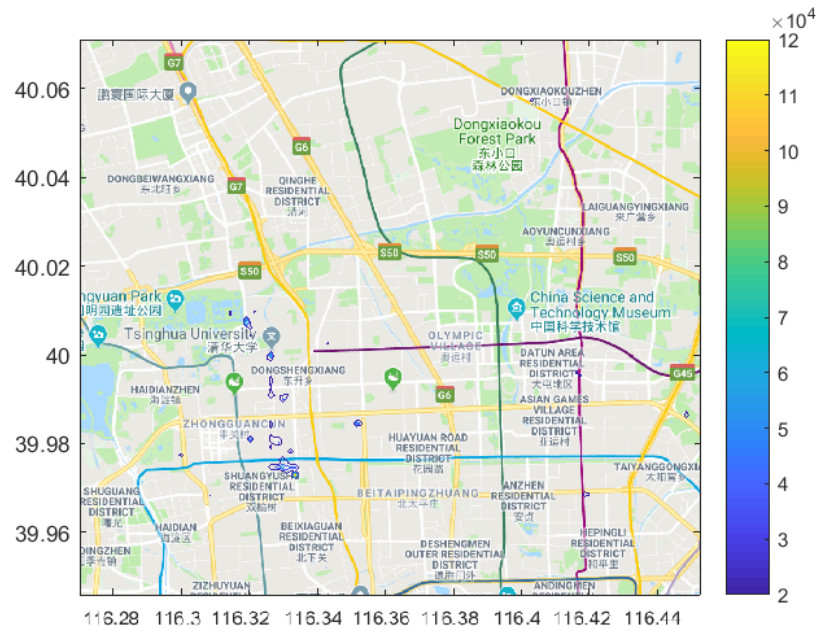


Figure 4.1: Contour plot created with grid with intervals of 100 meters, zoom on map.

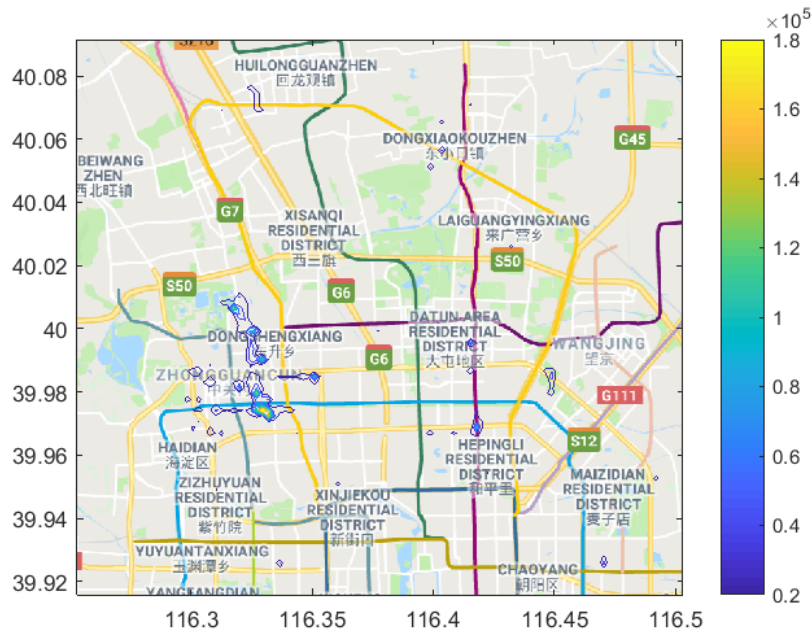


Figure 4.2: Contour plot created with grid with intervals of 200 meters, zoom on map.

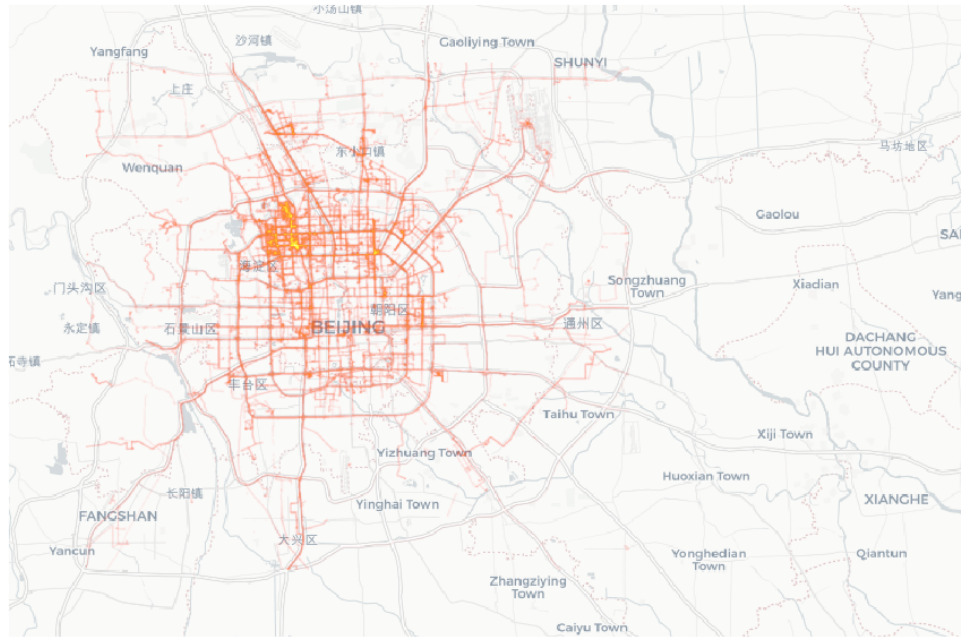


Figure 4.3: Trajectories in Beijing for time slot 1: [00:00:00, 03:59:59].

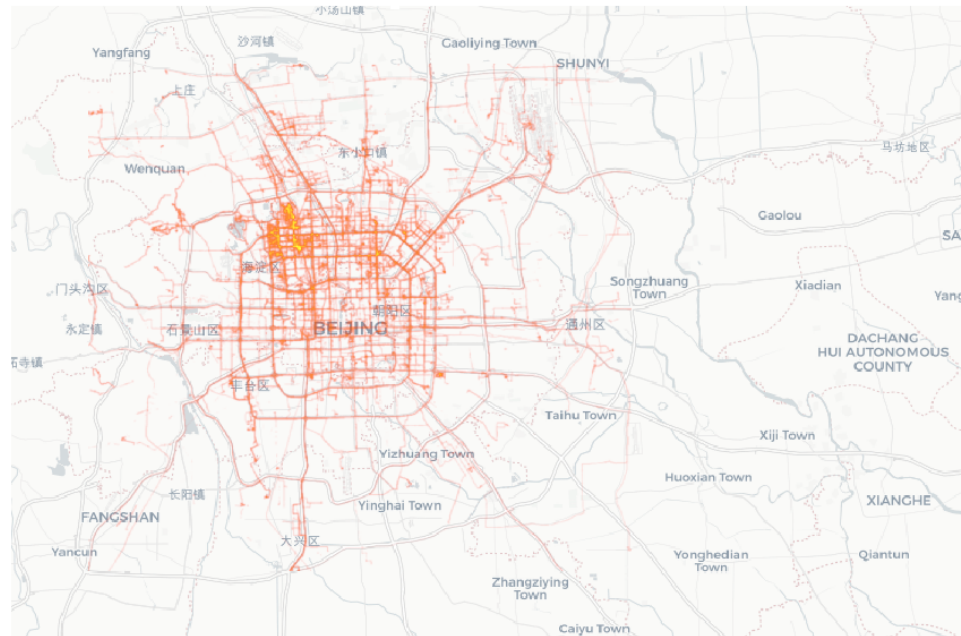


Figure 4.4: Trajectories in Beijing for time slot 2: [04:00:00, 07:59:59].

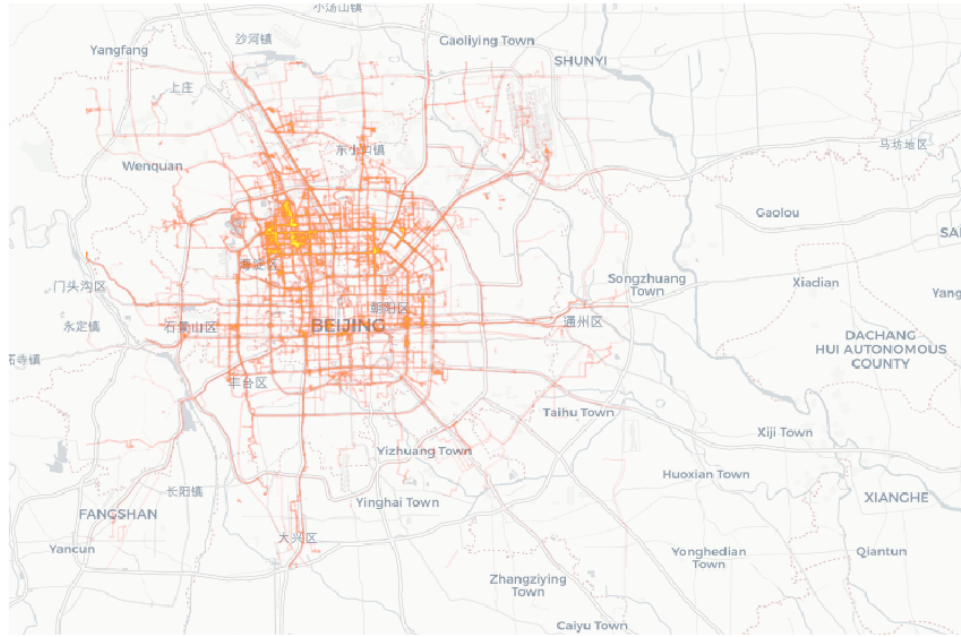


Figure 4.5: Trajectories in Beijing for time slot 3: [08:00:00, 11:59:59].

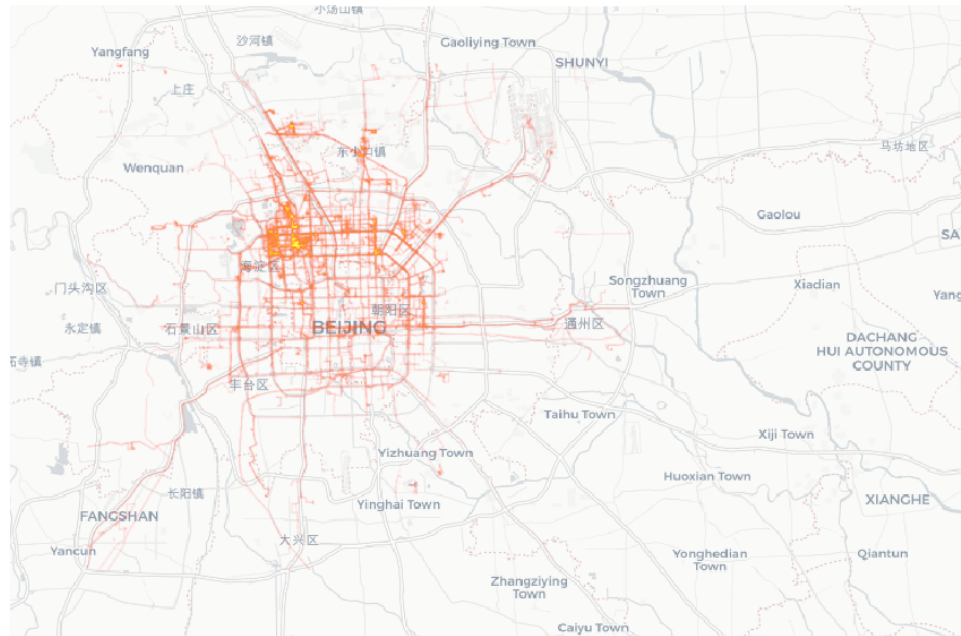


Figure 4.6: Trajectories in Beijing for time slot 4: [12:00:00, 15:59:59].

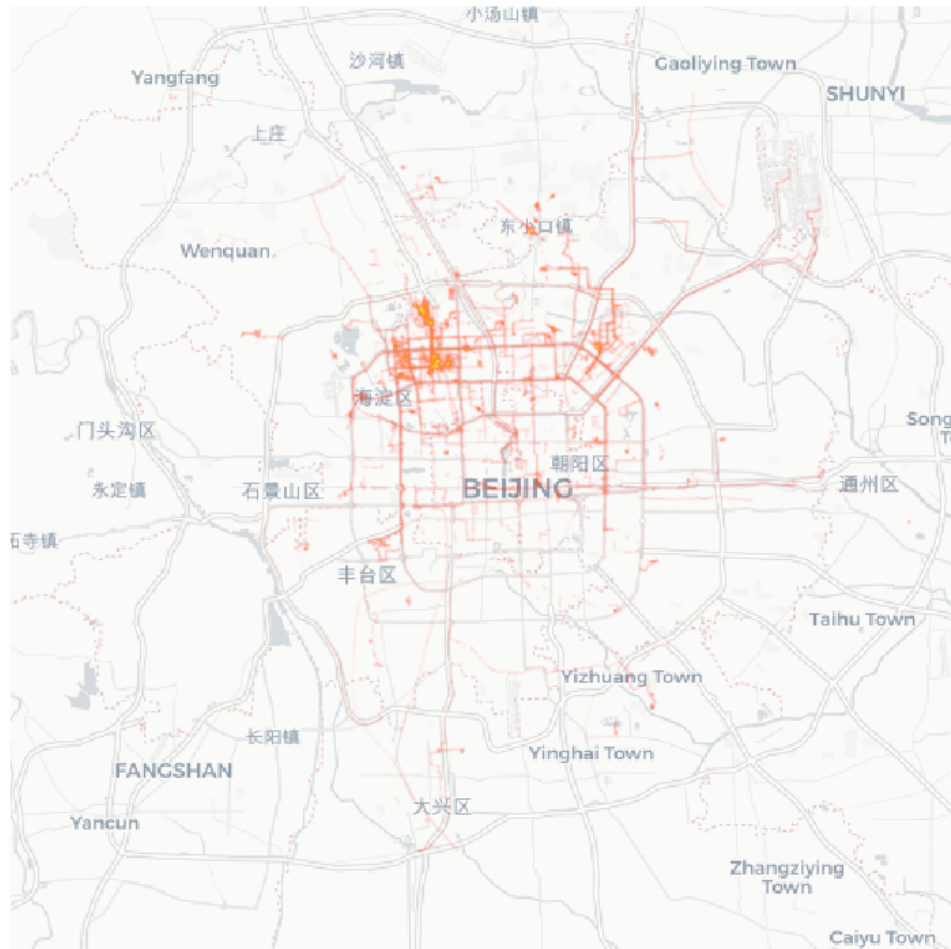


Figure 4.7: Trajectories in Beijing for time slot 5: [16:00:00, 19:59:59].

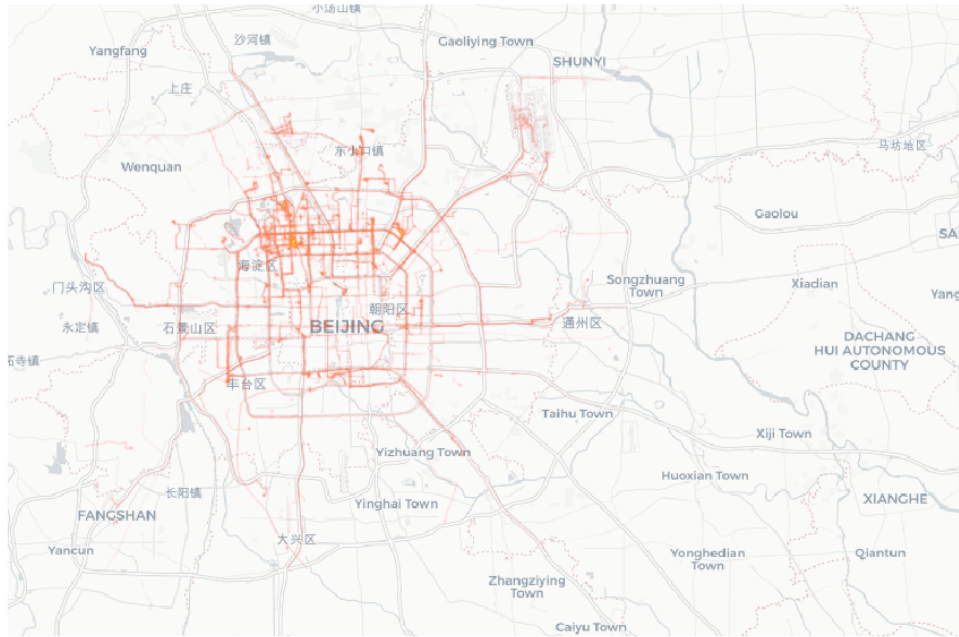


Figure 4.8: Trajectories in Beijing for time slot 6: [20:00:00, 23:59:59].

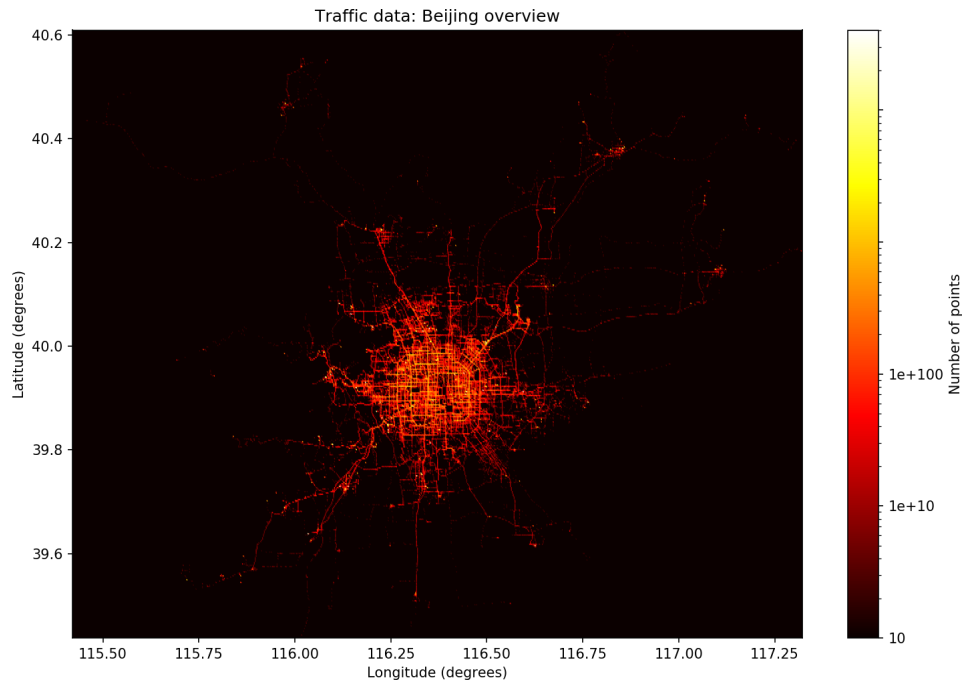


Figure 4.9: Heatmap of trajectories in Taxi dataset.

4.2.1 StayPoints extraction

1094 For every slot of time in GeoLife dataset, after grouping trajectories by users, the second step
1095 of our work was the StayPoints (*SPs*) detection [73]. When we find a region in which a user
1096 has spent a considerable time on its surroundings, the centroid (the mean of coordinates of the
1097 points belong to it) of this cluster represents an *SP*. The algorithm that we implemented for
1098 the *SP* detection (pseudocode in Algorithm 3) needs as input a time threshold (*TimeThr*) and
1099 a distance threshold (*DistThr*). For example, if an individual stays over 20 minutes (*TimeThr*)
1100 within a distance of 200 meters (*DistThr*), a *SP* is detected.

1101 The Algorithm 3 analyses each trajectory individually, scanning its points in the order of
1102 registration. For each, considering the pairs between it and its successors, calculate the spatial
1103 distance between the two points p_i ; and p_j ; and the time span between them. All points
1104 are located in an area enclosed by the *DistThr* distance threshold and for which the time span
1105 is greater than *TimeThr* form a *SP* region. The average of their coordinates determines the
1106 centre of the area, and represents a StayPoint (*SP*) of the analysed trajectory. This captures
1107 the behaviour of the user, who is probably going around some building for a while. We want to
1108 study mainly the parts of the path in which he dwells for some reason in a spatial environment.
1109 We also extract the *arrival time* information for it, i.e. the time relative to point p_i ; and the
1110 *leaving time*, i.e. the time relative to point p_j .

1111 With reference to the log routes, each *SP* has a particular semantic meaning such as the
1112 restaurant where we go and the places we visit, etc. Figure 4.10 shows two categories of Stay-
1113 Points. The first situation, such as *SP* 1, occurs when an individual remains stationary in a
1114 single point for a period of time greater than a threshold. In most cases, this state occurs when
1115 people enter a building and the satellite signal is lost for a certain amount of time until arrival
1116 of the outdoor signal again.

1117
1118 In the other situation, like *SP* 2, a user wanders within a given spatial region for a period
1119 of time greater than the threshold considered. Consequently, we have to calculate the average
1120 coordinates of the GPS points of that spatial region and the final *SP* point is represented by its
1121 arrival time and its left time. In most cases cases, this is a situation arises when people travel
1122 and are attracted by the surrounding outdoors.

1123 If we group through clustering some significant places (*SPs*) by setting a minimum number
1124 of points necessary to form the cluster sufficiently, the few GPS points will be excluded like
1125 point 1. Therefore, if they are not surrounded in the vicinity by other significant places, the
1126 density of the points recorded there cannot satisfy the condition for the formulation of a cluster
1127 and instead the *SPs* “type 2” will be extracted.

4.2.2 Solution for Points Of Interest detection

1128 Points of Interest, commonly abbreviated POI, are a well-known concept in literature [4], [52],
1129 [102]. A POI is defined as an object associated with a latitude and a longitude which at least
1130 one person would reasonably be expected to have an interest or an utility. POI recommendation
1131 suggests places to visit, which are taken from an automatic analysis of the three real dataset:

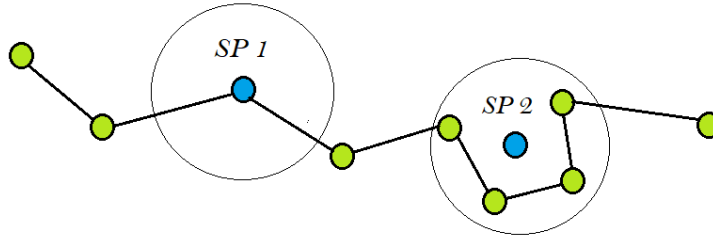


Figure 4.10: Example of the two types of StayPoint.

Algorithm 3: Pseudocode for SPs detection

Input A trajectory $T = \{p_i, p_{i+1}, \dots\}$, a distance threshold ($DistThr$) and time span threshold ($TimeThr$)

Output A set of StayPoints $SP = \{S\}$

$i = 0$, $cardinality_traj = |T|$

while $i < cardinality_traj$ **do**

$j = i + 1$;

while $j < cardinality_traj$ **do**

$dist = distance(p_i, p_j)$;

if $dist > DistThr$ **then**

$\delta time = time_{p_j} - time_{p_i}$;

if $\delta time > TimeThr$ **then**

$S.coords = MeanCoords(\{p_k | i < k < j\})$;

$S.arrive_time = time_{p_i}$;

$S.left_time = time_{p_j}$;

$SP.append(S)$;

break;

$j = j + 1$;

$i = j$;

return SP

1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152

Taxi trajectories data, GeoLife and Truck trajectories data.

After the SP detection, we focused on POIs that cluster together SPs of different users (at least 10), and checked if, in different time slots, the users that previously had a common POI move together to another one. We applied DBSCAN, to the SPs obtained, as it works well with large geographical dataset and likewise it can be adapted for any distance function.

For our data, the clustering algorithm DBSCAN has determined clusters for all SPs. For example in Geolife dataset we set MinPts equal to 10 or 15 and ϵ from a minimum of 200 meters to a maximum of 400 meters.

Geospatial clustering must depends on geographic information domain knowledge and the context of the users. For example 200 meters was chosen as a spatial parameter because it is the average space between an intersection or a square around an attraction.

Both parameters control the local neighbourhood of the points: making better use of geospatial and clustering knowledge to select suitable constraints and parameters is likely to yield better and more meaningful clusters.

For example, applying another partition clustering method such as K -Means to the whole GPS Geolife dataset, we obtain a result like Figure 4.12. The K spatial partitions obtained will depend on the number K chosen and we can rely on the Calinski-Harabasz index for this. As shown in Figure 4.11 on the left, this criterion suggests us to use the highest value of the plot, i.e. $K = 9$. This value was confirmed as optimal for the division of different zones also by the Elbow method, as indicated on the right of Figure 4.11; the execution time for this clustering algorithm on our data was 1 hour and 27 minutes.

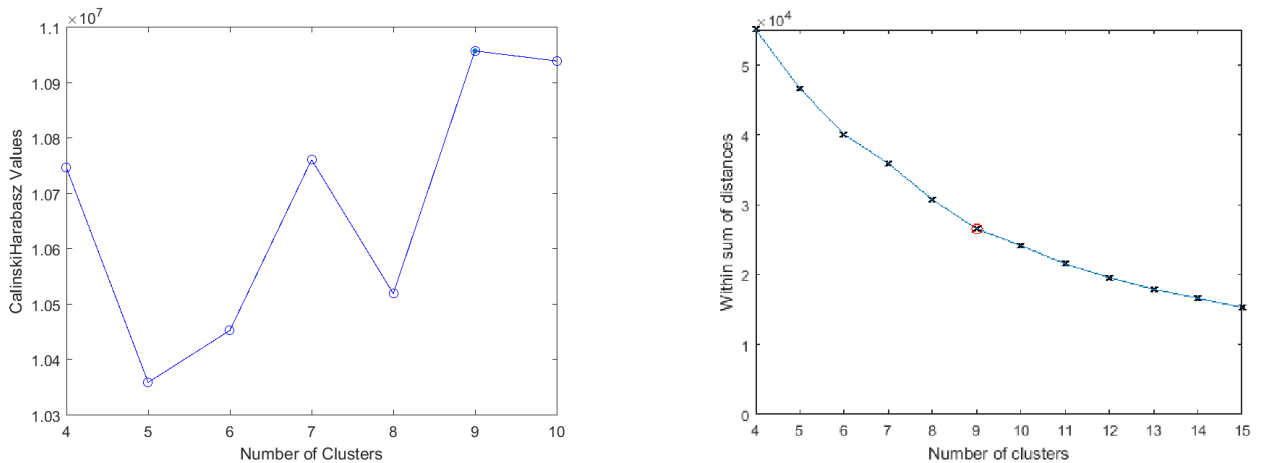


Figure 4.11: For Caliski-Harabasz criterion (shown with the plot on the left) the optimal number of clusters for the set of GeoLife GPS points is 9.

The Elbow method is a heuristic used in determining the number of clusters in a dataset. The method consists of plotting the explained variation as a function of the number of clusters, and picking the elbow of the curve (on the right) as the number of clusters to use, in this case $K = 9$.

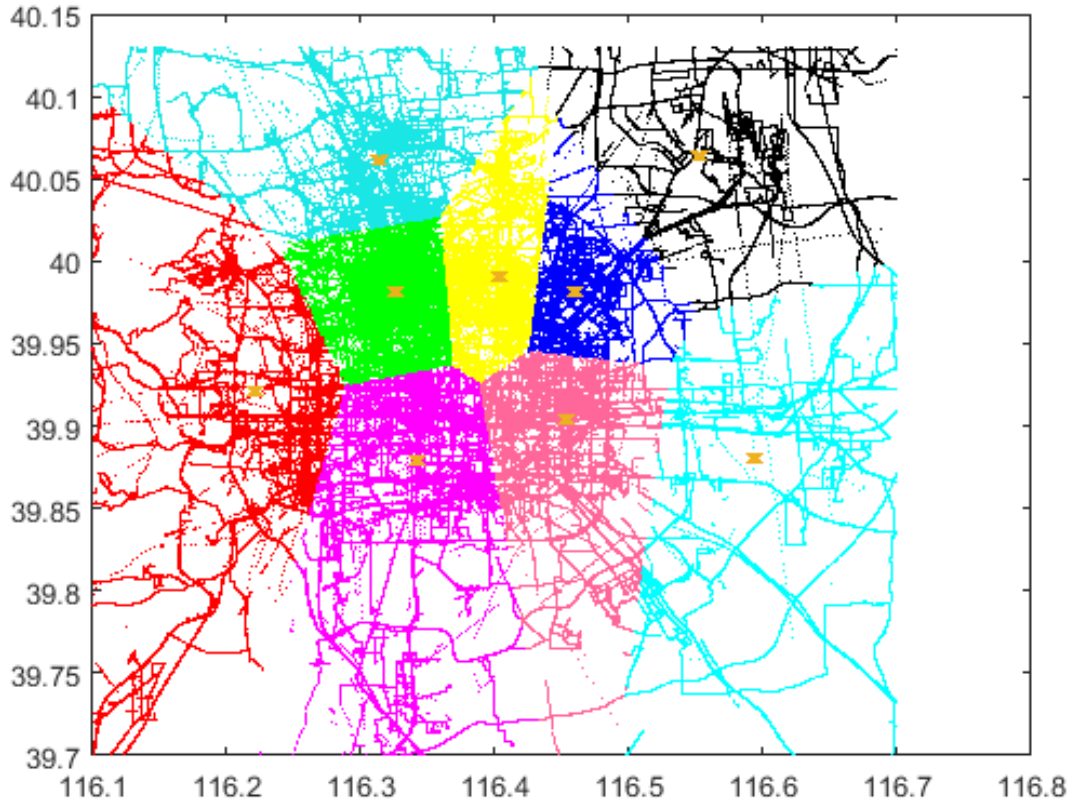


Figure 4.12: Cluster assignments and centroids, K -means with $K = 9$.

1153 So, we decided to cluster not directly on the total GPS points because it is computationally
 1154 expensive and we also want to consider only those of greater significance in the trajectories,
 1155 therefore we apply the DBSCAN algorithm based on density only on the StayPoints.

1156 The points found by our analysis were verified by matching the results with Google Maps
 1157 data. It was confirmed that they correspond to real POIs, i.e. parks, restaurants, etc., hence
 1158 validating our approach.

4.2.3 POIs Recommendation (proposed Multi-Agent system)

1159 In our proposed software architecture a centralized server collects and analyzes data from dif-
 1160 ferent agents, with the aim of offering suggestions to users and some real-time data on visiting
 1161 POIs. In addition, each user (with reference to an application) was modeled as an independent
 1162 agent who communicates with the centralized server. In this approach, an agent is an applica-
 1163 tion installed on an Android device that tracks the user movements around Points Of Interest,

1164 in order to suggest new places to visit. Therefore, our recommendation on POIs is based on a
1165 multi-agent system that performs the following steps:

1166 STEP 1: the server sends to agents the list of known POIs for a city. Each POI was
1167 previously determined by the DBSCAN algorithm discussed in Section 3.5.2. Each POI has
1168 the following information: the most visited time slots; the number of agents present on the site
1169 (POI) in real time; a series of feedback on the site, created by a Sentiment Analysis algorithm
1170 that analyzes the comments released by users; a rating estimated from previous information
1171 that recommends (or dismiss) the POI to the user. The union of this information is used to
1172 drive the recommendation system.

1173 STEP 2: thanks to the evaluation, the agent chooses a place of interest that the user can
1174 visit. As soon as the GPS coordinates are within a radius of less than one kilometer from
1175 the coordinates of a POI, then the location will be sent to the centralized server, which can
1176 determine the number of users present. GPS coordinates are sent only in controlled areas to
1177 preserve user privacy.

1178 STEP 3: once the visit is over, the user can use the agent to comment on the place visited.
1179 This information will be sent to the server, and analyzed there using Sentiment Analysis algo-
1180 rithms [61], which in turn allows the server to determine a score that identifies whether the POI
1181 was satisfactory for the user.

1182 For each POI visited, the data forwarded to the server are: the inaccurate GPS position of
1183 the device, visiting hours (through a timestamp), the name of a place, an evaluation score of
1184 the place (i.e. a number from 0 to 5). To forward such data, the user needs to authorise the
1185 application to share his GPS position.

1186 STEP 4: the agent will receive the list of POIs again adding information for a specific POI
1187 related to the site already visited based on previous users experiences. One of the aspects
1188 considered is the sequentially: we take into account not only the places visited, but also the
1189 sequence of them. Therefore if a group of users who are considered similar to Claire have visited
1190 POIs *A-B-C-D* and she is on *C*, then the system will advise her to visit *D* or *B*.

1191 Finally, these points are sorted according to the user current distances. E.g. if there is a POI
1192 near her in which the number of people is adequate (i.e. according to the preferences of the user),
1193 the time slot matches the current time, and users have appreciated the visit (having assigned
1194 a positive score), then such a POI is suggested to the user by alerting her agent. Otherwise, if
1195 the place has not been appreciated by other users or the current time is not within the time slot
1196 referred by available data, then the place is labelled as not recommendable.

1197 Agents accept suggestions from the server and filter them, since the exact location of the user
1198 is known only to the local agent. The suggested points, also given to the experiences acquired
1199 by other agents, are sorted based on the user's current location and filtered for the time slot
1200 corresponding to the current time.

1201 Moreover, points not yet visited by the user are marked. As for the time slots, the hours
1202 in a day were organised into six 4-hour slots; Slot1 = [00:00:00, 03:59:59], Slot2 = [04:00:00,
1203 07:59:59], Slot3 = [08:00:00, 11:59:59], etc. Therefore, data gathered for a POI are associated
1204 with a time slot based on the start time of the visit. E.g., if the visit starts at 10 a.m. this point
1205 will be labeled in the 4th time slot.

1206 Such cooperation can generate benefits for groups of people who share the same interests
1207 (tourists, students, etc.). Thanks to the exchange of information between agents and servers, it
1208 is possible to define the rules for our POI-based system recommendation.

1209 The most important properties of the agents in this proposed recommendation system are
1210 the following:

- 1211 • Agents are independent: if an agent stops working the others continue their work without
1212 consequences.
- 1213 • Agents are mobile: they are easily transportable thanks to their integration in the mobile
1214 application.
- 1215 • Agents are reliable: given that the GPS coordinates are taken directly from the device,
1216 avoiding to obtain false data.
- 1217 • Agents work by preserving user privacy: there is no continuous sending of the users posi-
1218 tions, instead the server only receives data when agents are near POIs, and the coordinates
1219 sent to the server are displaced by a bounded random amount, thus preserving the users
1220 privacy. More specifically, two important steps are taken to offer greater security to users:
1221 (i) users share the location only if they are close to a well-known POI, i.e. their position
1222 falls within 300 *m* from the POI); (ii) the position is processed through data masking
1223 techniques, i.e. a variable and unknown displacement will be added before being sent to a
1224 central server storing it. Having altered the position does not affect the recommendation
1225 system and guarantees better protection for the user.

1226 In this context, security and robustness are features embedded in all the agents, working
1227 together by means of a centralised server. Thus, thanks to a system based on multi-agents, and
1228 a centralised server, it is possible to create a recommendation system based on the experience
1229 and appreciation of users, able to suggest points of interests to users, while ensuring user privacy.

1230 The agent takes some decisions on behalf of the user, i.e. the agent determines:

- 1231 • when geolocation data are sent,
- 1232 • which displacement is added to geolocation data to preserve user privacy,
- 1233 • how to rank incoming data.

1234 For the latter, since the agent shields the user preferences to the outside world, hence the agent
1235 is entrusted with the task of muting some alerts when the user has low interest for it, or they
1236 are not related to the actual flow of the user, and with the ranking of incoming comments and
1237 selected places of interests according to match user preferences. Each agent sends data to a
1238 centralised server, and is identified by an app id that the server has provided. User identity
1239 is preserved since: the server does not disclose such ids to the agents; agents can not directly
1240 communicate with each other; the user is not providing his personal data.

4.2.4 Privacy Preserving

1241 In the literature, user privacy is a topic much discussed. The unconditional use of smartphones
1242 and the multiple downloaded applications put a strain on the protection of user rights. Existing
1243 systems have mainly taken three approaches to improving user privacy: (a) introducing uncer-
1244 tainty or error into location data [41], [69], (b) relying on trusted servers or intermediaries to
1245 apply anonymisation [55],[57], and (c) using cryptography techniques [70], [72]. However, each
1246 system has weaknesses: the first approach loses accuracy because the uncertainty rate degrades
1247 the data, the second approach can be risky, since private data is exposed to proxies that could be
1248 violated and the third approach is often computationally expensive. Hence, hybrid approaches
1249 are often a proper compromise. This work offers a hybrid system while preserving privacy in
1250 two ways:

- 1251 1. by using a centralized server the agents do not directly know each other and can not
1252 exchange information (as for systems (b) above);
- 1253 2. by only extracting information on POIs, and an approximate position of the user, only
1254 when he expresses an opinion on a POI (as for systems (a) above).

1255 In addition, this approach is the first approach in the tourism sector that deals with preserving
1256 user privacy. There is no user ID, no precise position is provided to the other participants. Users
1257 interact with the system a centralized and not directly with other users. Allowing the direct
1258 exchange of information between the various devices, which are in a certain area, would make
1259 known the user position that we do not want. To ensure privacy, therefore, the details -as the
1260 precise position- are obscured, We considered more disadvantages to leaking private information
1261 than the advantages.

1262 Data collected are processed and cleaned up in the user device, in order to respect anonymity,
1263 and then processed data are sent the server. In detail, each agent: analyses only user data located
1264 in the vicinity (about one kilometre) of known points of interest and adds a displacement to the
1265 data before making it available to the server. In this way, it preserves the user's privacy in two
1266 ways: it does not always share the user's position to the server, and when sharing the position,
1267 it is affected by a variable uncertain amount.

4.2.5 People movements identification

1268 It is difficult to find significant user behaviors from the available data. In this section and in
1269 the following one, great attention is paid to identifying significant paths for users. The Section
1270 4.2.7 aims to predict user behavior.

1271 Many people flows are discovered on the datasets analyzed before and these are useful,
1272 e.g. to find common itineraries, to suggest places that can be reached by users or to present
1273 improvements on public infrastructure [22].

1274 The proposed approach was validated by three experiments on location data, to arouse
1275 significant places, starting from the data available on: taxi, truck and people movements.

1276 We use an alternative technique, by means of significant flows of people within the city, to
1277 identify POIs. It can detect people or taxis flow to extract POIs with a high flow of interest,

1278 starting from a series of GPS trajectories. The detailed results are presented in Sections 5.2 and
1279 5.3.

1280 We define two *spatially similar* trajectories if they cross the same n points (unless of a
1281 certain tolerance) in the same temporal order. The distance between two points is defined in the
1282 Haversine formula. A *flow* is a continuous and uniform movement of entities (such as people,
1283 cars or truck) in one direction and a *flow density* is the number of entities traveling through the
1284 same flow. Finally trajectories (spatially similar) are called temporal similar if entity movements
1285 occur in the same time slot. The first necessary step is the preprocessing of data, that is, each
1286 trajectory is cleaned up by removing the noise and outliers. We also cleaned the trajectory
1287 from points too close to each other, so you can work with others homogeneous trajectories.
1288 Outliers removal compares each point with the previous and next points, looking sharp speed
1289 changes. In case of abrupt changes, this point it is an abnormal label and is discarded. Once the
1290 data preprocessing phase is complete, a the algorithm that calculates all flows is executed their
1291 density, finding spatially similar trajectories. The search for comparison was done in parallel in
1292 order to speed up the procedure.

1293 Using this approach, relevant flows are extracted in the city. Starting from these shared
1294 paths it is possible to extrapolate POIs by making a match between the POIs of the analysed
1295 city (see Section 5.3) and the points close to each detected flow (with a maximum distance of
1296 100 m).

4.2.6 Mining corridors from GPS trajectories, through Apriori

1297 In this part of work we performed different tests in order to detect shared routes, so we applied a
1298 grid of uniformly sized cells to discretize the trajectories. The collection of trajectories D of the
1299 original dataset is transformed into a new collection of trajectories D' , where each trajectory is
1300 represented as a sequence of grid cells. If two or more consecutive coordinates are mapped into
1301 the same grid cell, we report only the first instance, not allowing a trajectory to have consecutive
1302 points of the same grid cell.

1303 In this case, we have considered that trajectories are not characterized by traveling direction:
1304 going forward or backwards on a sequence of cells results in two instances of the same trajectory.
1305 Hence, trajectories can be seen as a path in an undirected graph.

1306 Given this trajectory representation, we can define what is a *corridor*:

1307 **Definition 4.2.1.** A (M,S) -corridor is a sequence of M cells in D' such that all cells are shared
1308 by at least S trajectories in a given time slot.

1309

A first naive approach to detect corridors is given by this pseudocode:

Algorithm 4: Brute force approach to find corridors

Procedure: Find_Corridors($(T_1, T_2, \dots, T_L), M, S$)

Input: A set T_1, T_2, \dots, T_L of L trajectories formed by N cells each, M : the minimum number of cells for a corridor, S : the minimum number of trajectories for a corridor.

Output: C the set of (M, S) -corridors.

$C = \{\}$

for T *from* 1 *to* $L - S + 1$ **do**

for i *from* 1 *to* $N - M + 1$ **do**

$A = (T[i], T[i + 1], \dots, T[i + M - 1])$

$s = 1$

for $T' > T$ **do**

for j *from* 1 *to* $N - M + 1$ **do**

if $A = (T[j], T[j + 1], \dots, T[j + M - 1])$ **or**

$A(T[j + M - 1], \dots, T[j])$ **then**

$s = s + 1$

end

end

end

if $s \geq S$ **then**

 add s to C

end

end

end

1310

1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

1325

1326

1327

1328

For every trajectory, this algorithm scroll down its elements to decide whether a pair of consecutive cells are shared by at least other $S - 1$ trajectories. If this happens, it must be checked whether the subsequent cell, together with the frequent pair of cells found, constitutes a frequent triple for at least L trajectories. And so on, until you find the largest set of frequent cells.

The complexity of this algorithm is a function of L and N . The number of performed operations is $\sum_{i=1}^{L-S+1} (N - M + 1) * (L - i) * (N - M + 1) * 2$, which results in a complexity $O(N^2L^2)$.

The proposed method comprises two phases: first the generation of *candidate corridors*, which are coarse traces of several distinct trajectories, and then filtering candidates for the selection of *fine-grained corridors*.

First, we analyze the database once to generate a mapping graph on real roads and get information on trajectories. To this end, we define a grid of uniformly sized cells on the map of interest, mapping the coordinates (in decimal degrees) of the points for every trajectory to discrete grid cells. For example with square cells 1 km wide, we get a 51 * 48 cell grid to apply to the selected dataset (see Figure 4.13).

In Figure 4.13, it is possible to notice more intense cells which indicate a greater concentration of trajectories that cross them.

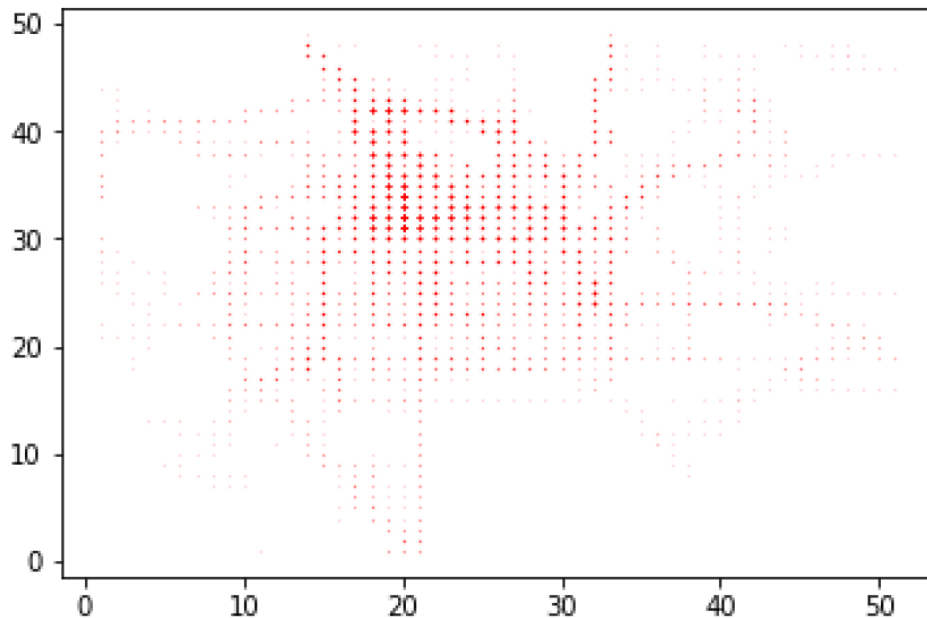


Figure 4.13: Trajectories mapping into a grid with cells of 1 *km* for side.

1329 Based on this trajectory representation, we can generate the candidate corridors to obtain
 1330 information on the areas that approximate the sub-routes common to several users with the
 1331 Apriori algorithm.

1332 In the second phase, this information is filtered using Radius Neighbors Graph and used to
 1333 find fine-grained frequent corridors.

1334
 1335 Apriori algorithm is currently used for example for marketing strategies in supermarket
 1336 chains, so you want to find out which products are often bought together by the customers.
 1337 This frequency-based approach that allows us to find sub-paths of the road network which are
 1338 shared by at least a fixed number of trajectories.

1339 In our scenario this algorithm can be used to reduce computational complexity. The key
 1340 point of its application is to consider GPS points as a set of items and trajectories as a set of
 1341 baskets.

1342 Then the minimum support to be set as input in the algorithm corresponds to the minimum
 1343 number of trajectories that must have in common tuple of cells. The application of the Apriori
 1344 algorithm in the case under consideration can be improved by setting constraints: for example,
 1345 by considering that K -tuples can only be formed by groups of adjacent cells. Extending the
 1346 search for other frequent cells to be aggregated only to the neighborhood located at the edge of
 1347 these cells, since a cell can have at most eight adjacent cells, would lead to the reduction of the
 1348 execution time in the generation of the candidate cells.

1349 We propose three different strategies for the application of Apriori (Step 1, Step 1.2, Step 2

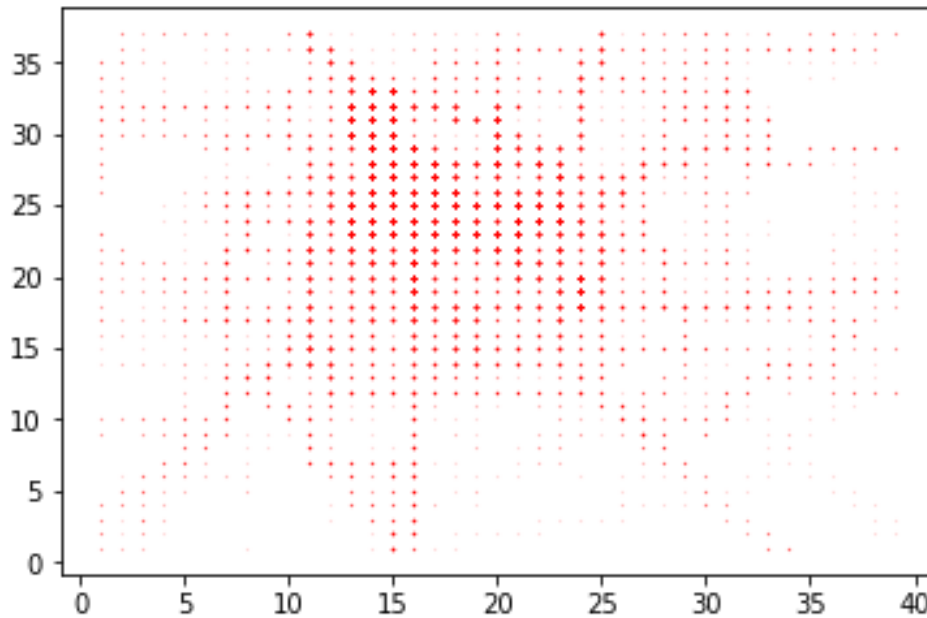


Figure 4.14: Trajectories corresponding to a grid with cells of 1300 meters.

1350 below) to improve execution time significantly. The results are then compared, the experiments
 1351 performed in GeoLife dataset show that for different approaches they lead to overlapping paths.
 1352 We worked on a random sample of 1000 trajectories on the first slot by setting min_support
 1353 $= 50$ trajectories to 1000 or considering all of them, the lowest value was chosen at 15 out of
 1354 5878.

1355 In our approach Apriori algorithm allows the progressive identification of cells that are
 1356 common to several trajectories and the identification of subsets of cells with low support of
 1357 trajectories.

Step 1

1358 The first step of the algorithm is to discretize the trajectories. In our case, we used a grid
 1359 consisting of 1,443 square cells (39 horizontally per 37 vertically) which result in cells whose
 1360 side is 1,300 meters (see grid in Figure 4.14) and about half of these are crossed by at least one
 1361 trajectory.

1362 Each trajectory, that was first considered as a sequence of GPS points, now becomes a list
 1363 of consecutive and distinct cells.

1364 The application of the Apriori algorithm performed in this step is as follows:

- 1365 • First the 1-frequent cells L_1 set is found, scanning the dataset to count number of occur-
 1366 rences of each item. All the individual cells that satisfy the minimum support are counted.

1367 That is, they are crossed by at least S trajectories (one of the parameter of the corridor
1368 detection algorithm).

1369 • From the elements of L_1 , the set C_2 is formed by considering all possible pairs of cells
1370 (candidate frequent pairs of cells).

1371 • To find the L_2 set, the algorithm applies a pruning step, considering among the pairs of
1372 C_2 only those that satisfy the *min_sup* of trajectories (2-frequent cells).

1373 • By increasing the size of the cells at each K iteration of Apriori the sets C_K and L_K are
1374 formed and C_k is generated from L_{k-1} .

1375 • The algorithm stops at job K when it is no longer possible to find a set of frequent cells
1376 of size $K + 1$.

1377 As a result of this step, we have identified the most frequent cell sets in trajectories. It is
1378 important to stress the fact that those cell sets may not be necessarily close in geographical
1379 terms, and for this reason they could represent disconnected paths.

Step 1.2

1380 Once we have identified in the previous step the most frequent cell sets in trajectories, we can
1381 proceed to process the data by reversing the roles of items and bags in the Apriori algorithm.
1382 In this case, we will consider that trajectories correspond to “items” and geographical cells
1383 correspond to “bags”. In this way, at the end we will get the most frequent trajectory sets, that
1384 is, the sets of trajectories that share at least a given number of cells.

1385 In this step we apply the Apriori algorithm on a grid of 500 meters (102 * 96 cells), visible
1386 in Figure 4.15.

1387 By setting the minimum support parameter equal to 51 cells (which make a total length
1388 of 51 * 500 $m = 25.5 km$) and taking into account that the “populated” cells in the first slot,
1389 for example, are 4,603 (cells crossed by trajectories), we have obtained a percentage of support
1390 greater than that of Step 1.

1391 The description of the Apriori algorithm performed in Step 1.2 is as follows:

1392 • The algorithm starts by first identifying all the individual trajectories which satisfy the
1393 minimum support, that is, those which are at least M cells long (one of the parameter of
1394 the corridor): this is the set L_1 .

1395 • From the elements of L_1 with the Join step the set C_2 is formed. That is, all the possi-
1396 ble pairs of trajectories (candidate pairs of trajectories). Candidates are the itemsets
1397 containing all potentially frequent itemsets.

1398 • To form L_2 the algorithm does the pruning, considering among the pairs of C_2 only those
1399 that satisfy the *min_sup* of cells (2-frequent trajectories).

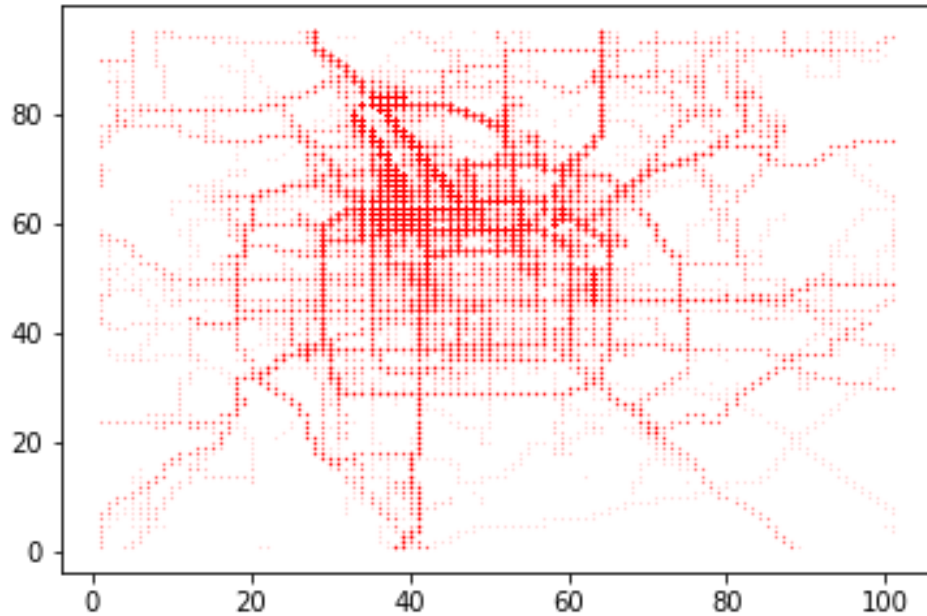


Figure 4.15: Trajectories mapping into a grid with cells square of 500 m wide.

- 1400 • By increasing the number of trajectories at each K^{th} iteration ($K \geq 2$) of Apriori, it
 1401 generates candidate K -trajectories C_K from the frequent $(K - 1)$ -itemsets L_{K-1} of the
 1402 last iteration. Now only the frequent sets L_{K-1} and candidates C_K reside in memory,
 1403 whereas other itemsets of previous iterations are discarded.
- 1404 • The algorithm stops at job K when it is no longer possible to find a set of frequent
 1405 trajectories of size $K + 1$, (that is, when there are no $K + 1$ trajectories which cross at
 1406 least M cells together).

1407 Compared to the previous strategy, in this step, fixed the desired minimum length of the
 1408 corridors in terms of cells, Apriori finds those with size greater than or equal to this and shared
 1409 by an increasing number of trajectories at each iteration of the algorithm.

Step 2

1410 Since we only considered the last level of the output of the Apriori algorithm, now we want
 1411 to analyze the subsets of cells that not appear in the last result of Apriori, the trajectories
 1412 which do not participate in frequent itemsets of the longest cells. In addition, to optimize the
 1413 execution time and vary the results, instead of using all the trajectories of every time each time,
 1414 we considered a subset of 1,000 random trajectories. We chose to build a grid of cells 1 km
 1415 wide, so the total number of cells are 51 (horizontal) * 48 (vertical).

1416 We applied Apriori by following this convention: the role of “baskets” is played by trajectories
1417 and the role of “itemsets” is played by cells. As for the input of Apriori we chose a minimum
1418 support of 50 trajectories on the total number of the random set: the support of a cell is the
1419 percentage of trajectories in which that cell occurs.

1420 This support being greater than that chosen for Step 1 and Step 1.2 leads to an improvement
1421 in the execution time of the algorithm. In fact, the output in this case is the tuples of cells shared
1422 by at least 50 trajectories and so there are less candidate corridors that contain 50 subtrajectories
1423 compared to Step 1 where the minimum support was 15. Moreover, the reduction in the number
1424 of the sample of trajectories also contributes to bringing this advantage. By varying the sample
1425 several times, different results will be obtained.

4.2.7 Probability of movement for GPS points

1426 We aim to determine the probability that users are moving from one point somewhere in a
1427 city to another point. Such a probability is then used to provide recommendations accord-
1428 ingly. Recommendations, alerts, or user requests, are communicated by means of a smartphone
1429 app. Therefore, our proposed solution comprises two parts: (i) an algorithm for the movement
1430 prediction and (ii) an app on the user’s device to track movements and suggest destinations.

A. Determining the movement prediction

1431
1432 To determine the probability of movement, we perform two steps. The first step is to deter-
1433 mine the flows of people shared by the points recorded during each user’s previous movements.
1434 Therefore, the second step consists in obtaining statistics on the amount of people who, having
1435 stopped at point A , move subsequently on to point B .

1436 By analyzing each GPS trajectory, that is the set of recorded GPS points ordered in time,
1437 we have extracted its StayPoints (SP). They are the centers of the areas within which a user
1438 stays for more than a certain time: for some reason that area is of interest.

1439 Then, the geographical area where the SP s of all the datasets are located has been discretized
1440 by means of a grid, composed of equal Square Cells. Each given SP was associated with a single
1441 square cell if contained in that space. Of course, a cell could contain multiple SP s if these are
1442 close enough together, depending on the width of the cell. Cells that do not contain SP s were
1443 not considered.

1444 We then determined the subset of **frequently visited cells** consisting of all cells that
1445 have at least one SP within them and that were visited by at least 10% of people. For the
1446 sake of reliability, we only calculate the statistics between frequently visited cells and consider
1447 *Confidence* as used by the Market Basket Analysis. **Confidence** indicates the percentage of
1448 trajectories that frequently visit a cell B that also frequently visit cell A . For a Confidence value
1449 above a threshold (set as 60% in our experiments), we can say that a large group of people who
1450 have visited cell A moves together with cell B . Confidence is an estimate of the conditional
1451 probability. Two or more cells for which there is greater than 60% Confidence that have been
1452 visited by a large group of people are called **co-visited cells**.

1453 Then, we check the reliability of the association rules obtained ($A \Rightarrow B$) via *Lift*, which will
1454 confirm that the transition of a user from the SP in A to the SP in B has a positive correlation.

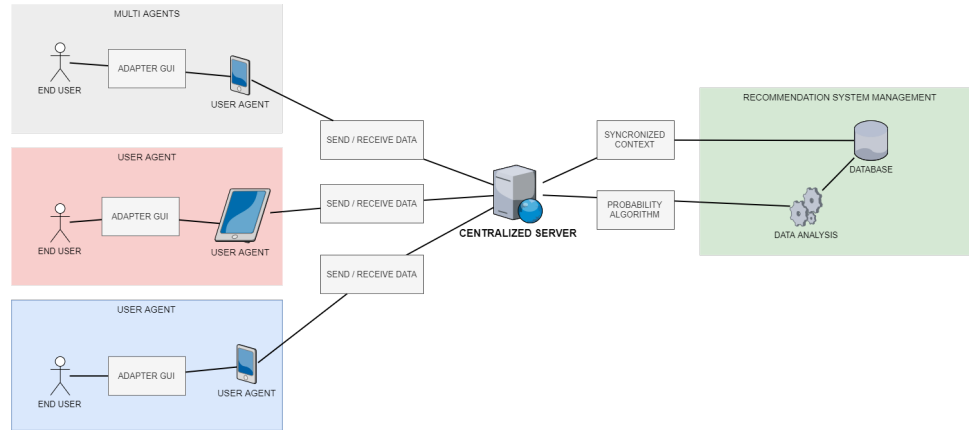


Figure 4.16: System architecture showing the interaction between agents and server: each agent sends his preference for crowded places and where he is, the server gathers data and creates the recommendation system.

B. A Multi-Agent Recommendation System

This section describes the multi-agent system that provides users with recommendations. In general, an agent, according to Wooldridge [97], is simply “a software (or hardware) entity that is in an environment and is able to react autonomously to changes in that environment”. Each agent has the basics to learn and communicate, and in our case learning takes place by acquiring the user’s GPS positions and, communication is achieved by connecting to a centralized server, which alerts all agents when necessary and stores the geographic coordinates of the points visited by users. Figure 4.16 shows the main components of the proposed multi-agent system. An agent runs on a smartphone like an app to get suggestions about the possible destination. The agent offers advice by highlighting any “hot”, which is very busy or A “cold” place, that is, not very crowded, using the collected statistics as described in the previous section. For this, the agent periodically reads the user’s location and checks if a known StayPoint (SP) is nearby. Then, the agent communicates to the server if it is close to an SP . This allows the server to determine the number of people close to an SP , rather than providing their actual location, thus preserving user privacy. In this context, the protection of privacy aims to prevent the disclosure of information relating to the exact location of the user. Figure 4.17 shows the app providing information to the user.

The server, having acquired the closest SP position from the user, returns the list of other SP s that could be visited according to the estimate of the probability of passing through that point (0 equals low probability, 1 equals high probability). In this way, we create a recommendation based on the Collaborative Filtering system [67], as it is based on the choices of other users.

Finally, the user via an administration panel, can set with a flag, if he prefers “hot” or “cold” places. Therefore, the agent, based on the choice of users and the list received from the server,

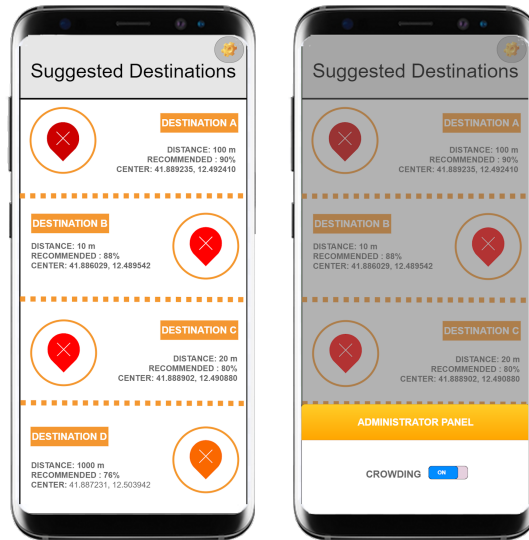


Figure 4.17: User communicates with agents via application GUI. The left panel shows the list of destinations suggested by the multi-agent system and the right panel is the administration view where the user gives his preferences on (un)crowded places. The colour of the icons represents the intensity of crowding, that is, more (less) red equals more (less) crowded.

1479 decides what information to show and then suggests to the user. Destinations are displayed via
 1480 a map or a list of suggested destinations.

1481 All agents are independent of each other and since they extrapolate data directly from the
 1482 device they are reliable, making the architecture stable and trustworthy.

1483
 1484 Figure 4.18 summarizes the methods and the contributions of this thesis.

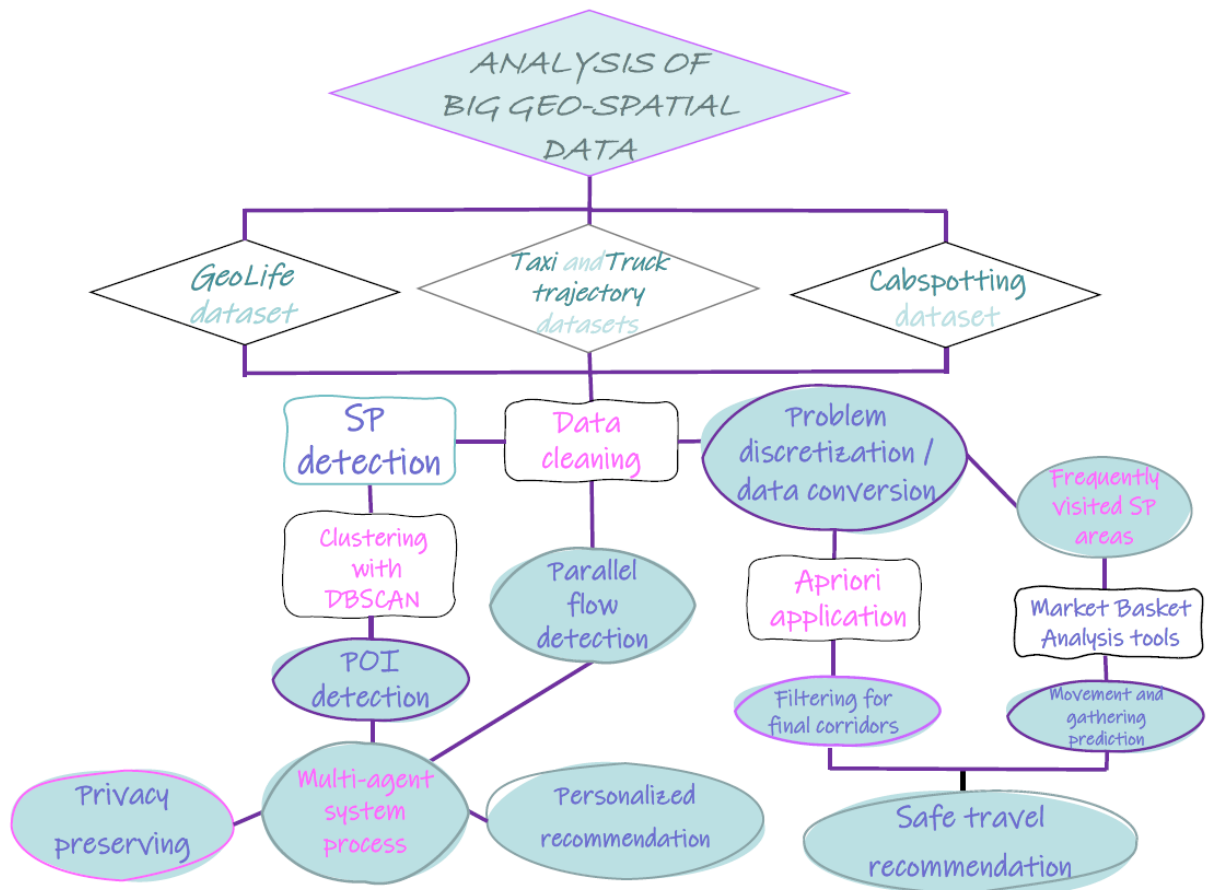


Figure 4.18: The analysis of the datasets is shown through the used methods (rectangular shapes) and the contributions (ovals).

Chapter 5

Evaluation

1485 This chapter describes the experiments carried out using a dataset that collects real movements
1486 from one part of the city to another by taxis and/or people. Positions have been gathered by
1487 periodically reading geo-coordinates from tablets or smartphones. The experiments in Section
1488 5.5 focus on data analysis for determining the probabilities of moving from one StayPoint to
1489 another. This approach is used in our centralized server in order to select the list of suggestions
1490 to send to the agents. The used dataset allows us to simulate the behaviour of a reasonable
1491 number of users, showing the usefulness of the app. Over time, data are updated as provided by
1492 agents. Below we describe the tests carried out and the results that have been found.

5.1 SPs and POIs obtained

1493 The approach presented in this work finds the POIs (starting from a series of user trajectories)
1494 and uses the exchange of information with a centralized server to improve city services and user
1495 knowledge by creating software of content filtering: this creates personalized recommendations
1496 specific to the user to help him in his choices.

1497 The execution time of the SPs detection algorithm (see Algorithm 3 for its pseudo-code) in
1498 GeoLife dataset for 100 trajectories is about 16 minutes. We obtained many SPs for every time
1499 slot, as shown in Table 5.1.

1500 Applying DBSCAN to SPs detected in trajectory of Geolife, we obtained on average 20
1501 clusters for every time slot (120 POIs in total, for a minimum of 9 POIs to a maximum of 29)
1502 that represent significant places for users, i.e. the centroids of these clusters are POIs. E.g.
1503 when we considered time Slot 3 and we set that the minimum number of SPs necessary to make
1504 a cluster as 15 and the ϵ equal to 200 meters, we obtained 29 clusters, hence 29 POIs (Figure
1505 5.1).

1506 The next step of our work was to filter the POIs detected according to popularity. We
1507 considered only POIs with a number of users greater than 10 (called *Popular POIs*), in order to
1508 understand users interaction and similarity. E.g., for time Slot 3 we obtained 9 POIs shared by
1509 a minimum of 11 individuals to a maximum of 80 individuals (see Figure 5.2).

1510 Our experiments in GeoLife dataset have shown that in different time slots a set of different

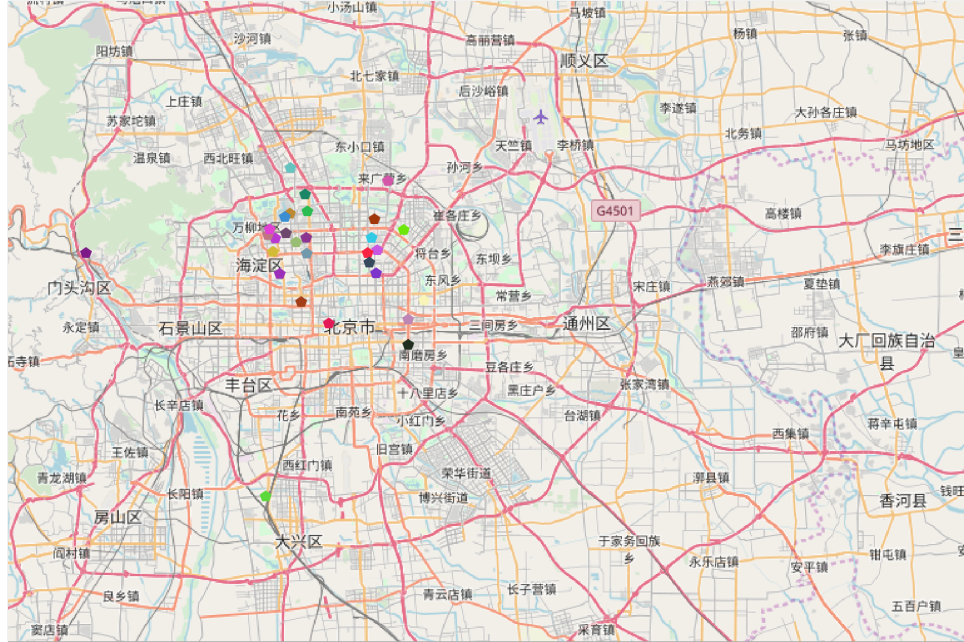


Figure 5.1: POIs obtained for the trajectories on time Slot 3.

Table 5.1: Results about SPs and POIs obtained for GeoLife’s trajectories.

Time Slot	Total number of				DBSCAN	
	<i>SPs</i>	<i>Users</i>	<i>POIs</i>	Popular <i>POIs</i>	<i>Eps (km)</i>	<i>MinPts</i>
1	2966	122	18	8	0.3	15
2	3772	124	27	8	0.25	15
3	4146	145	29	9	0.2	15
4	1899	130	24	6	0.2	15
5	751	84	13	4	0.4	10
6	545	84	9	1	0.4	10

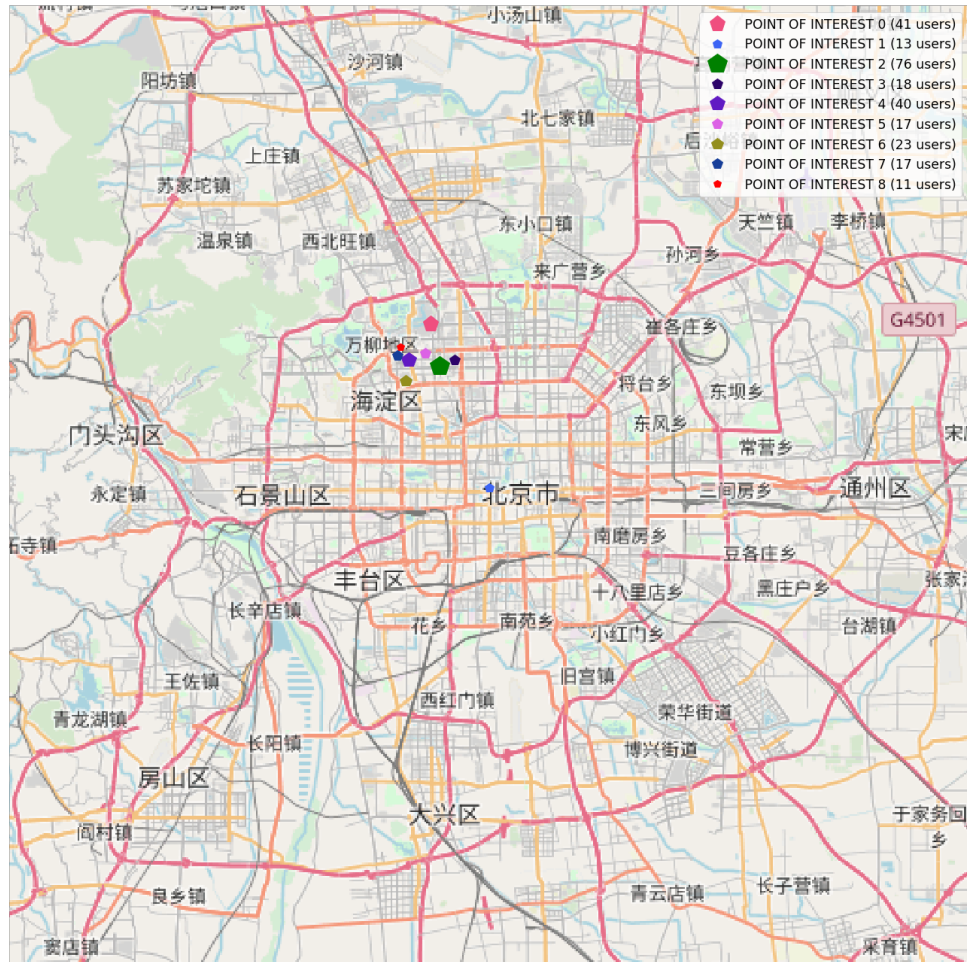


Figure 5.2: Popular POIs in time slot 3 on the map, Geolife dataset.

1511 individuals move together to the same POIs, like parks, departments of Universities, shopping
1512 centres, hostels, parking spaces, libraries, stadiums, banks, Metro and bus stops. This suggests
1513 us a similarity between users.

1514 The total of Popular POIs on all 6 time slots was 36; looking for the most distant pairs
1515 of points they have 8 km of longitude difference, 25 km of latitude difference on this area
1516 [39.908309, 116.262296] x [116.368098, 40.128495]. These places of interest in question started
1517 from the north in Yangyang Paradise (amusement park), up to the Cultural Palace of Nation-
1518 alities in Fuxingmen Inner Street in south, crossing Changping District, Haidian District with
1519 Tsinghua Park, Beijing Shi and Zhongguancun.

1520 For detected POIs we can further say that our experiments show a correlation of people
1521 moving from one POI to another: users remain in these areas in certain common time slots. In
1522 our experiments the execution time of DBSCAN on the 6 time slots for GeoLife ranges from
1523 a minimum of 240 ms to a maximum of 1.44 s. Our implementation uses Python 3, and the
1524 experiments were run in a host having an Intel Xeon CPU E5-2620 v3 2.40GHz, with RAM 32
1525 GB.

1526 Figure 5.3 shows a list of POIs close to the user. Thus, for each POI the user can access the
1527 scores collected from other people's comments, as well as his comments. As mentioned in STEP
1528 1, the POIs are detected by the implemented algorithm discussed in the next section. We can
1529 see four nearby points labeled as POI:

- 1530 • Chaofan Weiye Kejiao Library with coordinates: 39.98405510061326, 116.3204636235443;
- 1531 • Haidian Stadium with coordinates: 39.987213527969644, 116.30248430595732;
- 1532 • Beijing Zhongguancun Branch Commercial Rural Bank with coordinates: 39.980016801082485,
1533 116.30856309688643;
- 1534 • Beihang University with coordinates: 39.98011363182701, 116.34218061609567.

1535 Another nearby POI has been associated with a parking area (however, it is not listed in Figure
1536 5.3:

- 1537 • Parking of the satellite building with coordinates: 39.97673497237701, 116.33137904408086.

1538 To validate the results of our algorithm in finding POIs, each discovered site has been verified
1539 on Google Maps. So, the list above consists of real sites, which are POIs according to Google
1540 Maps, which are located within a radius of 100 meters from the POIs detected by our algorithm.

1541 In *Taxi trajectory data*, for the SP detection algorithm we chose a distance threshold of 200
1542 meters and a time threshold of 5 minutes.

1543 The obtained SPs were 31,621, with an execution time of 51 minutes and 9 seconds. An
1544 example of SPs related to trajectories of a taxi is shown in Figure 5.4.

1545 Relatively to the POI detection in this dataset, the parameters set in the DBSCAN were ϵ
1546 equal to 200 meters and minimum points (SPs) equal to 8. This algorithm produced 560 clusters
1547 whose centroids are the POIs, with an execution time of 10.2 seconds. Thirdly, resulting POIs
1548 were filtered to find the popular POIs, i.e. the POIs shared by at least 8 taxis. The total

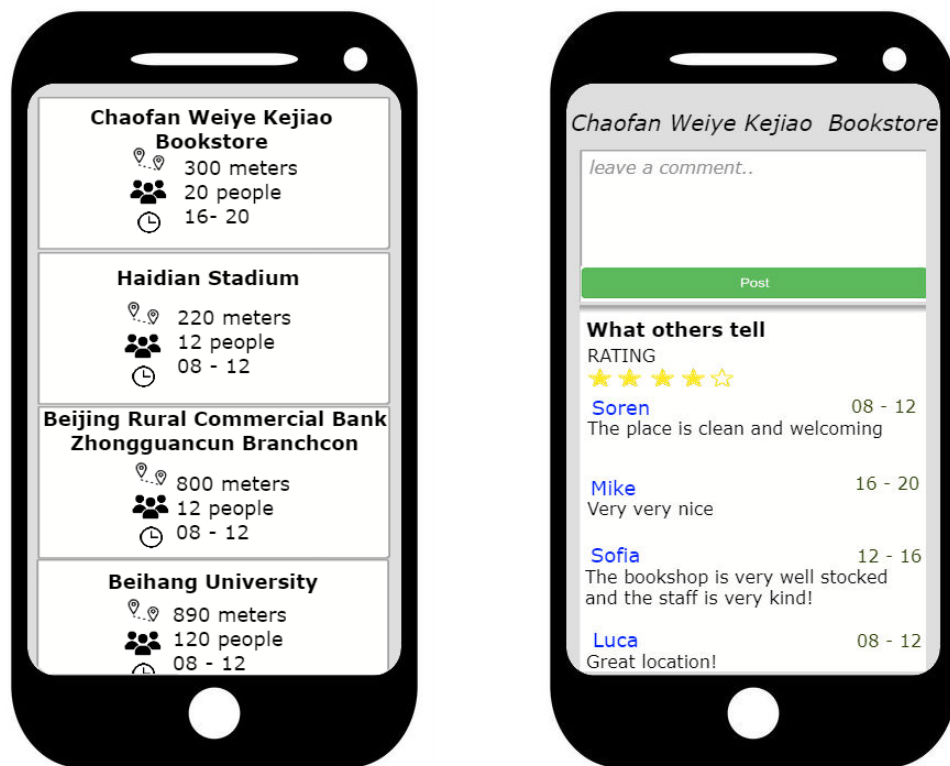


Figure 5.3: A list of POIs and associated dynamic data presented by an agent.

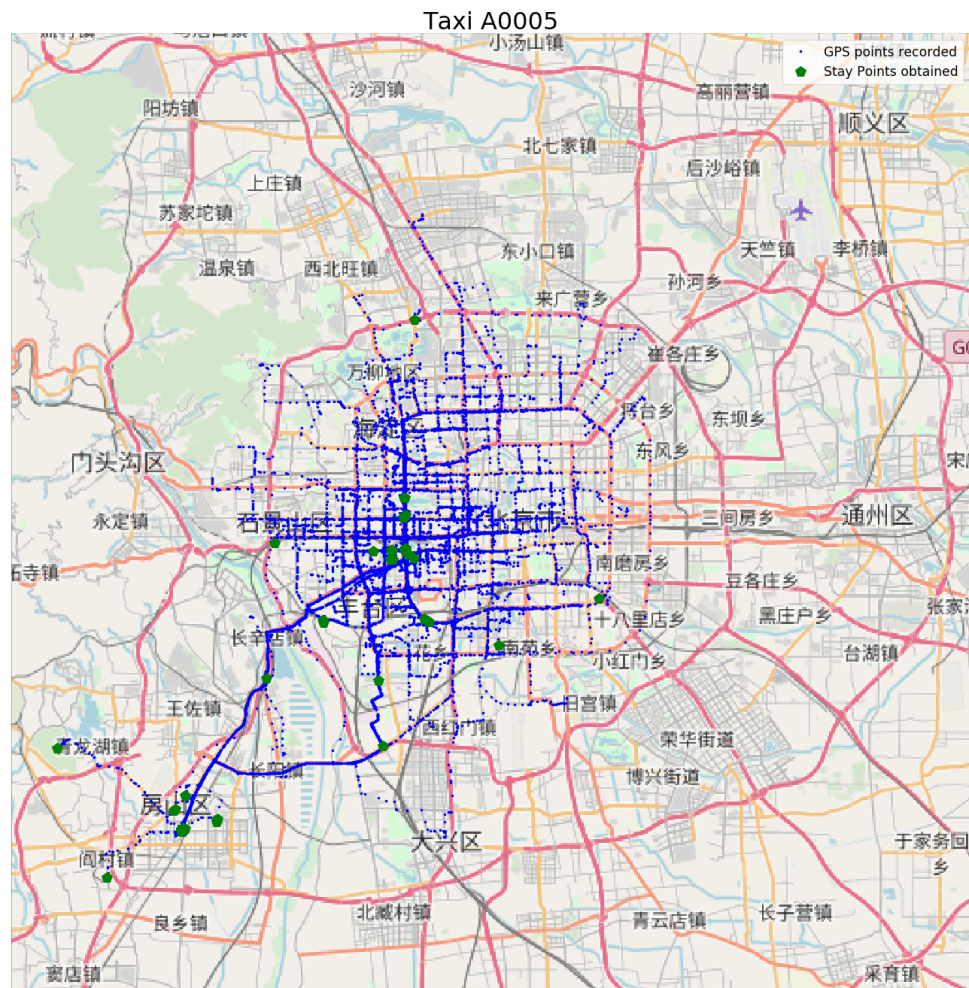


Figure 5.4: StayPoints of a taxi on its trajectories.

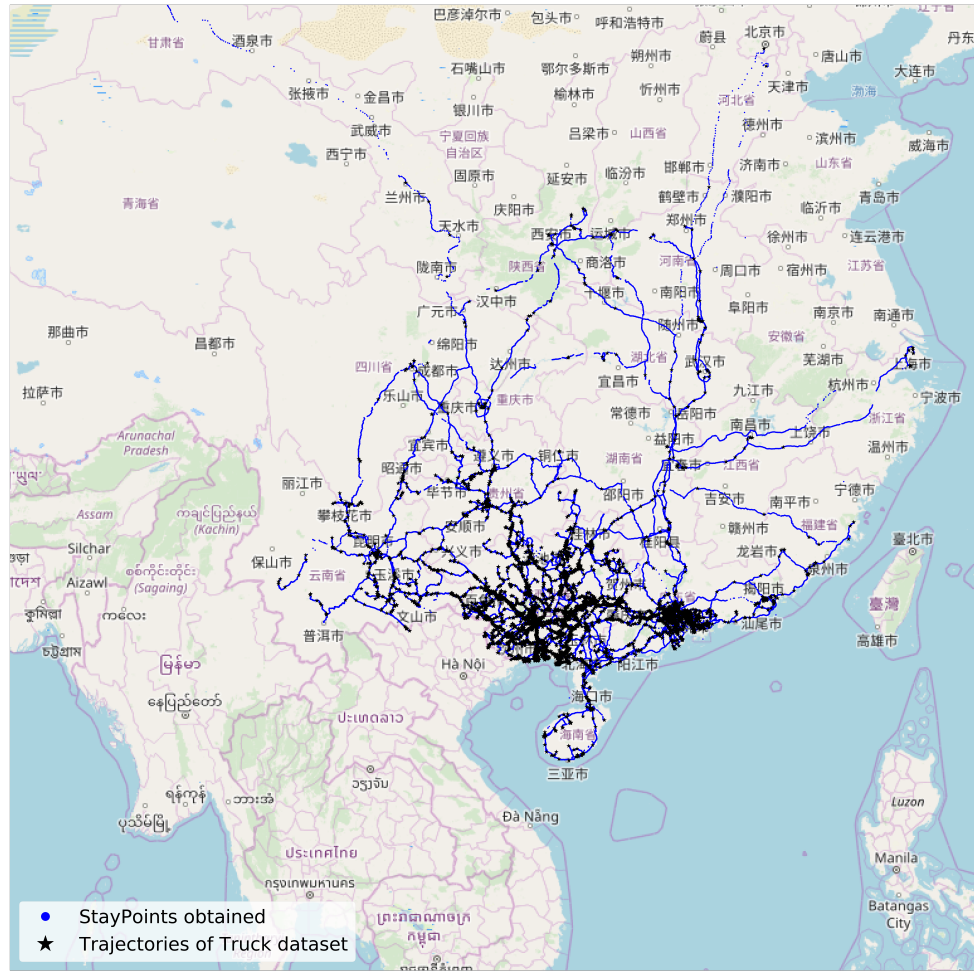


Figure 5.5: StayPoints of trucks on their trajectories.

1549 number of such popular POIs was 257. Out of a total of 101 taxis, the popular POIs obtained
 1550 were visited by 8 to 95 taxis.

1551 The identified popular POIs can be found in the area $[116.093765, 39.79077] \times$
 1552 $[116.608337, 40.092962]$. They are in Beijing and the most far apart pairs of points have
 1553 distances 45 *km* for the longitude and on 35 *km* for the latitude.

1554 Looking for SPs in *Truck* dataset the parameters $\text{DistThr} = 200$ meters and $\text{TimeThr} = 10$
 1555 minutes were chosen. A total of 54,962 SPs were identified, in a time of 5 hours, 1 minute and
 1556 5 seconds.

1557 In Figure 5.5 it is possible to see the trajectories recorded by trucks and their SPs, they
 1558 pass through these provinces: Shānxī Shěng, Shānxī Shěng, Gansu, Henan, Hubei, Hebei, Bei-
 1559 jing Municipality, Hunan, Sichuan, Guizhou, Yunnan, Guangxi Zhuang Autonomous Region,
 1560 Guangdong, Jiangxi, Anhui, Fujian, Zhejiang and Shanghai Municipality.



Figure 5.6: Popular POIs in Truck dataset.

1561

The choice of different values of time threshold set in the SP detection is due to the different nature of the three datasets. For taxis, a reasonable time of stay is 5 minutes, for trucks 10 minutes if it has to consider the bays, for Geolife dataset, which includes routes of users on foot, DistThr was chosen equal to 20 minutes. These parameters were validated by the average speed value relative to the flows found, in the vicinity of the Popular POIs.

1562

1563

1564

1565

1566

For POIs detection the spatial threshold in DBSCAN remained unchanged ($\epsilon = 200\text{ m}$) and the MinPts = 8 as in the previous case. The clustering execution time was 5.31 seconds with 1,065 POIs detected.

1567

1568

1569

According to the minimum number of taxis (8 out of 101), the popular POIs obtained were 127: they are shared by a minimum number of 8 trucks and a maximum number of 24 trucks. They cross the counties: Shangsi, Longzhou, Tiandong, Long'an, Mashan, Pingnan, Teng, Yunan, the prefecture cities: Chongzuo, Wúzh-ou, Yunfu, the districts Jinchengjiang and Yun'an, the Luoding city and the Kunming Subdistrict (see Figure 5.6).

1570

1571

1572

1573

1574

Popular POIs obtained covered the area $[106.880591, 21.609655] \times [113.670971, 24.685619]$, with 700 km of difference in longitude and 350 km of latitude difference between the two most distant pairs of points. From west to east touched the cities: Chongzuo, Nanning, Guigang, Qinzhou, Fangchenggang, Beihai, Yulin, Wuzhou, Zhaoqing, Foshan, Canton and Dongguan.

1575

1576

1577

5.2 Movement of people detected

1578

The analysis of Geolife dataset has found 25 flows, shared by a minimum of 20 users to a maximum of 53 users and with a minimum distance equals to 150 meters between the trajectories. The length of paths that covered these flows is between 2 and 5 kilometres. More details are visible in Table 5.2. One flow for 21 walking users is shown in Figure 5.7.

1579

1580

1581

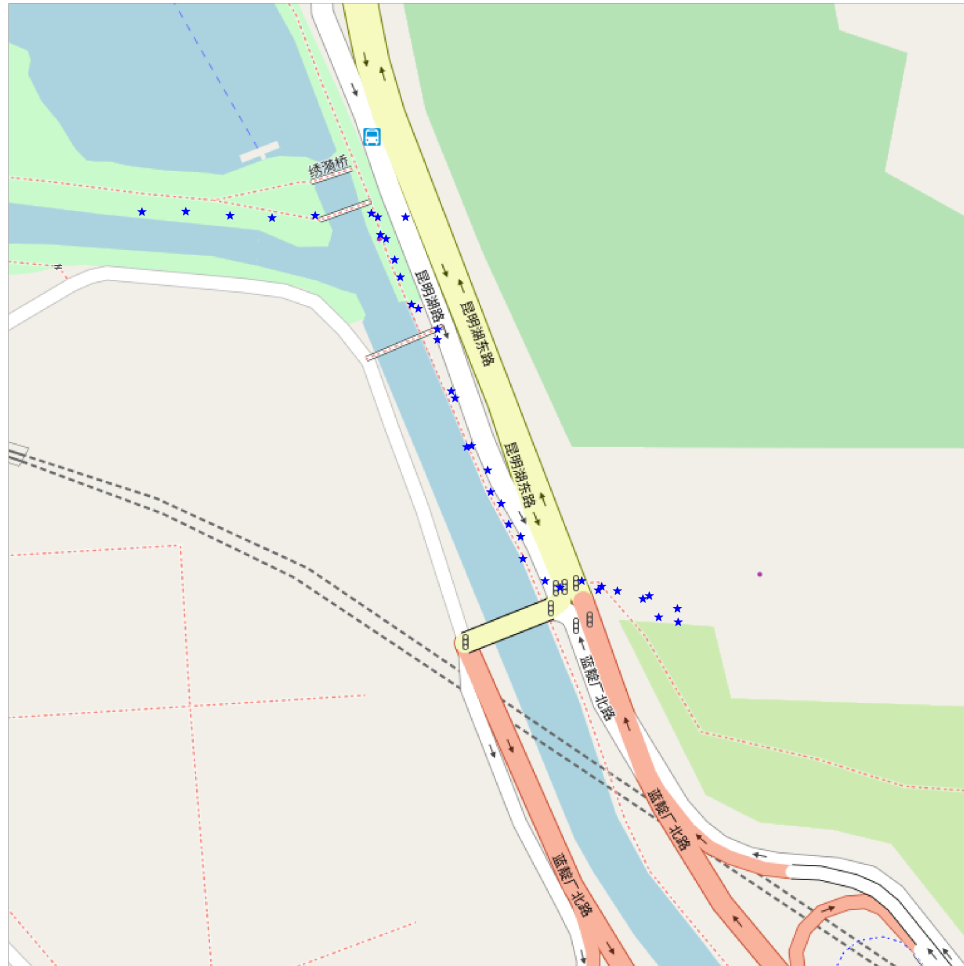


Figure 5.7: A flow for walking users, detected in GeoLife dataset.

Table 5.2: Details about flows obtained in GeoLife dataset.

Starting point	End point	Distance (km)	Total number of users
(39.932367, 116.385852)	(39.931153, 116.36233)	2.010	20
(39.881835, 116.269233)	(39.910339, 116.268027)	3.171	20
(39.988628, 116.483763)	(39.972852, 116.466223)	2.304	29
(39.986877, 116.399344)	(39.987707, 116.429464)	2.567	44
(39.90712, 116.350367)	(39.930045, 116.349922)	2.549	29
(39.979295, 116.284629)	(39.93731, 116.267806)	4.883	21
(39.971537, 116.274198)	(39.950695, 116.269099)	2.357	27
(39.93856, 116.268389)	(39.962849, 116.268952)	2.701	24
(39.931233, 116.321381)	(39.96773, 116.314436)	4.101	23
(39.972308, 116.327603)	(39.951754, 116.329522)	2.291	29
(39.984809, 116.333326)	(39.986664, 116.392776)	5.069	41
(39.986462, 116.375304)	(39.985521, 116.340627)	2.956	53
(39.976682, 116.453254)	(39.987662, 116.419631)	3.114	33
(39.985884, 116.361805)	(39.986612, 116.388662)	2.289	50
(39.984303, 116.326575)	(39.985593, 116.356151)	2.523	53
(39.986183, 116.433872)	(39.968712, 116.464021)	3.220	25
(39.948103, 116.345853)	(39.974371, 116.333989)	3.090	28
(39.943232, 116.34601)	(39.961544, 116.339601)	2.108	37
(39.98915, 116.331707)	(39.96691, 116.338599)	2.541	20
(40.023843, 116.347193)	(40.040263, 116.333505)	2.166	38
(40.032671, 116.313707)	(40.00121, 116.331161)	3.800	35
(39.931233, 116.321381)	(39.949994, 116.317978)	2.106	22
(39.944018, 116.318824)	(39.96041, 116.316824)	1.830	41
(39.969003, 116.313939)	(39.988966, 116.3099)	2.246	35
(39.991503, 116.333198)	(39.990766, 116.310343)	1.948	35

1582 In Taxi dataset, thanks data preprocessing, about 400 outliers were removed. A minimum
1583 distance equals to 150 meters has been chosen for the algorithm and 77 flows were found with
1584 a minimum length equals to 1 *km* and a minimum number of taxis equal to 9. More details in
1585 Table 5.3.

1586 Then, an estimate was made of the useful time slot to suggest to potentially interested agents.
1587 According to available data, it has been found that the preferred time slot was between 00:00:00
1588 and 07:59:59.

1589 The algorithm takes a trajectory as a reference to find flows. In this case, according to data,
1590 checks to find flows were performed taking about 10% of the trajectories (randomly chosen)
1591 and all the results were assembled together. Table 5.5 shows the details for all the datasets
1592 considered, and the first row of data is about taxis. In the table, each row contains: the name of
1593 the dataset used (Dataset); the time slot most chosen by users (Time Slot); the number of flows
1594 found (Number of flows); the maximum density detected among flows (Maximum Density); the
1595 maximum distance in meters of the longest flow found (Maximum Distance).

1596 Subsequently, the algorithm was started in Truck dataset by taking 10 trajectories at random
1597 and using them as a reference for a comparison between all other trajectories. With data
1598 preprocessing, data removed were approximately between 0.001% and 0.007%, by eliminating
1599 outliers. In this case, a minimum distance between two points was set equal to 250 meters for
1600 elimination of noise. This choice was made due to the high number of points present in each
1601 trajectory and almost allows us to halve the response times of the algorithm.

1602 Therefore, it has obtained about 114 flows with a minimum distance of 1 km and with a
1603 minimum number of users equal to 8%. For such experiments, it has been checked that the
1604 preferred time slot is between 16:00 and 19:59, though each POI has been associated with its
1605 own time slot. Table 5.4 gives the details of the experiments.

5.3 Corrispondence between flows and POIs

1606 With reference to the performed analyses using the Taxi dataset, a 68% correspondence was
1607 found between the popular POIs and the shared flows, as can be seen in Figure 5.8 that illus-
1608 trates both flows and POIs (black diamonds). Such places correspond to real useful sites and
1609 attractions, like: World park, Grand View Garden, Muzhi Gongdian amusement park, Ance-
1610 tral Temple, Dongcheng Chongwen Science & Technology Museum, financial buildings, airport
1611 area, Long-distance Passenger Transport Terminal, Chaoyangmen SOHO shopping center, UK
1612 Embassy and Silver Bridge on Shichaha river.

1613 The density of flows detected in this test was between a minimum of 9 taxis to a maximum
1614 of 70 taxis (among a total of 101 taxis). The intersections between the identified popular flows
1615 and POIs are found in the Wanliu, Haidian, Fengtai Districts, in the Tianzhu, Cuige, Lai Guang
1616 Yin, Shibalidianxiang Villages, in the Dongfeng area and in Nanmofang residential district.

1617 Experiments on the Truck dataset have shown that there are 56 occurrences of Popular
1618 POIs near the 114 identified flows. It indicates a meaningful 50% correspondence, which gives a
1619 proximity within 100 *m* on the latitude and longitude of the points with respect to the identified
1620 flows. Points and flows were calculated on a minimum number of 8 trucks sharing them.

Table 5.3: Details about flows obtained in Taxi dataset.

<i>Starting point</i>	<i>End point</i>	<i>Distance (km)</i>	<i>Total number of users</i>
(40.01391983, 116.4700165)	(40.01517487, 116.4484177)	1.302	12
(40.05770493, 116.5829163)	(40.07411575, 116.5801697)	1.654	48
(39.94765091, 116.3908005)	(39.94740677, 116.3716965)	1.510	10
(39.96813583, 116.4283142)	(39.96772766, 116.4118347)	1.199	10
(39.88315964, 116.4419022)	(39.88715744, 116.4552155)	1.028	15
(39.94997406, 116.365387)	(39.93468475, 116.3665466)	1.143	14
(39.93462753, 116.4555817)	(39.94881439, 116.455368)	1.262	23
(39.87185669, 116.3431015)	(39.8875885, 116.3426819)	1.604	11
(40.02097321, 116.4004517)	(40.02082825, 116.4183197)	1.505	9
(39.94997406, 116.365387)	(39.9397049, 116.3663635)	1.081	21
(39.90634155, 116.330368)	(39.89759445, 116.3229828)	1.051	10
(39.98168945, 116.287735)	(39.98122025, 116.3092651)	1.642	10
(39.91970825, 116.3504028)	(39.90489197, 116.3504639)	1.350	28
(39.95053864, 116.4556808)	(39.96003342, 116.4426346)	1.283	9
(39.85622406, 116.3926163)	(39.85575867, 116.4096146)	1.404	9
(39.94730377, 116.3659515)	(39.94137955, 116.3498306)	1.520	23
(39.89347839, 116.4391632)	(39.90650558, 116.4296646)	1.618	18
(39.91085434, 116.4288025)	(39.92374802, 116.4281006)	1.127	28
(39.97738647, 116.28302)	(39.98391342, 116.2976837)	1.343	14
(39.86820984, 116.2659988)	(39.87461853, 116.2800522)	1.020	18
(39.89577484, 116.2631836)	(39.89577484, 116.2804489)	1.209	11
(39.90711975, 116.462883)	(39.90679169, 116.4413147)	1.559	12
(39.91335297, 116.3504181)	(39.89775085, 116.3500671)	1.500	21
(39.97241974, 116.295433)	(39.96090317, 116.301918)	1.213	15
(39.9569931, 116.3484344)	(39.9441452, 116.3498993)	1.348	14
(39.88085175, 116.4550858)	(39.89391327, 116.4554977)	1.398	29
(39.8991394, 116.455452)	(39.91464996, 116.4552536)	1.624	17
(39.90274048, 116.4475021)	(39.89246368, 116.4366989)	1.468	16
(39.87311172, 116.3432999)	(39.88673782, 116.3427353)	1.436	20
(39.89627075, 116.3370667)	(39.90134811, 116.3506317)	1.233	12
(39.9216423, 116.3496017)	(39.90307236, 116.350502)	1.781	15
(39.94994736, 116.3995514)	(39.93841171, 116.402298)	1.069	10
(39.98218918, 116.2887192)	(39.9839325, 116.3079834)	1.452	17
(39.91148758, 116.4289856)	(39.92631912, 116.4281158)	1.379	30
(40.02723312, 116.3072968)	(40.0164032, 116.3104858)	1.024	12
(40.0593605, 116.5826035)	(40.07794189, 116.5806198)	1.992	23
(39.93490982, 116.3380661)	(39.94364929, 116.353653)	1.424	12
(39.91411209, 116.3503036)	(39.89731216, 116.3494492)	1.589	24

<i>Starting point</i>	<i>End point</i>	<i>Distance (km)</i>	<i>Total number of users</i>
(39.9209938, 116.3500671)	(39.90653992, 116.3503342)	1.323	70
(39.92617798, 116.4555359)	(39.9408226, 116.455513)	1.502	15
(39.90365601, 116.4556198)	(39.91807556, 116.4553299)	1.286	57
(40.06372833, 116.5824814)	(40.07794189, 116.5806198)	1.580	51
(39.94777298, 116.3972855)	(39.94748688, 116.3790817)	1.269	21
(39.89669418, 116.344902)	(39.91262817, 116.3507309)	1.553	10
(39.85474014, 116.3630981)	(39.85559464, 116.3847809)	1.848	12
(39.91291809, 116.4304504)	(39.92630386, 116.4309845)	1.214	25
(39.90378571, 116.455513)	(39.91992569, 116.4580688)	1.548	17
(39.94758224, 116.3987656)	(39.94828033, 116.4191971)	1.569	22
(39.92295456, 116.426033)	(39.92560577, 116.4419327)	1.049	11
(39.85756302, 116.4300995)	(39.85603333, 116.4152679)	1.194	12
(39.93951797, 116.4056015)	(39.93938828, 116.3879166)	1.379	11
(40.05770493, 116.5829163)	(40.07120132, 116.5812683)	1.461	53
(39.9569397, 116.4121323)	(39.94761276, 116.4052505)	1.002	9
(39.96665955, 116.3703537)	(39.96731186, 116.3871002)	1.040	14
(39.9621048, 116.2922821)	(39.97661209, 116.2945175)	1.099	9
(39.94763565, 116.3919983)	(39.94746399, 116.372818)	1.219	22
(40.00532532, 116.2823639)	(40.01837158, 116.2777863)	1.495	11
(39.92563629, 116.3669815)	(39.90939331, 116.3677673)	1.398	12
(39.91980743, 116.3500671)	(39.90393448, 116.3505173)	1.403	50
(39.96784592, 116.4162979)	(39.96609116, 116.4334335)	1.240	22
(39.98101807, 116.3206863)	(39.96736145, 116.319397)	1.211	13
(39.90114594, 116.4554825)	(39.91568375, 116.4553986)	1.207	42
(39.89970398, 116.4366837)	(39.90280914, 116.4532013)	1.418	12
(39.88513565, 116.4499969)	(39.89596939, 116.4547653)	1.083	27
(39.85876846, 116.4561691)	(39.85668564, 116.454567)	1.003	14
(39.97903824, 116.2846527)	(39.98388672, 116.3008804)	1.266	28
(40.01387787, 116.4676666)	(40.00965118, 116.4824829)	1.218	14
(39.96792221, 116.4197693)	(39.96332169, 116.4377518)	1.365	15
(39.96615219, 116.3343353)	(39.96644974, 116.3530502)	1.232	19
(39.94676208, 116.363266)	(39.93883514, 116.3483658)	1.227	43
(39.88864517, 116.4552002)	(39.90485001, 116.4554672)	1.473	15
(39.94742584, 116.383316)	(39.94783401, 116.4070129)	1.928	9
(39.8993988, 116.3507004)	(39.91819, 116.3501816)	1.804	15
(39.86477661, 116.4526367)	(39.87872314, 116.4550018)	1.425	23
(39.88868332, 116.4384689)	(39.90473557, 116.4297638)	1.598	20
(39.95137024, 116.3627014)	(39.9397049, 116.3663635)	1.125	22
(39.91291809, 116.4304504)	(39.92887497, 116.4280472)	1.497	31

Table 5.4: Details about flows obtained in Truck dataset.

<i>Starting point</i>	<i>End point</i>	<i>Distance (m)</i>	<i>Total number of users</i>
(23.46398, 110.1776)	(23.45764, 110.1636)	1329.87	9
(23.19097, 109.7451)	(23.18311, 109.7294)	1490.506	13
(23.38019, 110.0936)	(23.39321, 110.1007)	1230.484	11
(23.32668, 110.0111)	(23.3253, 109.9933)	1492.372	9
(23.36797, 110.0626)	(23.37064, 110.0579)	1213.874	11
(23.35973, 110.0617)	(23.37147, 110.0706)	1570.259	9
(23.53082, 110.399)	(23.51622, 110.3932)	1293.346	9
(22.16008, 108.1142)	(22.15369, 108.1282)	1282.383	12
(22.15512, 107.9694)	(22.16267, 107.9869)	1696.31	9
(23.1298, 109.3957)	(23.13622, 109.4109)	1419.325	12
(22.73798, 109.3333)	(22.72276, 109.3332)	1476.157	10
(22.8669, 108.3745)	(22.87254, 108.3894)	1503.326	28
(23.11833, 109.5236)	(23.11325, 109.5407)	1451.289	17
(23.15197, 109.4091)	(23.13818, 109.4114)	1469.86	9
(23.1204, 109.5085)	(23.11804, 109.5263)	1693.694	11
(22.87388, 108.3946)	(22.86811, 108.3807)	1226.46	20
(23.11759, 109.5282)	(23.11086, 109.5443)	1722.323	18
(23.13085, 109.3981)	(23.13505, 109.4145)	1402.94	12
(23.11982, 109.5211)	(23.1212, 109.5031)	1465.044	17
(23.12318, 109.483)	(23.12212, 109.4974)	1204.83	13
(22.83292, 108.3841)	(22.82915, 108.3879)	1118.652	9
(22.8275, 108.3957)	(22.82808, 108.4114)	1094.745	9
(22.90489, 108.2954)	(22.91195, 108.2957)	1640.532	11
(22.66003, 108.3862)	(22.67287, 108.3863)	1163.431	18
(22.88428, 108.2993)	(22.87448, 108.3139)	1459.984	16
(22.88168, 108.3084)	(22.88609, 108.2946)	1456.035	16
(22.73972, 108.3798)	(22.75565, 108.3777)	1390.175	11
(22.76152, 108.3703)	(22.75287, 108.3796)	1205.48	15
(22.8559, 108.3662)	(22.86585, 108.3544)	1216.196	23
(22.83669, 108.3734)	(22.85397, 108.3667)	1785.581	10
(22.84362, 108.3709)	(22.85662, 108.3656)	1466.063	12
(23.43506, 110.3099)	(23.44214, 110.2959)	1625.982	9
(23.39799, 111.1887)	(23.39769, 111.1798)	1165.062	11

<i>Starting point</i>	<i>End point</i>	<i>Distance (m)</i>	<i>Total number of users</i>
(23.53223, 110.3998)	(23.51722, 110.3931)	1293.346	9
(23.37135, 110.0665)	(23.37696, 110.0852)	1915.182	11
(22.83352, 108.3868)	(22.8329, 108.3744)	1200.724	9
(22.93619, 108.2913)	(22.92397, 108.2971)	1161.448	10
(22.96408, 108.2891)	(22.97635, 108.2969)	1308.752	10
(22.71571, 108.3417)	(22.72903, 108.3304)	1392.192	20
(22.7945, 108.2684)	(22.80962, 108.2667)	1364.938	12
(22.86854, 108.3314)	(22.88037, 108.3378)	1162.064	15
(22.88345, 108.3021)	(22.88009, 108.3198)	1693.3	15
(22.88591, 108.2952)	(22.88088, 108.3118)	1525.455	19
(22.87914, 108.3368)	(22.86574, 108.332)	1553.411	12
(23.14492, 108.2352)	(23.15158, 108.247)	1412.635	9
(22.85945, 108.3627)	(22.86626, 108.3734)	1057.851	29
(22.86559, 108.3569)	(22.8696, 108.3838)	2492.366	10
(22.92994, 108.5176)	(22.93651, 108.5323)	1501.728	9
(22.92137, 108.4967)	(22.91548, 108.4802)	1500.757	10
(22.87822, 108.2934)	(22.89584, 108.2946)	1624.175	29
(22.88225, 108.3057)	(22.89031, 108.2942)	1340.602	13
(21.75837, 108.6056)	(21.75968, 108.5937)	1261.628	9
(22.14857, 108.5792)	(22.14441, 108.5931)	1107.592	9
(22.11363, 108.6227)	(22.10348, 108.6327)	1152.363	14
(22.3357, 108.399)	(22.32056, 108.4011)	1450.596	18
(22.08649, 108.6327)	(22.1009, 108.6329)	1411.097	14
(22.2349, 108.4161)	(22.25158, 108.4173)	1664.687	13
(22.09361, 108.6313)	(22.10851, 108.6261)	1319.371	11
(22.59992, 108.385)	(22.61542, 108.3855)	1432.261	13
(22.24552, 108.417)	(22.26123, 108.4164)	1646.367	12
(22.61952, 108.3845)	(22.63658, 108.3863)	1495.875	18
(22.28093, 108.4189)	(22.29208, 108.4063)	1541.909	19
(22.26574, 108.4183)	(22.28073, 108.4192)	1603.248	15
(22.27401, 108.4218)	(22.28527, 108.4135)	1399.102	17
(22.36413, 108.3838)	(22.37866, 108.3773)	1450.092	17
(22.09086, 108.6327)	(22.10606, 108.6308)	1659.565	10
(22.39749, 108.3841)	(22.38282, 108.3774)	1430.2	19
(22.74021, 108.3797)	(22.75572, 108.3777)	1390.175	12
(22.78894, 108.2468)	(22.7846, 108.261)	1135.41	9
(22.0058, 108.6383)	(22.00215, 108.6533)	1454.266	15
(21.96763, 108.5858)	(21.97193, 108.5985)	1213.778	9
(21.74269, 108.6015)	(21.73132, 108.5916)	1119.86	12
(22.07253, 108.6324)	(22.05612, 108.6317)	1414.589	14
(22.00949, 108.6315)	(22.00352, 108.6457)	1364.895	24
(22.06896, 108.6328)	(22.05241, 108.6305)	1513.593	15
(21.96581, 108.6823)	(21.95049, 108.6838)	1439.659	17

<i>Starting point</i>	<i>End point</i>	<i>Distance (m)</i>	<i>Total number of users</i>
(21.95528, 108.6852)	(21.97146, 108.6771)	1539.229	15
(22.02115, 108.6338)	(22.00407, 108.6414)	1850.013	13
(22.0795, 108.6326)	(22.06653, 108.633)	1198.458	9
(22.047, 108.6287)	(22.0311, 108.631)	1426.291	19
(23.0508, 112.6255)	(23.04628, 112.6094)	1494.967	9
(23.05515, 112.6359)	(23.04925, 112.6203)	1348.628	14
(22.7636, 108.2572)	(22.76369, 108.2725)	1553.032	26
(22.85393, 108.2529)	(22.84287, 108.248)	1278.944	16
(22.88176, 108.3012)	(22.87944, 108.3169)	1312.716	17
(22.88609, 108.3395)	(22.87303, 108.331)	1235.096	10
(22.86431, 108.3297)	(22.87877, 108.335)	1442.315	15
(22.89781, 108.2826)	(22.88789, 108.2942)	1539.527	10
(22.91355, 108.4698)	(22.90682, 108.4566)	1323.153	10
(21.60608, 108.3474)	(21.61972, 108.3512)	1212.667	14
(22.86503, 108.369)	(22.87001, 108.3843)	1459.871	31
(22.8648, 108.3592)	(22.86687, 108.376)	1540.272	22
(21.60355, 108.3447)	(21.6161, 108.3486)	1368.788	14
(22.88035, 108.416)	(22.87374, 108.4025)	1150.239	12
(22.8604, 108.3175)	(22.86467, 108.3284)	1170.169	17
(21.76312, 108.3744)	(21.74969, 108.3731)	1387.631	10
(21.6694, 108.3656)	(21.67165, 108.362)	1244.098	12
(21.61372, 108.3506)	(21.63112, 108.3576)	1966.214	10
(22.87388, 108.3946)	(22.8696, 108.3839)	1378.008	12
(22.96362, 108.2673)	(22.95115, 108.2609)	1199.331	9
(23.62062, 107.0523)	(23.6268, 107.0366)	1291.845	9
(23.61226, 107.0659)	(23.61055, 107.0808)	1192.61	10
(22.92252, 107.9986)	(22.92764, 107.9814)	1493.632	12
(22.8472, 108.1149)	(22.84776, 108.1292)	1321.734	16
(22.73798, 109.3333)	(22.72276, 109.3332)	1218.709	10
(22.80511, 108.2687)	(22.81823, 108.2577)	1419.84	9
(22.82901, 108.2062)	(22.83938, 108.2102)	1100.438	11
(22.84139, 108.1561)	(22.83088, 108.1669)	1258.154	18
(22.76363, 108.2608)	(22.7749, 108.2531)	1307.189	12
(22.76152, 108.3703)	(22.7492, 108.3799)	1385.364	15
(22.83818, 108.161)	(22.82777, 108.1718)	1600.854	14
(22.76432, 108.2508)	(22.76374, 108.2681)	1548.94	27
(22.7491, 108.2928)	(22.73492, 108.29)	1183.701	12
(22.82743, 108.1799)	(22.82482, 108.1947)	1435.011	16

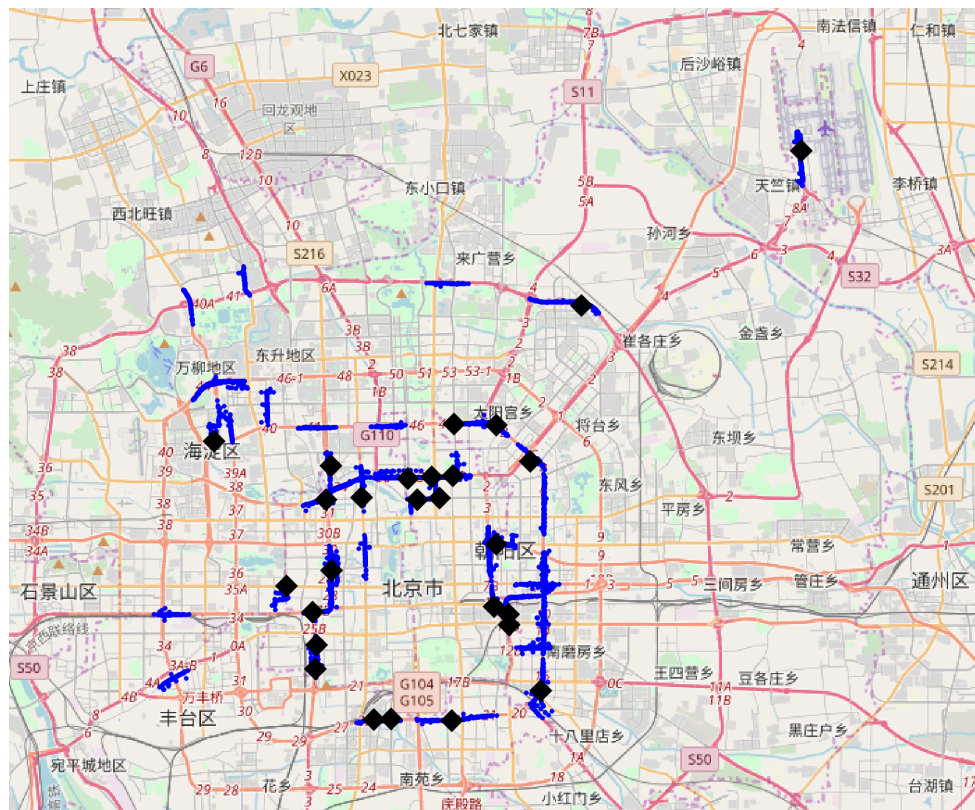


Figure 5.8: POPULAR POIs in correspondence of Taxi dataset's flows.

Table 5.5: An overview of the tests performed about flow found.

Dataset	Time Slot	Number of flows	Maximum Density	Maximum Distance
Taxi	00:00-07:59	77	70	1716
Geolife	08:00-11:59	25	53	3052
Truck	16:00 - 19:59	114	31	1718

Popular places of interest found near the flows correspond to: service areas, parking lots, grocery stores, tollgates, ATMs, auto parts stores, lottery ticket dealer, Driver Examination Center, hotel, restaurants, repair service area, gas station, lubricating oil sale service and Port of Fangcheng. They belong to Nanning, the capital of the Guangxi Zhuang Region, to the prefecture cities: Fangchenggang, Qinzhou, Guigang, Zhaoqing, to the Guiping county-level city and to the district Wuming.

As for the Geolife dataset, on 25 flows found, there were 13 matches with popular POIs (hence 52% of correspondence).

Figure 5.9 shows the longest route found, recorded by 35 users and also the Popular POIs discovered, which are in the same area. Such places are found on 5 of 6 different time slots, from 00:00:00 to 19:59:59 (near flow number 21), and they correspond to the sites: the Tsinghua Park for time Slot 1, Tsinghua University West Stepping Classroom for time Slot 2, the Tsinghua University Human Resources Service & Employment Center for time Slot 3, the Tsinghua University Biomedical Library for time Slot 4 and the High School Attached to Tsinghua University for time Slot 5.

Other popular POIs near the flows detected in Geolife dataset are: China Academy of Space Technology, China Aerospace Zhongguancun Astronautics Community, Satellite Building Parking Lot, Jade Palace Hotel Office Building and Kangtuo Science e Technology Mansion.

The above analysis shows that the correspondence among flows and popular destinations is 50% to 68%. Such values are high when we consider that we matched the flows of common people, rather than simply tourists. For the considered flows, many people would go to their homes and working places, many of which are not popular destinations.

We have used three datasets in order to show the generality and robustness of our approach in finding POIs, and flows of people. Finally, the analysis let us reveal how close the main flows of people are to the real points of interests, which were taken from a curated list. Hence, we can show that the suggestion of Points Of Interest can actually be convenient for users who need not travel long distances if they accept the suggestion.

5.4 Relations with the proposed multi-agents system

POIs provide the initial knowledge base to the agents. As, at the beginning, when agents have not yet collected and shared the user opinions, the multi-agent system is based on data obtained from the detection of POIs and flows, as explained above. Therefore, at the beginning all POIs have been computed by means of the approach detailed in Section 4.2.2 and used as a knowledge base for the multi-agent system. The POIs are defined by the name and GPS coordinates. Such



Figure 5.9: POIs of 5 time slots near Geolife's flow 21.

1653 points have been identified by using the SPs, and other details (i.e. the name of the place) can
1654 be obtained, via APIs, since available on the web using services such as e.g. OpenStreetMap or
1655 Google Maps.

1656 Experiments have shown how POIs vary according to the needs of the users. For available
1657 data, by analysing the flows of people for the different datasets, we could determine that users
1658 moving on foot mainly want to visit banks, bars, restaurants; whereas users moving by taxi
1659 (generally tourists) mainly want to visit popular places like churches and museums.

1660 This allows us to offer a more efficient recommendation system, since the POI emerges from
1661 the gathered user positions and their trajectories, hence are related to the users.

1662 In addition, thanks to the GPS position taken by the agent, the POIs are recommended
1663 based on the users distance from it and near the POI, the user's position is monitored, taking
1664 care to protect his privacy, as explained in the multi-agents section 4.2.3 and Section 4.2.4.

1665 Finally, data and suggestions offered by means of the proposed multi-agents system are
1666 enriched by using the flow detection.

1667 Such an approach has allowed us to validate the POIs previously found, and then to highlight
1668 the time slot that is generally preferred by users. The analysis has been performed for six time
1669 slots, in this way the recommendation system improves its efficiency, offering information on
1670 time slots. Suggestions to users are always updated and improved by the exchange of data
1671 between the independent agents, guaranteeing a reliable and updated service to the end user.

1672 About the multi-agent architecture proposed by us, summarizing:

- 1673 1. The assistants in the literature have been used as a personal assistant (via app) and we
1674 have proposed an assistant who proposes the place to visit (first characteristic of the agents:
1675 *ASSISTANCE*);
- 1676 2. each agent has an *AUTONOMY* (second characteristic of the agents) and we have a
1677 (mobile) device in which the individual application does not depend on the applications
1678 of the other devices (*INDEPENDENCE*);
- 1679 3. there is a *COLLABORATION* (third characteristic of the agents) through the central
1680 server with the mediator and they communicate with a client-server approach.

1681 5.5 Test performed to investigate probability of users' movement

1682 We considered the trajectories of Cabspottingdata and the data cleaning was done in order to
1683 eliminate noise, due for example to GPS errors. It was performed by computing the instantane-
1684 ous speed of each point of the rides recorded on the taxi. The maximum acceptable speed
1685 threshold has been set for $150 \frac{km}{h}$. We considered 6 time slots of 4 hours each, to visualize the
1686 movement of the vehicles at different times and the trajectories were therefore split according
1687 to the 6 time slots.

1688 For each trajectory we apply the StayPoint detection algorithm (see its pseudocode in Algo-
rithm 3) with time threshold, TimeThr, equal to 10 minutes, and distance threshold, DistThr,

1689 100 meters. Such thresholds should suffice to select the positions in which a user dwells (in sev-
1690 eral *SPs*) as he finds the place interesting, and removes the locations where a user is stopping
1691 because e.g. he is blocked at the traffic light.

1692 The execution time for the *SP* detection algorithm on the whole Cabspotting dataset (536
1693 taxis and more than 11M points) was 36 minutes and 54 seconds. We obtained a total of 4261
1694 *SPs*, which is an average of 8 *SPs* per vehicle journey. The results show that 98% of users have
1695 at least one StayPoint associated with their trip (523 users out of 536). The implementation of
1696 StayPoints detection algorithm used Python and the experiments were executed in a host having
1697 an Intel Xeon CPU E5-2620 v3 2.40GHz, with RAM 32 GB. Figure 5.10 shows the recorded
1698 trace for each trip in blue, and the detected *SPs* in yellow.

1699 To predict people’s movements, firstly a grid was built, covering the map of all *SPs*: it
1700 consists of square cells with a side of 1 *km*. This grid allows us to discretize the data and
1701 estimate the probability of movement from a cell that has one *SP* inside it in another cell
1702 also having at least one *SP*. There are two distinct geographical areas that include some *SPs*
1703 represented as two square cells without intersection, therefore a spatial partition is formed.
1704 Figure 5.11 shows such a grid, having dimensions 80 × 46 cells (latitude by longitude) and the
1705 *SPs* obtained are mapped into the grid and displayed as a red diamonds.

1706 Some cells have much more *SPs* than others, so some red diamonds are denser in some areas
1707 than others, such as shown in said figure. To determine if cell *A* is a frequent destination, the
1708 support for each *SP* cell was calculated. The **Support** of a cell is the ratio between the number
1709 of trajectories that contain the cell and the total number of trajectories present in the grid. If
1710 this ratio exceeds a certain threshold, i.e. if cell *A* is crossed by a certain number of different
1711 trajectories (a value of 10% has been chosen for **Minimum Support**, which is 0.1), then the
1712 cell (containing one or more *SP*) will be a *frequently visited cell*.

1713 Our experiments on the taxi dataset above showed that there are 43 cells visited by a number
1714 of users greater than or equal to 52. That is, we can say that in the dataset there are 43 frequently
1715 visited *SPs* cells. This means that there has been a probable meeting in that cell, as users have
1716 remained stationary in the same time slot in the same cell. The data are updated in real time
1717 through the agents running on smartphones as an app, therefore the *Minimum Support* is fixed,
1718 however the number of frequent cells in output and the position of these frequently visited cells
1719 will vary over time. By lowering the Minimum Support, i.e. the threshold of the minimum
1720 amount of people sharing the same cell, the number of cells considered as having a sufficient
1721 amount of people will increase and then the number of cells considered overcrowded will increase.

1722 In order to compute *Association Rules* only between the frequently visited cells in the dataset,
1723 we considered the **Confidence** of the *Market Basket Analysis* for our approach. Given two cells
1724 called *A* and *B* we have that:

$$Support(A \Rightarrow B) = \frac{Frequency(A, B)}{N}$$

1725 the Support of the association rule ($A \Rightarrow B$) denotes the percentage of trajectories containing
1726 *A* which contain also *B*, where *N* is the total number of trajectories.

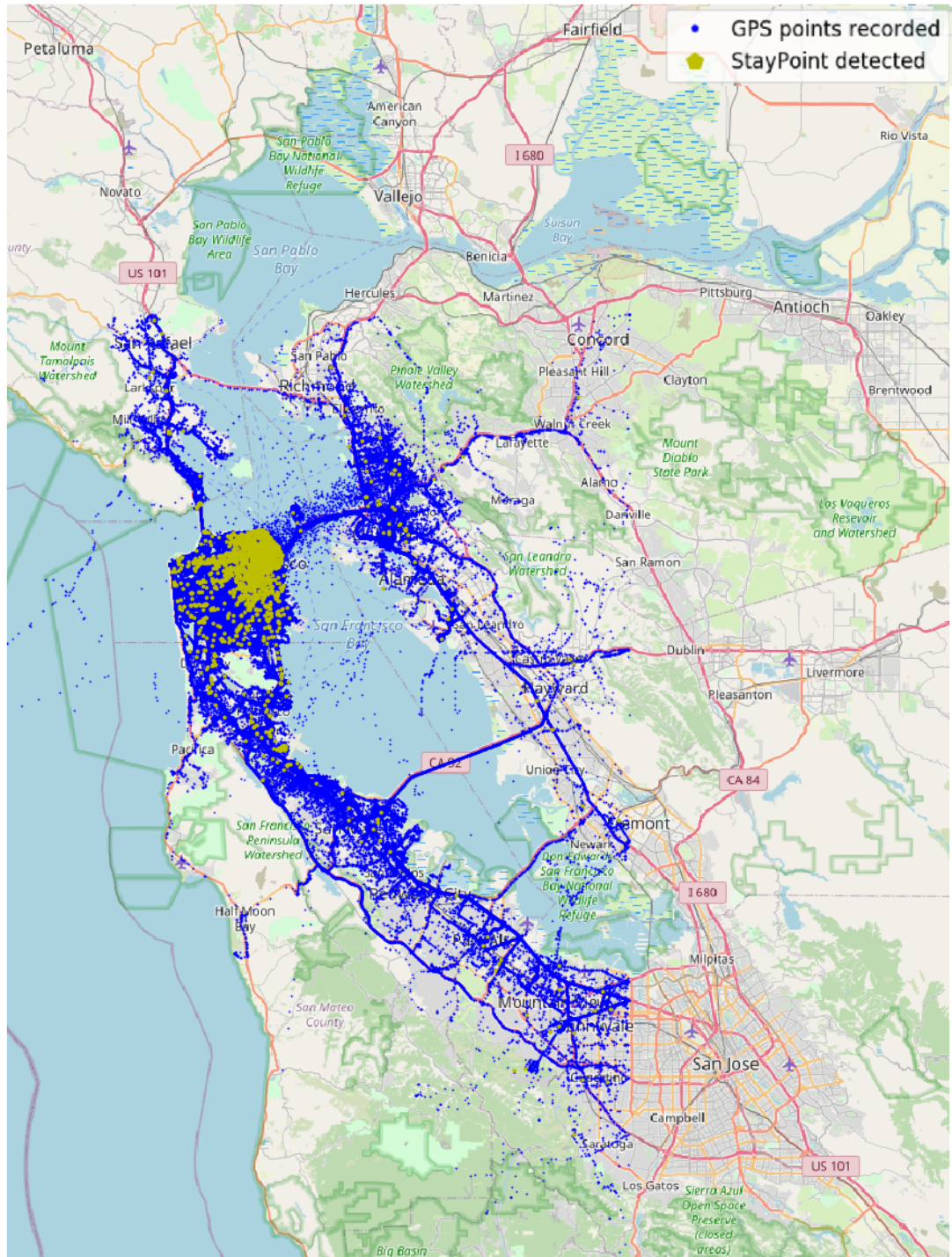


Figure 5.10: Blue points for trajectories and StayPoints obtained in yellow.

$$Confidence(A \Rightarrow B) = \frac{Support(A \Rightarrow B)}{Support(A)}$$

Hence, confidence is an estimation of conditioned probability, which can be expressed as follows:

$$Confidence(A \Rightarrow B) = \frac{p(A \cap B)}{p(A)} = P(B|A).$$

In Market Basket Analysis, *Confidence* is the probability of purchasing item *B*, said *consequent*, given the purchase of object *A*, said *antecedent*, within the same transaction. The higher the *Confidence*, the greater the reliability of the ($A \Rightarrow B$) rule (more details can be found in [45]). In our context, the value computed as the $Confidence(A \Rightarrow B)$ gives the probability that a user is in a *SP* in cell *B* moving there together with at least 10% of the total number of users, if he has already been in cell *A* and dwelling in one of its *SPs*.

Going forward along this procedure, we compute $Confidence(A, B \Rightarrow C)$ and after that $Confidence(A, B, C \Rightarrow D)$, in order to determine a common path that crosses several cells having highly visited *SPs*. We compute:

$$Confidence(A, B \Rightarrow C) = \frac{Frequency(A, B, C)}{Frequency(A, B)}$$

and so on.

Therefore, the results obtained are useful to predict the number of gatherings on some place. Moreover, given that there is knowledge about an infected person on some area, our results can be used to predict whether a user can be potentially infected (as his trajectory is estimated), and predict who else he will infect (i.e. people whose trajectories are expected to pass through the same areas).

The Confidence limit is due to the fact that it does not consider the Support of the item on the right side of the Association Rule and therefore does not provide a correct evaluation in case the groups of items are not stochastically independent.

A measure that takes this eventuality into account is $Lift(A \Rightarrow B)$, defined as:

$$Lift(A \Rightarrow B) = \frac{Confidence(A \Rightarrow B)}{Support(B)} = \frac{p(A \cap B)}{p(A) * p(B)}$$

$Lift(A \Rightarrow B)$ takes into account the importance (the Frequency) of *B*. Using such an amount, then we can say

- if $Lift > 1$ the events are positively correlated;
- if $Lift \leq 1$ the events are negatively correlated or independent.

Therefore, **Lift** indicates how the occurrence of one event raises the occurrences of the other. At this point, the setting would be a Minimum Confidence (0.6 i.e. 60%) to skim the results and obtain only the Association Rules that had a higher Confidence and also a Support higher

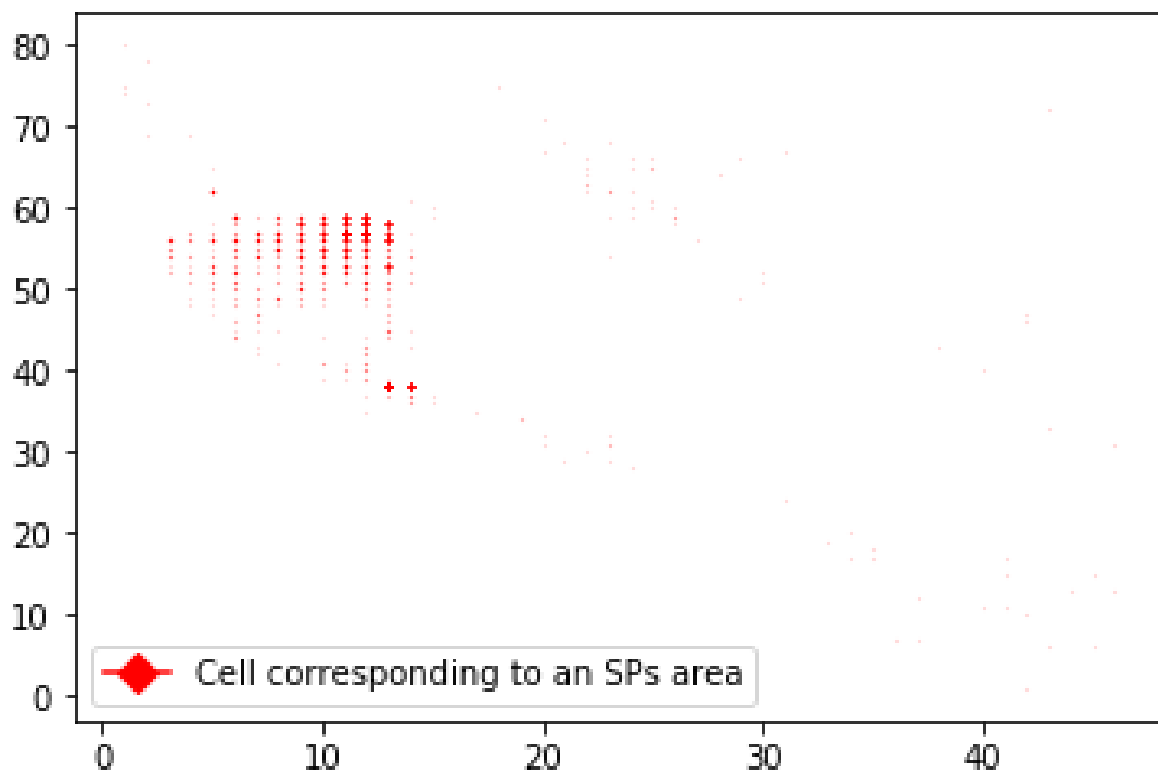


Figure 5.11: Grid formed by squared cells of 1 km per side, each red dot represent a cell having at least one *SP*.

1755 than the Minimum Support (0.1 chosen), such a setting is named *Strong Rules*. Finally, these
 1756 rules were checked with the Lift, the last column of Table 5.6. Then, for the Association Rule
 1757 ([2587] \Rightarrow [2588]) in row 6 the events of movement from cell *A* to cell *B* are negatively correlated.

1758 Figure 5.12 shows the plot of every Strong Rule obtained as a point, as a value for its Support
 1759 and Confidence (the latter according to the Support).

1760 For Association Rules with higher support the Confidence, that is the probability of moving
 1761 to the frequently visited cell *B*, decreases.

1762 The Lift and the Confidence of the Strong Rules obtained are directly proportional, as we
 1763 can see in Figure 5.13. The Pearson correlation coefficient between them is 0.9999999999999999
 1764 and this implies an exact linear relationship.

1765 At the end, in order to find common subtrajectories between different users, the trajectories
 1766 of Cabspotting have been summarized for the comparison of the distance: we considered for
 1767 each path two successive points in temporal order only if they were at a minimum distance of
 1768 140 m. This was done to decrease the size of the dataset and therefore will allow a reduction
 1769 in the execution time of the algorithm. To carry out a statistical analysis, 90% of trajectories

Table 5.6: Strong Association Rules: Minimum Confidence 60%, Minimum Support 10%.

$(A \Rightarrow B)$	A	B	$Support(A)$	$Support(B)$	$Support(A \Rightarrow B)$	$Confidence(A \Rightarrow B)$	$Lift(A \Rightarrow B)$
1	[1715]	[2588]	0.395793	0.671128	0.284895	0.719807	1.072533
2	[1716]	[2588]	0.313576	0.671128	0.223709	0.713415	1.063008
3	[2405]	[2588]	0.151052	0.671128	0.108987	0.721519	1.075084
4	[2541]	[2588]	0.137667	0.671128	0.099426	0.722222	1.076132
5	[2543]	[2588]	0.242830	0.671128	0.175908	0.724409	1.079391
6	[2587]	[2588]	0.435946	0.671128	0.292543	0.671053	0.999888
7	[2633]	[2588]	0.145315	0.671128	0.103250	0.710526	1.058704
8	[2634]	[2588]	0.281071	0.671128	0.206501	0.734694	1.094715
9	[2635]	[2588]	0.235182	0.671128	0.177820	0.756098	1.126607
10	[2679]	[2588]	0.202677	0.671128	0.147228	0.726415	1.082379
11	[2680]	[2588]	0.242830	0.671128	0.175908	0.724409	1.079391
12	[1715, 1716]	[2588]	0.149140	0.671128	0.112811	0.756410	1.127073
13	[2587, 1715]	[2588]	0.156788	0.671128	0.112811	0.719512	1.072094
14	[2634, 2587]	[2588]	0.141491	0.671128	0.101338	0.716216	1.067183

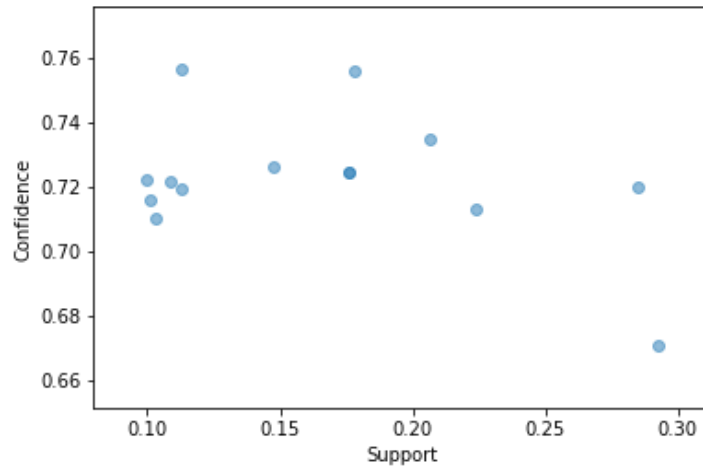


Figure 5.12: Support vs Confidence for Strong Rules obtained by our analysis and given as rows in Table 4.1. The points show a pair of cells (A and B), or a triple of cells (for the last three rows of Table 4.1).

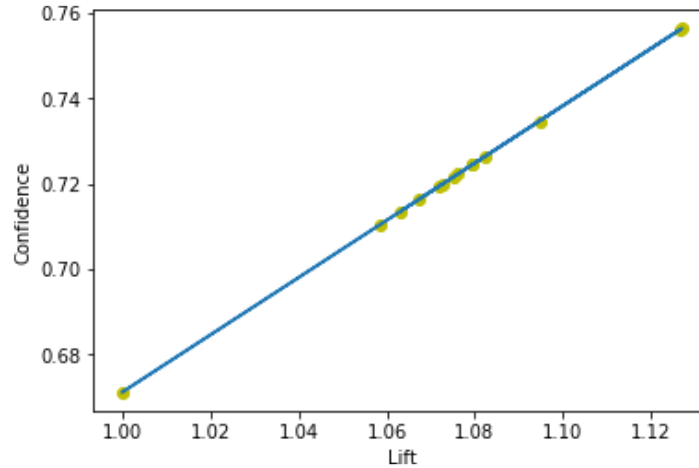


Figure 5.13: Lift vs Confidence in this test, as for the points shown in Figure 5.12.

1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793

were randomly selected, and this set was the *Train* set for the the flow detection algorithm. The complementary set, that is the remaining 10% of the trajectories, consists in the *Test* set, that is the verification set.

To identify the sub-trajectories common to different users in the same time slot, a maximum tolerance distance was set between two different points of different users as 280 *m*. The distance between two points was computed by using the Haversine distance.

We define **flows** as close sub-trajectories, belonging to different users, spatially similar and recorded in the same time slot. The density of a flow is the number of users that pass through it. To detect flows in this dataset, the minimum density threshold was set to 25. According to these parameters, 12 flows were identified, ranging from 1 to 2 *km* in length. The minimum density of the flows found is 26 taxis, while the maximum density found is 192 taxis. Then, by taking the complement of the trajectory sample (10% of the taxis, as the *Test* set) we checked where their GPS points were compared to the previous train set. We found that the points of the *Test* set intersect with the 12 paths identified on the *Train* set. Another check was carried out by confirming the correspondence of the points of the flows on a map. It consists of the process of matching the coordinates of the obtained flows and the road segments, and assessing that there are no external points with respect to road segments (see Figure 5.15).

Figure 5.14 shows the detected flows in magenta and the *SPs* in that area in green.

Moreover, the results tell us that the probability of transitioning from one cell with *SPs* to another is high even in correspondence with the indicated flows and that different highly visited cells having *SPs* belong to different flows. For each cell we have checked which taxis passed there and which passed at a later time on other flows passing through other frequently visited cells.

Our proposed work uses a multi-agent system communicating with a server that acts as

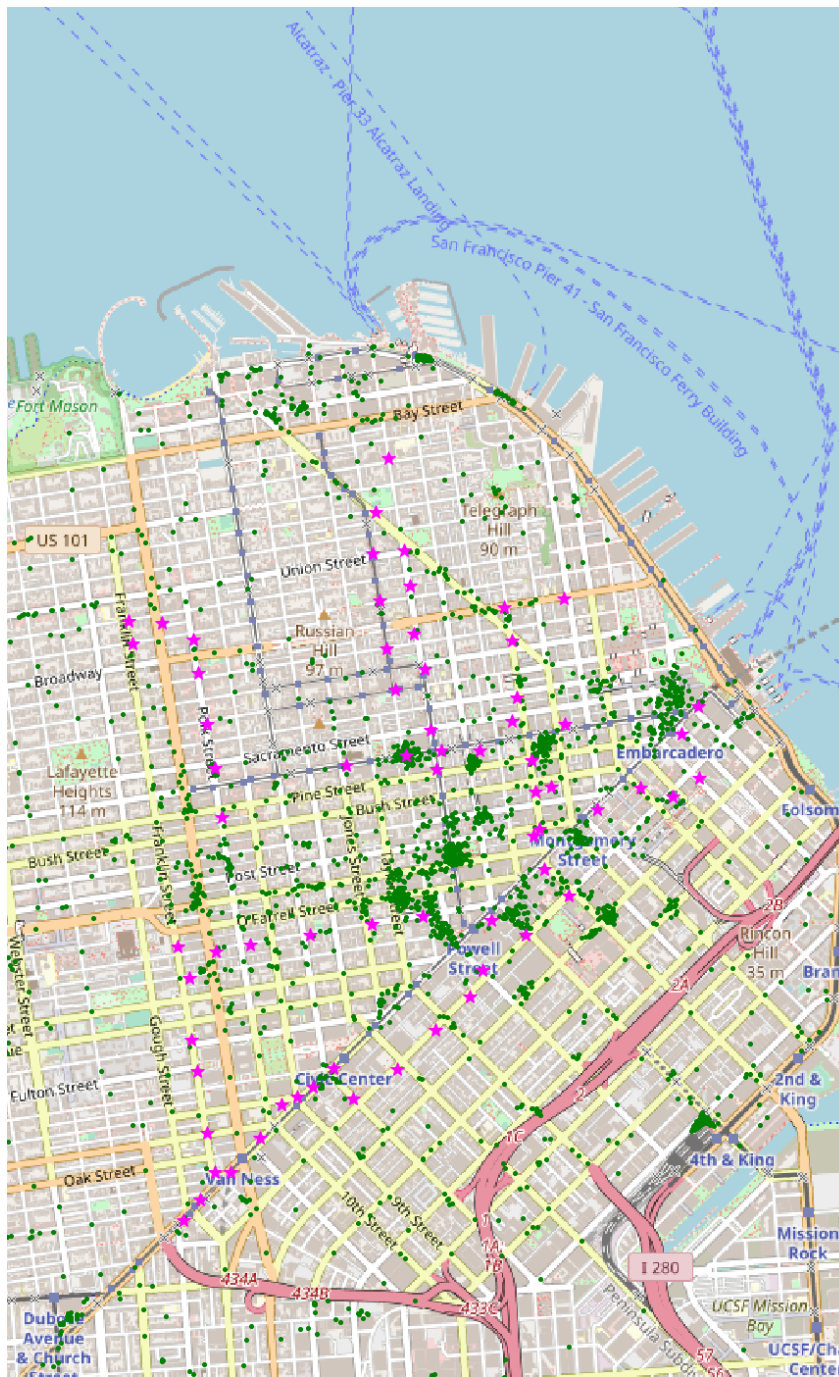


Figure 5.14: A zoomed in map of an area in Figure 5.10, showing flows in magenta and the nearby StayPoints in green.

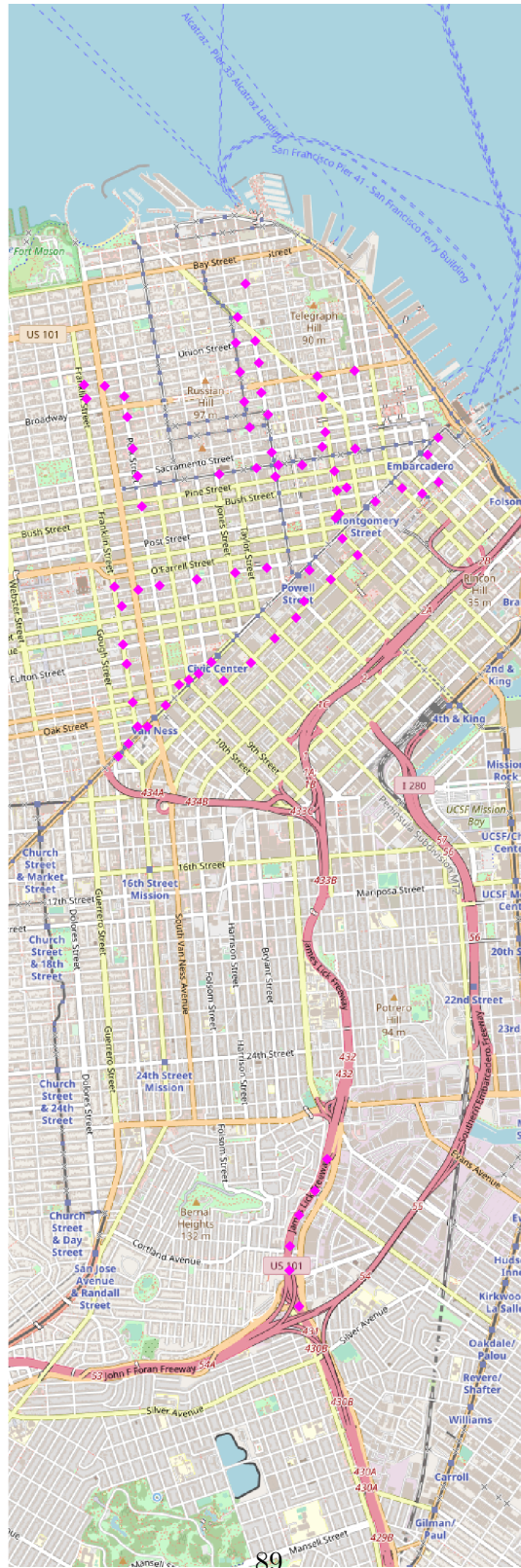


Figure 5.15: Flows detected for the Cabspting dataset.

1794 an intermediary and predicts movements by means of an innovative and reliable mathematical
1795 solution. Unlike the work presented in [11], our goal is to detect the foreseeable routes by
1796 computing their probability, whereas their method determines the probability that a group of
1797 users is moving together.

1798 The paper [32] presents personalised recommendations for guiding tourists through the city
1799 of Melbourne by observing their actions. This system is modelled as a Markov decision process
1800 that recommends the user in sequence the next place to visit. However, unlike the StayPoint
1801 analysis presented here, it does not consider the stationary nature of visitors over a period of
1802 time and this is a key element in avoiding overcrowding.

5.6 Different methods by varying the corridor parameters

1803 The output obtained is frequent K -tuples of cells, that are shared by at least $n \geq \text{min_support}$
1804 (number of trajectories) Apriori uses a bottom-up approach, starts from the smallest subsets
1805 (the single cells) and checks whether they are *Frequently*.

1806 For $K \geq 2$ since if an element (set of cells) is frequent, so too are its subsets, it generates the
1807 cardinality candidates $K + 1$ only those formed from the frequent K -tuples. This allows you to
1808 greatly reduce counting calculations.

1809 The execution time is 1 hour and 22 minutes for 15 trajectories, considering instead the
1810 inverted problem: cells = baskets that contain/are crossed by the trajectories the algorithm
1811 ends in about 22 minutes, with a min_support of 55 cells of 500 meters each. So the advantage
1812 is the reduction of execution times significantly. The algorithm is very fast on the sample of
1813 1,000 trajectories and with a minimum support of 50 trajectories it ends in 297 *ms*.

1814 We compared our proposed method with the computation of Discrete Fréchet distances
1815 between curves, which can be used to find those close to each other within a certain fixed
1816 tolerance. Considering the same random set of trajectories in the first time slot (a total of 1000
1817 curves), the DDF between pairs of trajectories (one at a time) is determined in a total time of
1818 6 hours and 24 minutes. Given that the total average number of trajectories for each time slot
1819 is about 4500 this strategy would be overly time-consuming (about 4 times and half more) than
1820 the three strategies we implemented (in the worst case 1 hour and 22 minutes have passed).
1821 Furthermore, determining the matrix of the distances for the trajectories to be compared would
1822 be impractical. In the machine available for our experiments we get a memory error, it is not
1823 possible for us to memorize a matrix so large.

Experimental results of Step 1

1824 The application of Apriori results in the list of the K -frequent itemsets (i.e. subsets of frequent
1825 cells of cardinality K , the second parameter of the corridors), where frequent means that they are
1826 shared by a number of trajectories greater than or equal to the minimum support of trajectories
1827 chosen as Apriori parameter in input.

1828 In Step 1 we chose 15 trajectories for each time slot as the minimum support threshold. For
1829 example, the total number of trajectories in the first time slot is 5,878, so the set percentage of

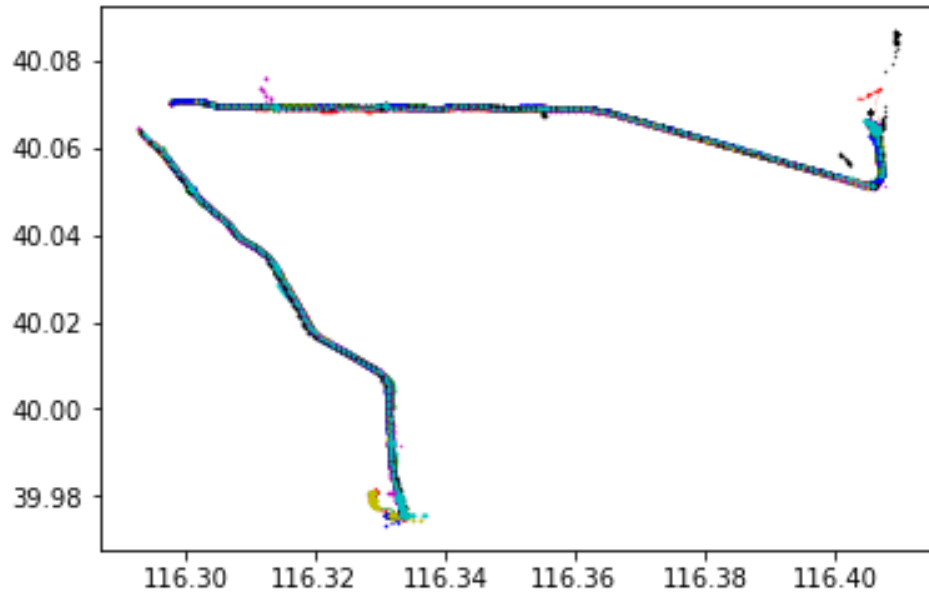


Figure 5.16: A candidate corridor detected with Step 1, with parameters $M = 19$ cells and $S = 15$ trajectories.

1830 *min_sup* on the total was low but the advantage is to be able to find corridors even for a limited
 1831 number of trajectories. Setting the minimum number of trajectories, the length of the corridor
 1832 candidates obtained at each iteration of the Apriori algorithm increases.

1833 In our case, Apriori returns a list of over 4 million K -frequent itemsets of cells, for K from
 1834 1 to 21. The largest subset of cells, which is the last iteration of the algorithm, corresponds to
 1835 a path of about 25 km. At this point we considered all the trajectories passing from the first
 1836 cell of the 21-tuple of cells. We deleted all the trajectories that did not pass for the second cell,
 1837 then we deleted the trajectories that did not pass for the third cell and so on, until the 21st cell.
 1838 At the end we found 15 trajectories, that are exactly the trajectories passing for all the 21 cells.
 1839 After verifying that these trajectories cross these cells in the same or in the reverse order and
 1840 consecutively, we considered the 21-tuple as a candidate corridor. Figure 5.16 shows the longest
 1841 candidate corridor obtained in Step 1: the interruption in the upper left part of the graph is
 1842 probably due to the loss of the GPS signal during the recording.

Experimental results of Step 1.2

1843 The output produced by the algorithm consists of 55,109 sets of K elements of frequent trajec-
 1844 tories (in the time slot Δt_1), for K from 1 to 14: these trajectories share candidate corridors of
 1845 over 25 km.

1846 The advantage is that, compared to the previous step, it improves the execution time. Indeed
 1847 the scanning of the trajectories to check if they are frequent or not is faster because they have

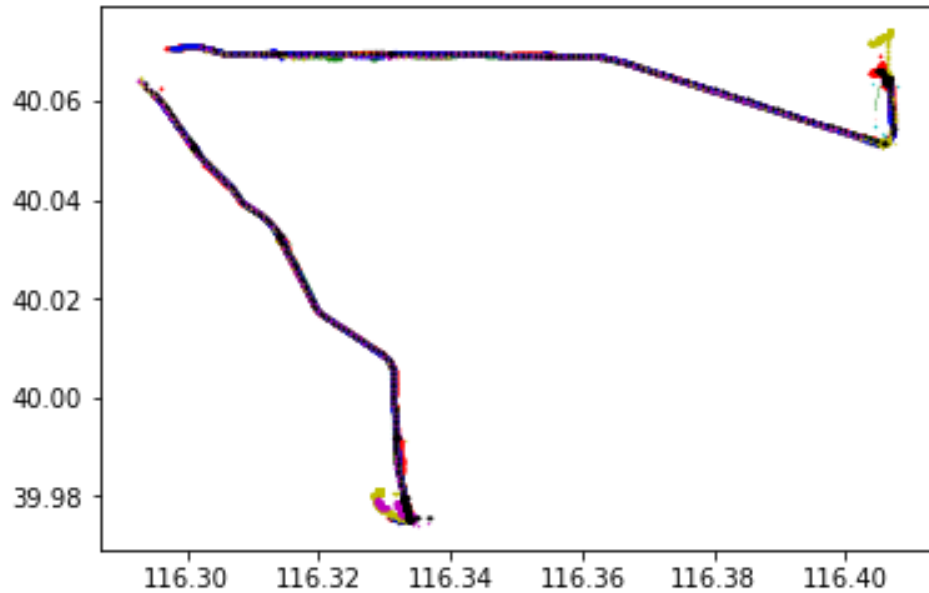


Figure 5.17: Candidate corridor detected with Step 1.2.

1848 few elements (on average they have a length of about 16 cells of 500 *m* each) and also the
 1849 support threshold has been raised. In fact, if the *min_sup* is larger, then the cardinality of the
 1850 frequent set L_K and therefore also those of C_{K+1} are much smaller than in Step 1. Moreover,
 1851 the number of iterations of the algorithm is also lower (it made 14 iterations instead of 21).
 1852 Another improvement was the use of a thinner grid which gives a greater precision of the road
 1853 results.

1854 The longest corridors of these 14 trajectories are similar to the ones found in the previous
 1855 case (Step 1). Among these 14 trajectories, 7 are in common with the 15 obtained in the Step
 1856 1. The plot of trajectories crossing 51 cells consists of about 25 *km* visible in Figure 5.17, and
 1857 like the previous case we verify that this sequence is a candidate corridor of parameters $M = 51$
 1858 and $S = 14$.

1859 Step 1.2 is a verification of Step 1, also allowed to refine the grid and therefore to have a
 1860 greater precision (for instance we used a grid with side cells 500 meters wide instead of a grid
 1861 with cells 1,300 meters per side as in Step 1 - see Figure 4.15 -).

1862 A further advantage is the reduction of execution times significantly. We found two candidate
 1863 corridors, in the same street shared by a union of 22 trajectories, found in Steps 1 and 1.2.

1864 The grid is out of phase, because 500 (meters) is not a multiple of 1,300 (meters), but this
 1865 gives a greater number of trajectories that are recorded on the same road, with the union of the
 1866 results obtained.

Experimental results of Step 2

The output obtained for data in Δt_1 was a list of frequent itemsets of cells, a table of 1,398 K -frequent itemsets, for K from 1 to 9. The longest sets of cells were three 9-frequent cells:

$$F = \{1956, 1703, 1804, 1550, 1905, 1906, 1652, 1754, 1855\},$$

$$G = \{1956, 1703, 1550, 1905, 1906, 1652, 2007, 1754, 1855\} \text{ and}$$

$$H = \{1956, 1703, 1804, 1905, 1906, 1652, 2007, 1754, 1855\}.$$

We set $A = F$. Among the 8-frequent itemsets we looked for those that were not a subset of A . If there was a 8-frequent itemset that shared at most one cell with F , then we saved it as a potential corridor. At the end, we added to A all the elements of the 8-frequent itemsets that were not in A and we repeated the procedure with the 7-frequent itemsets. Continuing in this way until the 2-frequent itemsets, we found 7 potential corridors with at least 2 cells.

Doing the same starting from G or H we found the same potential corridors of parameter M from 2 to 5:

$$C_1 = [1601, 1602, 1603, 1604, 1656], \text{ shared by 53 trajectories;}$$

$$C_2 = [1753, 1702, 1703], \text{ shared by 58 trajectories;}$$

$$C_3 = [1548, 1549], \text{ shared by 58 trajectories;}$$

$$C_4 = [1548, 1599], \text{ shared by 64 trajectories;}$$

$$C_5 = [1550, 1551], \text{ shared by 91 trajectories;}$$

$$C_6 = [1600, 1601], \text{ shared by 74 trajectories;}$$

$$C_7 = [1804, 1805], \text{ shared by 53 trajectories.}$$

The performance of this procedure depends on the number of K - and $(K - 1)$ -cells and size K for every iteration and the number of new cells found in until level K .

Starting for example from the candidate corridor C_1 obtained, we carried out the filtering which is the final phase of our work. The plot of this candidate corridor is shown in Figure 5.18.

5.6.1 Filtering process

Given the problem of approximation of the grid, the previous results does not assure that all sub-trajectories have a similar initial and final part. Because of this, we must restrict the search for the corridor to a sub-area, starting from the determined corridor candidate.

Given a radius r set by the user, the Radius Neighbors Graph algorithm constructs a graph in which all the points of a trajectory, seen as nodes of the graph, are connected to all the other points of the other trajectories that are close to it within distance r . The adjacency matrix allows to quickly distinguish which sub-trajectories are in the neighborhood and which sub-trajectories are to be discarded.

Every corridor selected with this filtering phase (applying Radius Neighbors Graph to the cells containing the candidate corridors) is a road segment that can be considered as a “hot” route. The result of the filtering process on candidate corridor detected above is shown in Figure 5.19. Finding the frequently repeated patterns of sub-routes can help to analyze and predict the movements of objects; starting from the results obtained we could propose an estimate of the similarity between users (who share corridors) using the Jaccard Similarity index.

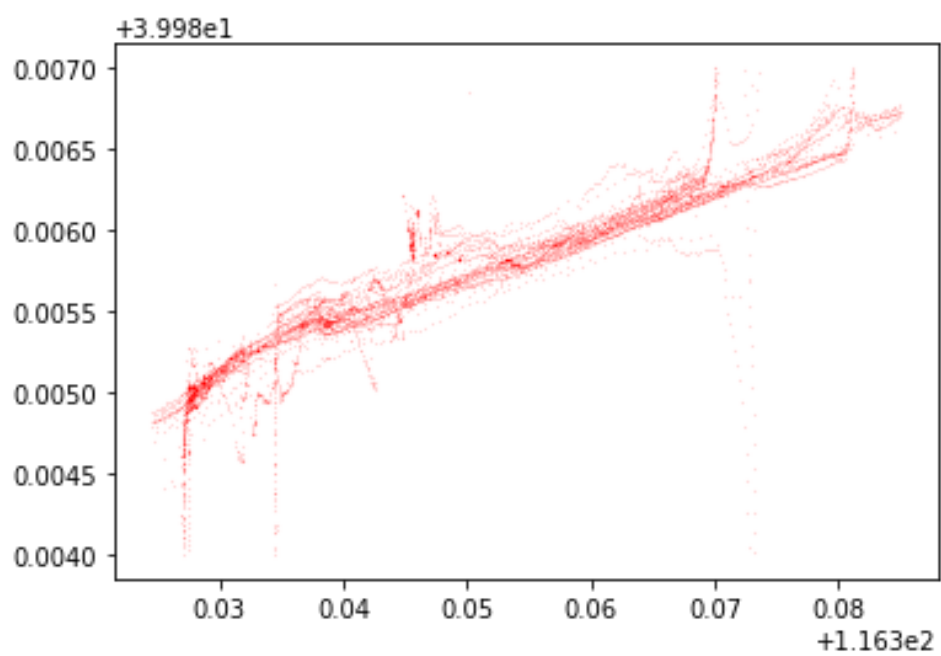


Figure 5.18: C_1 candidate corridor detected in Step 2, with parameters $M = 5$ cells and $S = 52$ trajectories.

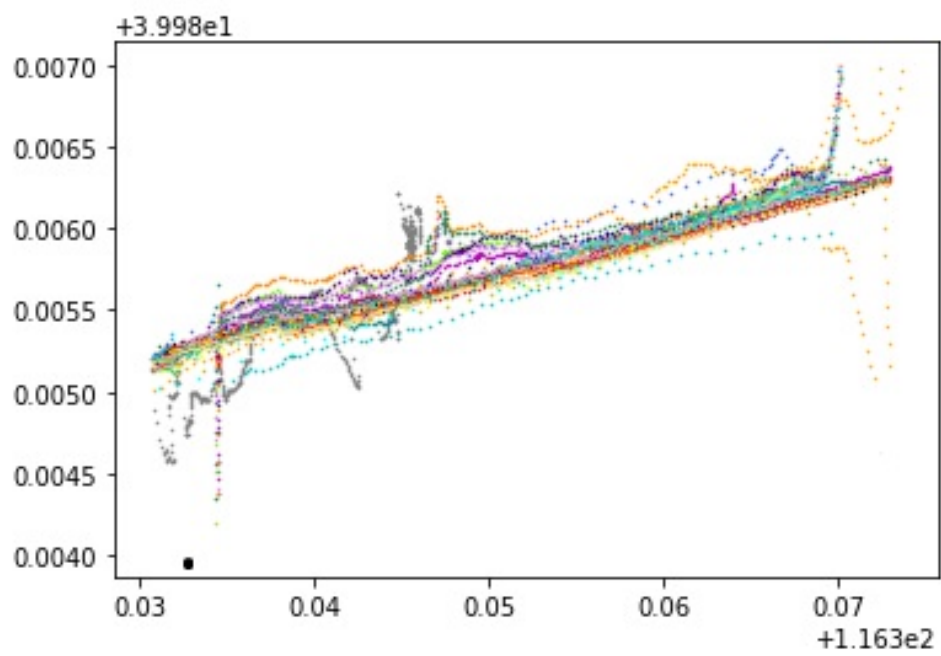


Figure 5.19: A final corridor obtained by the filtering process, shared by $S = 50$ trajectories and $3,7 \text{ km}$ long.

Chapter 6

Final discussion, remarks and future work

1905 Given the widespread use of mobile devices and various technologies that track their geographical
1906 position, it is nowadays easier to acquire information about users' GPS in real time. This
1907 availability has triggered numerous studies based on user positioning, such as analyzing flows of
1908 people [93] in cities or predicting people's movements. This has also led to the improvement of
1909 the services that identify the points of interest of a city to offer advantages to users who want
1910 to reach a place but do not have sufficient knowledge for an immediate choice.

1911 The problem of "hot routes" detection is important for transport services, which are inter-
1912 esting for finding corridors for public transport of multiple routes (e.g. by bicycle, subway, car
1913 sharing) taking advantage of GPS sub-trajectory datasets.

6.1 Epilogue

1914 The first proposed approach provided a solution for sharing collective knowledge about popular
1915 and useful points about a city. Users can share their comments on a place they visit thanks to
1916 an app available on their smartphones and can receive suggestions on the next place to visit.

1917 For the proposed multi-agent system, a setting in our agent application on the smartphone
1918 allows to collect GPS coordinates. This is useful for a geographical position where there are no
1919 previously collected trajectories. Thus, we can continuously extract user trajectories and update
1920 both user tips and POI recommendations. Using this setting, each agent periodically releases
1921 both the GPS coordinates and his identity to the server, however, to take into account privacy
1922 issues, the user's identity is masked and the GPS position is randomly moved up to 300 meters.
1923 So, we get a trajectory that is not very precise, but still useful. As trajectories are dynamically
1924 collected, POIs are also determined dynamically as data arrive to the server.

1925 The experiments in GeoLife dataset confirm that users stop in the same areas for some
1926 common reasons, such as visiting a tourist attraction or taking advantage of the same service
1927 and remain in certain areas in common time slots. Additionally, these tests reveal that people
1928 move together from one POI to another.

1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965

The main contribution of the first part of our work is:

- automating the extraction of the list of places of interests;
- updating such a list in real time;
- characterizing places of interests with time slots;
- providing to the user comments on places of interests;
- selecting comments and places of interests according to user preference and position;
- preserving privacy for the user.

In the second study we carried out we analyzed two other geo-datasets real thunderstorms, Taxi and Truck Trajectories with different characteristics, applying the algorithms explored but with parameters adapted to the cases in question.

By analysing the available people positions, it has been shown that it is possible to find: Points Of Interests and people flows.

To validate the new Points Of Interest found, flows shared by different users (within a fixed radius) were identified by similar sub-trajectories. The density of the corridors that we were looking for (it is equal to the number of users that passed through the same corridor) has been fixed and a great match between these and the points of interest identified was found.

Experiments confirm the possible hypothesis that the vehicles have been booked to reach a site of interest or after visiting one of these to leave another destination, also being a popular place.

The POIs found were verified by means of Google Maps and compared with a curated list brought to us by a native of the analysed places. For each dataset used, consisting of data referring different categories of people and movements, it was shown that the points of interests are compatible with the nature of the dataset and the behaviour of the users who collected them. In fact, Points Of Interests have been found in service areas, supplies and spare parts sales outlets with regard to users with trucks; service buildings and companies for businessmen travelling by taxi; and places of study and work for participants of the Microsoft Research Asia Geolife project.

A third path followed during my PhD, that in particular I started during a period of study and research at the Universitat de Barcelona in Spain under the supervision of Professor J. Vitrià, concerns the identification of the corridors, which are shared by a considerable number of users.

The dataset has been discretized in order to solve the problem by applying the Apriori algorithm. Having an iterative layer-by-layer approach, it allowed us to check more quickly which common routes were frequent, starting only from frequent sub-trajectories smaller than the ones we wanted to identify and excluding the others. The new approach has reduced execution times compared to a comparison between pairs of trajectories, in which the discrete Fréchet distance or other similarity measures between curves are used.

1966 The total run time for Apriori algorithm in every step is defined as (see [89] Suneetha and
 1967 Krishnamoorti, 2011):

$$\sum_{i=1}^n (t_s * m_k + t_c + p_{k+1} * k + 1/2 * t_s * n_k / B)$$

1968 where, t_s is the time cost of a single scan of the database, t_c is the time cost of generating
 1969 C_{k+1} from L_k , m_k is set to be the amount of itemsets in C_k , the variable $p_k + 1$ is set to be
 1970 the amount of itemsets in C_{k+1} and the variable n_k is set to be the amount of itemsets in L_k .
 1971 B is the number of records in the database and n represents the dimension of the data. The
 1972 Apriori execution time, with minimum support of 15 trajectories in Step 1 is 1 *h* and 22 *min*;
 1973 for *min_sup* of 51 cells in Step 1.2 the execution time is 21 *min* and 44 *s*; for 1,015 random
 1974 trajectories, with *min_sup* of 50 trajectories in Step 2: 297 *ms*.

1975 A subsequent filtering of the results, according to a desired tolerance, returned us outputs
 1976 cleaned from initial discretization errors.

1977 The application of the Apriori algorithm and the subsequent filtering of the results have
 1978 this interpretation: If different users start at the same point P in the city with \pm a certain
 1979 radius r , with r smaller than or equal to the side of the cell of the grid, and they share a second
 1980 consecutive cell, or more than one, then there is a high probability that such individuals move
 1981 together on the same road to reach another point in the same time slot. The corridor will end
 1982 when any of these N trajectories will not share the next cell (when the number of trajectories
 1983 “traveling” close together becomes less than the *min_sup* threshold).

1984 Finally, we have proposed an approach for predicting how many people move to certain
 1985 destinations when it is known that a certain amount of people is in some other place. We use
 1986 an app that detects the location of people and sends such data to a server. So, being able to
 1987 calculate how likely a group of people will move, along with the previous stats, is useful for
 1988 estimating the amount of people in another specific place to track later. Having an educated
 1989 guess on the amount of people that will gather in some place before planning a trip can be very
 1990 useful to avoid overcrowded places and to keep with the current regulations.

6.2 Conclusions

1991 It is possible to calculate the similarity between users simply by checking the number of times
 1992 each pair of individuals has visited the same region (POI). For this purpose, we can use the
 1993 Jaccard similarity between users, by calculating the set of intersections and the union set of
 1994 the regions visited for each couple of users. Computing the similarity of people based on their
 1995 position history is useful to understand the interaction of the individuals in a particular region.
 1996 Since GeoLife is a dataset composed primarily by GPS logs of academic individuals, here POIs
 1997 can be a university building, a residence, a restaurant or a research center. In fact, for the city
 1998 of Beijing, most of the POIs are being built around Peking University, Tsinghua University,
 1999 Microsoft Research Asia and some leisure places like Yuanmingyuan Park.

2000 For individuals with a high similarity, they are likely to know and have each other similar
 2001 visit / travel preferences. For the GeoLife dataset, these people can be, for example, univer-

2002 sity colleagues. This approach can be useful for suggestions from friends in social networks.
2003 Facebook, for example, suggests a friend based on the number of friends that two users have in
2004 common. Likewise, this approach can be applied in the same way in a social network based on
2005 geographic location about. In addition, this strategy can be applied for different GPS logs or
2006 other sources of geographic information than by users, in order to find the similarity of people
2007 based on their geographic location record, StayPoints or Points Of Interest.

2008 The use of spatio-temporal information aims to understand the movement patterns (corri-
2009 dors) of voluntary users in GeoLife and it can help to make an appropriate recommendation
2010 (to know, for example, in which time slot you mostly visit a certain urban area through certain
2011 roads).

2012 Therefore, the discovery of the flows in cities is very important as it has a wide range of
2013 applications: for example, it can facilitate the planning and optimization of transport services
2014 even in mass events to avoid the traffic.

6.3 Open Issues

2015 Another strategy is to use QuickBundles (QB) for clustering the trajectories for each user. With
2016 QB we could find for each individual few prototypes of trajectories (at each trajectory it assigns
2017 its closest trajectories) and finally compare pairs of clustered trajectories with the Discrete
2018 Fréchet Distance.

2019 QuickBundles, which is used in neuroimaging, finds dense road segments (i.e., with high
2020 number of trajectories) and merges them into dense routes (with similarity and minimizing the
2021 distance. At the end we will have to filter the trajectories that have low Fréchet distance (or
2022 also with Radius Neighbors Graph, with the help of the adjacency matrix). These methods will
2023 be the topics of future research.

2024 This analysis could have many interesting applications: for instance it would be possible to
2025 create a social network to connect people in the same travel location, give travel recommendation
2026 or point out the lacking of public transport lines in certain areas to the city council.

2027 The experiments that we have performed on previously gathered geographical locations have
2028 shown the viability and reliability of our approach. The more people use the app (our agent) the
2029 more the approach would give a correct estimate. To make the approach more robust, it could
2030 be extended in order to include data available online from other services that give indications
2031 on queues, road traffic and gatherings.

2032 Future work can take into account the geometry of some stations, museums, etc. of some
2033 popular destinations to compute the average distance of people given the estimated amount of
2034 people gathering. Moreover, alerts about overcrowding could be sent to both the people present
2035 in some place for a peak of incoming people, occurring later on, and the people that are moving
2036 towards there.

Bibliography

- [1] Agrawal, R.; Srikant, R. *Fast Algorithms for Mining Association Rules in Large Databases*. 20th International Conference on Very Large Data Bases. Morgan Kaufmann, Los Altos, CA, 1994, pp. 478–499.
- [2] Agrawal, R.; Srikant, R. *Privacy-preserving data mining*. ACM Sigmod Record, 29(2):439–450, 2000.
- [3] Ahas, R.; Silm, S.; Järv, O.; Saluveer, E.; Tiru, M. *Using mobile positioning data to model locations meaningful to users of mobile phones*. Journal of urban technology, 17(1):3–27, 2010.
- [4] Aliannejadi, M.; Crestani F. *Personalized context-aware point of interest recommendation*. ACM Transactions on Information Systems (TOIS), 36(4):45, 2018.
- [5] Alt, H.; Godau, M. *Measuring the resemblance of polygonal curves*. In Proc. 8th Annu. ACM Sympos. Comput. Geom., pages 102–109, 1992.
- [6] Ankerst, M.; Breunig, Markus M.; Kriegel, Hans-Peter; Sander, J. *OPTICS: Ordering Points To Identify the Clustering Structure*. ACM SIGMOD international conference on Management of data. ACM Press. pp. 49–60, 1999.
- [7] Batet, M.; Moreno, A; Sánchez, D.; Isern, D.; Valls, A. *Turist@: Agent-based personalised recommendation of tourist activities*. Expert Systems with Applications, 39(8):7319–7329, 2012.
- [8] Berndt, D. J.; Clifford, J. *Using dynamic time warping to find patterns in time series*. In KDD Workshop, pp. 359–370. AAAI Press, 1994.
- [9] Berzi, C.; Gorrini, A.; Vizzari, G. *Mining the social media data for a bottom-up evaluation of walkability*. In Proceedings of International Conference on Traffic and Granular Flow, pages 167–175. Springer, 2017.
- [10] Bian, J.; Tianb, D.; Tangc, Y.; Taod, D. *A survey on trajectory clustering analysis*. arXiv preprint, [arXiv:quant-ph/1802.06971].
- [11] Bicocchi, N.; Mamei, M. *Investigating ride sharing opportunities through mobility data analysis*. Pervasive Mobile Computing 14 (2014) 83–94.

- [12] Bicocchi, N.; Mamei, M.; Sassi, A.; Zambonelli, F. *On Recommending Opportunistic Rides*. IEEE Transactions on Intelligent Transportation Systems 2017, 18, 3328–3338.
- [13] Bobadilla, J.; Ortega, F.; Hernando, A.; Gutiérrez, A. *Recommender systems survey*. Knowledge-based systems, 46:109–132, 2013.
- [14] Boutilier, C. *Sequential optimality and coordination in multiagent systems*. In Proceedings of International Joint Conference on Artificial Intelligence (IJCAI), volume 99, pages 478–485, 1999.
- [15] Breese, J. S.; Heckerman, D.; Kadie, C. *Empirical analysis of predictive algorithms for collaborative filtering*. In Proceedings of the Fourteenth conference on Uncertainty in artificial intelligence, pages 43–52. Morgan Kaufmann Publishers Inc., 1998.
- [16] Buchin, K.; Buchin, M.; Gudmundsson, J.; Löffler, M.; Luo, J. *Detecting Commuting Patterns by Clustering Subtrajectories*. International Journal of Computational Geometry & Applications 2011, 21, 253–282. doi:10.1142/S0218195911003652.
- [17] Harabasz, Calinski T.; Karoński, M. *A dendrite method for cluster analysis*. Communications in Statistics. Vol. 3. No. 1. 1974. 1-27.
- [18] Canny, J. *Collaborative filtering with privacy*. In Proceedings IEEE Symposium on Security and Privacy, pages 45–57, 2002.
- [19] Canny, J. *Collaborative filtering with privacy via factor analysis*. In Proceedings of international ACM SIGIR conference on Research and development in information retrieval, pages 238–245, 2002.
- [20] Castro, P.; Zhang, D.; Li, S. *Urban traffic modelling and prediction using large scale taxi gps traces*, in Pervasive Computing, pp. 57–72, 06 2012.
- [21] Cavallaro, C.; Verga, G.; Tramontana, E.; Muscato, O. *Multi-Agent Architecture for Point of Interest Detection and Recommendation*. In Proceedings of the 20th Workshop From Objects to Agents, Parma, Italy. CEUR Workshop Proceedings, 2019, Vol. 2404, pp. 98–104.
- [22] Cavallaro, C.; Verga, G.; Tramontana, E.; Muscato, O. *Eliciting Cities Points of Interest from People Movements and Suggesting Effective Itineraries*. Intelligenza Artificiale, 14 (1), 49-61, 2020.
- [23] Cavallaro, C.; Verga, G.; Tramontana, E.; Muscato, O. *Suggesting Just Enough (Un)Crowded Routes and Destinations*. In Proceedings of the 21th Workshop From Objects to Agents, Bologna, Italy. CEUR Workshop Proceedings, 2020. Vol. 2706, pp. 237-251.
- [24] Cavallaro, C.; Vitrià, J. *Corridor Detection from Large GPS Trajectories Datasets*. Applied Sciences. 2020; 10(14):5003.

- [25] Chee, S. H. S. ; Han, J. ; Wang, K. *Rectree: An efficient collaborative filtering method*. In Proceedings of International Conference on Data Warehousing and Knowledge Discovery, pages 141–151. Springer, 2001.
- [26] Chen, C.-C.; Tsai, J.-L. *Determinants of behavioral intention to use the personalized location-based mobile tourism application: An empirical study by integrating tam with issm*. Future Generation Computer Systems, 2017.
- [27] Chessa, S.; Girolami, M.; Foschini, L.; Ianniello, R.; Corradi, A.; Bellavista, P. *Mobile Crowd Sensing Management with the ParticipAct Living Lab*. Pervasive and Mobile Computing (2016), <http://dx.doi.org/10.1016/j.pmcj.2016.09.005>
- [28] Cranshaw, J.; Schwartz, R.; Hong, J. Sadeh, N. *The livelihoods project: Utilizing social media to understand the dynamics of a city*. In Proceedings of International AAAI Conference on Weblogs and Social Media, 2012.
- [29] Crociani, L.; Vizzari, G.; Gorrini, A.; Bandini, S., *Identification and Characterization of Lanes in Pedestrian Flows Through a Clustering Approach*. In AI*IA 2018 – Advances in Artificial Intelligence. AI*IA 2018. Lecture Notes in Computer Science; Springer Verlag, 2018; Vol. 11298, pp. 71–82. doi:10.1007/978-3-030-03840-3_6.
- [30] Cruz, J.; Silva, E.; Rossetti, R. J.; Silva, D. C.; Oliveira, E. C. ; Neto, J. *Application of multi-agent systems to shared transport services: A review*. In Proceedings of Iberian Conference on Information Systems and Technologies (CISTI), pages 1–6. IEEE, 2018.
- [31] Davies, D. L.; Bouldin, D. W. (1979). *A Cluster Separation Measure*. IEEE Transactions on Pattern Analysis and Machine Intelligence. PAMI-1 (2): 224–227.
- [32] De Nijs, F.; Theocharous, G.; Vlassis, N.; De Weerd, M. M.; Spaan, M. T. J. *Capacity-aware sequential recommendations*. In Proceedings of International Conference on Autonomous Agents and MultiAgent Systems (AAMAS), 2018, p. 416–424.
- [33] Devogele, T.; Etienne, L.; Esnault, M.; Lardy, F. *Optimized Discrete Fréchet Distance between trajectories*. In Proceedings of the 6th ACM SIGSPATIAL Workshop on Analytics for Big Geospatial Data; Association for Computing Machinery: New York, NY, USA, 2017; BigSpatial’17, p. 11–19. doi:10.1145/3150919.3150924.
- [34] Earl, R. W. *Book reviews : Models of man: Social and rational* by Herbert A. Simon. New York: John Wiley & sons, 1957. pp. xiv 287. Educational and Psychological Measurement, vol. 21, no. 1, pp. 244–246, 1961.
- [35] Eiter, T.; Mannila, H. *Computing discrete Fréchet distance*. Technical report, Technische Universität Wien, 1994.
- [36] Ester, M.; Kriegel, H.-P.; Sander, J.; Xu, X. et al. *A density-based algorithm for discovering clusters in large spatial databases with noise*. In Proceedings of International Conference on Knowledge Discovery and Data Mining, pages 226–231, 1996.

- [37] Fano, A. E. *Shopper's eye: using location-based filtering for a shopping agent in the physical world*. In Agents, volume 98, pages 416–421, 1998.
- [38] Forgy, Edward W. (1965). *Cluster analysis of multivariate data: efficiency versus interpretability of classifications*. Biometrics. 21 (3): 768–769.
- [39] Fréchet, M. *Sur quelques points du calcul fonctionnel*. Rendiconti del Circolo Matematico di Palermo, 22:1–74, 1906.
- [40] Gambs, S.; Killijian, M.-O.; del Prado Cortez, M. N. *Next place prediction using mobility markov chains*. In Proceedings of the First Workshop on Measurement, Privacy, and Mobility, page 3. ACM, 2012.
- [41] Gedik, B.; Liu, L. *Location privacy in mobile systems: A personalized anonymization model*. In Proceedings of IEEE International Conference on Distributed Computing Systems (ICDCS), pages 620–629, 2005.
- [42] GeoLife, Microsoft Research Asia, <http://www.microsoft.com/en-us/download/details.aspx?id=52367>
- [43] Giunta, R.; Pappalardo, G.; Tramontana, E. *Aspects and annotations for controlling the roles application classes play for design patterns*. In Proc. of IEEE Asia Pacific Software Engineering Conference (APSEC), pages 306–314, Ho Chi Minh, Vietnam, December 2011.
- [44] Goldberg, D.; Nichols, D.; Oki, B. M. ; Terry, D. *Using collaborative filtering to weave an information tapestry*. Communications of the ACM, 35(12):61–71, 1992.
- [45] Han, J.; Kamber, M.; Pei, J. *Data Mining: Concepts and Techniques*, 3rd ed.; Morgan Kaufmann Publishers Inc.: San Francisco, CA, USA, 2011.
- [46] Hernández, J. Z.; Ossowski, S.; Garcia-Serrano, A. *Multiagent architectures for intelligent traffic management systems*. Transportation Research Part C: Emerging Technologies, 10(5-6):473–506, 2002.
- [47] Hosseinpoor, A.H.; Ali Abbaspour, R.; Claramunt, C. *A Geometric Framework for Detection of Critical Points in a Trajectory Using Convex Hulls*. ISPRS Int. J. Geo-Inf 2018, 7, 14. doi:10.3390/ijgi7010014.
- [48] Hsieh, H.-P.; Li, C.-T.; Lin, S.-D. *Exploiting large-scale check-in data to recommend time-sensitive routes*. In Proceedings of the ACM SIGKDD International Workshop on Urban Computing, pages 55–62. ACM, 2012.
- [49] Hsueh, Y.-L.; Huang, H.-M. *Personalized itinerary recommendation with time constraints using gps datasets*. Knowledge and Information Systems, 60(1):523–544, 2019.
- [50] Hu, Y.; Gao, S.; Janowicz, K.; Yu, B.; Li, W.; Prasad, S. *Extracting and understanding urban areas of interest using geotagged photos*. Computers, Environment and Urban Systems, 54:240–254, 2015.

- [51] Hu, M.; Liu, B. *Mining and summarizing customer reviews*. In Proceedings of ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD-2004). 2004.
- [52] Hwang, J.-R.; Kang, H.-Y.; Li, K.-J. *Spatio-temporal similarity analysis between trajectories on road networks*. In Proceedings of International Conference on Conceptual Modeling, pages 280–289. Springer, 2005.
- [53] Inman, J. *Navigation and Nautical Astronomy, for the Use of British Seamen*. F. & J. Rivington: London, UK, 1849.
- [54] Isaacman, S.; Becker, R.; Cáceres, R.; Kobourov, S.; Martonosi, M.; Rowland, J.; Varshavsky, A. *Identifying important places in people’s lives from cellular network data*. In Proceedings of International Conference on Pervasive Computing, pages 133–151. Springer, 2011.
- [55] Jiang, T.; Wang, H. J.; Hu, Y.-C. *Preserving location privacy in wireless lans*. In Proceedings of International conference on Mobile systems, applications and services, pages 246–257, 2007.
- [56] Jiménez-Bravo, D. M.; Pérez-Marcos, J.; De la Iglesia, D. H.; Villarrubia González, G.; De Paz, J. F. *Multi-agent recommendation system for electrical energy optimization and cost saving in smart homes*. *Energies*, 12(7):1317, 2019.
- [57] Kalnis, P.; Ghinita, G.; Mouratidis, K.; Papadias, D. *Preventing location-based identity inference in anonymous spatial queries*. *IEEE transactions on knowledge and data engineering*, 19(12):1719–1733, 2007.
- [58] Lau, B. P. L.; Hasala, M. S.; Kadaba, V. S.; Thirunavukarasu, B.; Yuen, C.; Yuen, B.; Nayak, R. *Extracting point of interest and classifying environment for low sampling crowd sensing smartphone sensor data*. In Proceedings of IEEE International Conference on Pervasive Computing and Communications Workshops (PerCom Workshops), pages 201–206, 2017.
- [59] Lee, J.; Han, J.; Li, X. *A unifying framework of mining trajectory patterns of various temporal tightness*, *IEEE Transactions on Knowledge and Data Engineering*, vol. 27, no. 6, pp. 1478–1490, 2015.
- [60] Liu, B. *Sentiment Analysis and Subjectivity*. Handbook of Natural Language Processing, Second Edition, N. Indurkha and F.J. Damerau, Eds. 2010.
- [61] Liu, B. *Sentiment analysis and opinion mining*. *Synthesis lectures on human language technologies*, 5(1):1–167, 2012.
- [62] Liu, B.; Fu, Y.; Yao, Z.; Xiong, H. *Learning geographical preferences for point-of-interest recommendation*. In Proceedings of ACM SIGKDD international conference on Knowledge discovery and data mining, pages 1043–1051, 2013.

- [63] Liu, Q.; Gao, L.; Lou, P. *Resource management based on multi-agent technology for cloud manufacturing*, in 2011 International Conference on Electronics, Communications and Control (ICECC), pp. 2821–2824, IEEE, 2011.
- [64] Lops, P.; De Gemmis, M.; Semeraro, G. *Content-based recommender systems: State of the art and trends*. In Recommender systems handbook, pages 73–105. Springer, 2011.
- [65] Mamei, M.; Colonna, M.; Galassi, M. *Automatic identification of relevant places from cellular network data*. Pervasive and Mobile Computing, 2016, <https://doi.org/10.1016/j.pmcj.2016.01.009>
- [66] Mathew, W.; Raposo, R.; Martins, B. *Predicting future locations with hidden markov models*. In Proceedings of ACM conference on ubiquitous computing, pages 911–918, 2012.
- [67] Melville, P.; Sindhvani, V. *Recommender systems.*, Encyclopedia of machine learning, vol. 1, pp. 829–838, 2010.
- [68] Miyahara, K. ; Pazzani, M. J. *Collaborative filtering with the simple bayesian classifier*. In Proceedings of Pacific Rim International conference on artificial intelligence, pages 679–689. Springer, 2000.
- [69] Mohamed, F.; Chi-Yin, C.; Walid, G. *The new casper: a privacy-aware location-based database server*. In Proceedings of International Conference on Data Engineering, pages 1499–1500, 2007.
- [70] Narayanan, A.; Thiagarajan, N.; Lakhani, M.; Hamburg, M.; Boneh, D. et al. *Location privacy via private proximity testing*. In NDSS, volume 11, 2011.
- [71] Nasukawa, T.; Jeonghee Y. *Sentiment analysis: Capturing favorability using natural language processing*. In Proceedings of the KCAP-03, 2nd Intl. Conf. on Knowledge Capture. 2003.
- [72] Papadopoulos, S.; Bakiras, S.; Papadias, D. *Nearest neighbor search with strong location privacy*. volume 3, pages 619– 629. VLDB Endowment, 2010.
- [73] Pérez-Torres, R. ; Torres-Huitzil, C. ; Galeana-Zapién, H. *Full on device stay points detection in smartphones for location-based mobile applications*. Sensors, 16(10):1693, 2016.
- [74] Piorkowski, M.; Sarafijanovic-Djukic, N.; Grossglauser, M. *CRAWDAD dataset epfl/mobility(v. 2009-02-24)*, 2009. doi:10.15783/C7J010.
- [75] Polat H.; Du, W. *Svd-based collaborative filtering with privacy*. In Proceedings of ACM Symposium on Applied Computing (SAC), pages 791–795, 2005.
- [76] Quercia, D.; Schifanella, R.; Aiello, L. M. *The shortest path to happiness: Recommending beautiful, quiet, and happy routes in the city*. In Proceedings of ACM conference on Hypertext and social media, pages 116–125, 2014.

- [77] Reiter, M. K.; Rubin, A. D. *Crowds: Anonymity for web transactions*. ACM transactions on information and system security (TISSEC), 1(1):66–92, 1998.
- [78] Rolim, V.B.; Silva, M.R.; Flamarion, F.C.; Filho, C.C.G.F.; Braz, F.J. *A Method for Identifying Patterns of Movement of Trajectory Sets by Using the Frequency Distribution of Points*. GEOProcessing 2017 : The Ninth International Conference on Advanced Geographic Information Systems, Applications, and Services, 2017. ISBN: 978-1-61208-539-5.
- [79] Rousseeuw, P. J. *Silhouettes: a Graphical Aid to the Interpretation and Validation of Cluster Analysis*. Computational and Applied Mathematics. 20: 53–65, (1987).
- [80] Sakamane, P.; Phithakkitnukoon, S.; Smoreda, Z.; Ratti, C. *Methods for Inferring Route Choice of Commuting Trip From Mobile Phone Network Data*. ISPRS Int. J. Geo-Inf. 2020, 9, 306.
- [81] Salton, G.; Buckley, C. *Term-weighting approaches in automatic text retrieval*. In Information Processing & Management, 24(5), 513–523, 1988.
- [82] Sarwar, B. M.; Konstan, J. A.; Borchers, A.; Herlocker, J.; Miller, B.; Riedl, J. *Using filtering agents to improve prediction quality in the groupLens research collaborative filtering system*. In Proceedings of the ACM Conference on Computer Supported Cooperative Work (CSCW), GroupLens Research Collaborative Filtering System, 1998.
- [83] Schafer, J. B.; Frankowski, D.; Herlocker, J.; Sen, S. *Collaborative filtering recommender systems*. In The adaptive web, pages 291–324. Springer, 2007.
- [84] Sebastia, L.; Giret, A.; Garcia, I. *A multi agent architecture for tourism recommendation*. In Trends in Practical Applications of Agents and Multiagent Systems, pages 547–554. Springer, 2010.
- [85] Shen, W.; Norrie, D. *Facilitators, mediators or autonomous agents*. In Proceedings of the Second International Workshop on CSCW in Design, Bangkok, Thailand, pp. 119–124, 1997.
- [86] Sinnott, R.W. *Virtues of the Haversine*. Sky Telesc. 68, 159, 1984.
- [87] Sycara, K. P. *Multiagent systems*, AI Magazine, vol. 19, p. 79, Jun. 1998.
- [88] Su, X.; Khoshgoftaar, T. M. *Collaborative filtering for multi-class data using bayesian networks*. International Journal on Artificial Intelligence Tools, 17(01):71–85, 2008.
- [89] Suneetha, K.; Krishnamoorti, R. *Web Log Mining using Improved Version of Apriori Algorithm*. International Journal of Computer Application 2011, 29, 23–27. doi:10.5120/3569-4923.
- [90] Syverson, P. F.; Goldschlag, D. M.; Reed, M. G. *Anonymous connections and onion routing*. In Proceedings of IEEE Symposium on Security and Privacy, pages 44–54, 1997.

- [91] Taxi and Truck datasets, <http://mashuai.buaa.edu.cn/traj.html>
- [92] Tramontana, E. *Minimising changes when refactoring applications to run multiple threads*. In Proceedings of IEEE Asia-Pacific Software Engineering Conference (APSEC), Nara, Japan, 2018, pp. 713–714.
- [93] Tramontana, E.; Verga, G. *Get spatio-temporal flows from gps data*. In Proceedings of IEEE International Conference on Smart Computing (SMARTCOMP), pages 282–284, 2018.
- [94] Trasarti, R.; Rinzivillo, S.; Pinelli, F.; Nanni, M.; Monreale, A.; Renso, C.; Pedreschi, D.; Giannotti, F. *Exploring real mobility data with m-atlas*, in Machine Learning and Knowledge Discovery in Databases (J. L. Balczar, F. Bonchi, A. Gionis, and M. Sebag, eds.), pp. 624–627, Springer, 2010.
- [95] Ungar, L. H.; Foster, D. P. *Clustering methods for collaborative filtering*. In Proceedings of AAAI workshop on recommendation systems, volume 1, pages 114–129, 1998.
- [96] Wang, Z.; Lu, M.; Yuan, X.; Zhang, J.; Wetering, H. v. d. *Visual traffic jam analysis based on trajectory data*, IEEE Transactions on Visualization and Computer Graphics, vol. 19, no. 12, pp. 2159–2168, 2013
- [97] Wooldridge, M.; Jennings, N. R. *Intelligent agents and multi-agent systems. i*, Applied artificial intelligence, vol. 9, no. 4, 1995.
- [98] Xu R.; Wunsch D. *Survey of clustering algorithms*. Trans. Neur. Netw. 16, 3 (May 2005), 645–678.
- [99] Xuan, P.; Lesser, V. *Multi-agent policies: from centralized ones to decentralized ones*. In Proceedings of the first international joint conference on Autonomous agents and multiagent systems: part 3, pages 1098–1105. ACM, 2002.
- [100] Xuan, P.; Lesser, V.; Zilberstein, S. *Communication decisions in multi-agent cooperation: Model and experiments*. In Proceedings of the fifth international conference on Autonomous agents, pages 616–623. ACM, 2001.
- [101] Yoon, H.; Zheng, Y.; Xie, X.; Woo W. *Smart itinerary recommendation based on user-generated gps trajectories*. In Proceedings of International Conference on Ubiquitous Intelligence and Computing, pages 19–34. Springer, 2010.
- [102] Yuan, Q.; Cong, G; Ma, Z.; Sun, A.; Thalmann, N. M. *Time-aware point-of-interest recommendation*. In Proceedings of the international ACM SIGIR conference on Research.
- [103] Zheng, Y. *Trajectory data mining: An overview*. ACM Trans. Intell. Syst. Technol. 2015, 6, Article 29. doi:10.1145/2743025.

- [104] Zheng, Y.; Li, Q.; Chen, Y.; Xie, X.; Ma, W.Y. *Understanding Mobility Based on GPS Data*. In Proceedings of ACM conference on Ubiquitous Computing (UbiComp 2008), Seoul, Korea. Association for Computing Machinery, 2008, pp. 312–321. doi:10.1145/1409635.1409677.
- [105] Zheng, Y.; Xie, X.; Ma, W.Y. *GeoLife: A Collaborative Social Networking Service among user, location and trajectory*. IEEE Data(base) Engineering Bulletin (2010).
- [106] Zheng, Y.; Zhang, L.; Xie, X.; Ma, W.Y. *Mining interesting locations and travel sequences from GPS trajectories*. In Proceedings of International conference on World Wild Web (WWW 2009), Madrid, Spain. Association for Computing Machinery, 2009, pp. 791–800. doi:10.1145/1526709.1526816.
- [107] Zou, Y.; Chen, Y.; He, J.; Pang, G.; Zhang, K. *4D Time Density of Trajectories: Discovering Spatiotemporal Patterns in Movement Data*. ISPRS Int. J. Geo-Inf 2018, 7, 212. doi:10.3390/ijgi7060212.
- [108] Zygouras, N.; Gunopulos, D. *Corridor Learning Using Individual Trajectories*. 19th IEEE International Conference on Mobile Data Management (MDM). IEEE, 2018, pp. 155–160. doi:10.1109/MDM.2018.00032.