

DATA SOURCING TECHNIQUES FOR PERSONALIZED LOCATION BASED TOURIST SPOT RECOMMENDER SYSTEM

SUNITA TIWARI



**AMARNATH & SHASHI KHOSLA SCHOOL OF
INFORMATION TECHNOLOGY
INDIAN INSTITUTE OF TECHNOLOGY DELHI
JANUARY 2018**

©Indian Institute of Technology Delhi (IITD), New Delhi, 2018

DATA SOURCING TECHNIQUES FOR PERSONALIZED LOCATION BASED TOURIST SPOT RECOMMENDER SYSTEM

by
SUNITA TIWARI

Amarnath and Shashi Khosla School of Information Technology

Submitted
in fulfillment of the requirements of the degree of Doctor of Philosophy

to the



INDIAN INSTITUTE OF TECHNOLOGY DELHI

JANUARY 2018

Certificate

The thesis entitled “**Data Sourcing Techniques for Personalized Location Based Tourist Spot Recommender System**” being submitted by **Ms. Sunita Tiwari** to the Indian Institute of Technology Delhi for award of the degree of **Doctor of Philosophy** is a record of original bonafide research work carried out by her. She has worked under my guidance and supervision, and has fulfilled the requirements for the submission of this thesis, which has attained the standard required for a Ph.D. degree of this institute.

The results presented in this thesis have not been submitted elsewhere for the award of any other degree or diploma.

Dr. Saroj Kaushik

Professor

Department of Computer Science and Engineering

Indian Institute of Technology Delhi

New Delhi

Acknowledgements

I would thank the Almighty for giving me the courage and determination to pursue the higher education in the field of Computer Science and Engineering. I would like to tender my heartfelt thankfulness to my supervisor **Prof. Saroj Kaushik** for her priceless supervision, continuous inspiration and enthusiastic support in every stage of the work. Her deep insight into problems, imaginative ideas, technical guidance and sufficient encouragement was a constant source of inspiration and motivation. Her critical comments on my technical writing, eye for detail and persistent encouragement strongly motivated me to strive hard to achieve the high targets set by her in completion of this thesis. Her religious, humanitarian and positive outlook towards the life has been a great source of inspiration even in the non-academic side of life.

I am extremely grateful to the members of my Student Research Committee (SRC) – Prof. M P Gupta, Prof. Aditeshwar Seth, Prof. Punam Bedi, who have been very helpful by offering suggestions and advice. I thank my Ph.D. colleagues, especially, Priti Jagwani, Shivendra Prasad Tiwari, and Kuntal Dey for their help on reviewing my work and offering useful suggestions. I also extend the sincerest gratitude to the summer interns at IIT Delhi Chhavi and Akash and most importantly my students who contributed in conducting experiments by registering in the proposed system and providing their valuable feedback. My sincere gratitude to my colleagues at work place for providing the unconditional support.

My Grandparents, my parents Mr. O. P Sharma and Mrs. Girija Sharma and my husband Sourabh Tiwari and other family members showed immense patience and provided me great support during the course of my work. Their patience, sacrifice, inspiration and moral support are of special mention. I dedicated my thesis to them.

Sunita Tiwari

Abstract

Research in recommender system has matured in the last two decades. The emphasis in recent research is more focused on Location Based Services (LBS) due to the advancement of wireless communication devices and location acquisition technologies like Global Positioning Systems (GPS). The mobile and handheld devices are the primary mode of communication and information access in this era. The research issues in the field of recommender system for mobile devices are more challenging as compared to traditional web based recommender systems. The techniques applied to the web based recommender system cannot be straight away used for such devices because of several reasons such as dissemination of relevant information in small screen as the size of devices are small, dissemination of information in real time as the users are mobile , etc.

The proposed research work aims to design and develop various techniques required for personalized location based recommender systems. Such systems will increase the usability of mobile device to many folds during their journey. Tourism is one of the largest industries in the world. We restrict our focus to tourism domain in order to facilitate the users with the interesting tourist spots, enriched tour details and so on.

The research work broadly addresses the following challenges in the field of the location based tourist spot recommender systems.

Firstly, it addresses the problem of automatic discovery of new popular tourist spots in a geographical region and semantic annotation of point of interests. Location semantics are gathered by crawling the web in a focused way using domain ontology. Our experiment shows that 20-23% of the locations retrieved from the hierarchical graph based approach presented in (Zheng Y. a.-Y., 2009) and relational algebra based approach presented in (Khetarpaul, 2011) are not interesting location from tourist point of view.

Secondly, the problem of automatic discovery of user preferences and personalized location based recommendation is discussed. Past travel histories and genetic algorithm based approach is exploited for the personalization problem of the tourist spot recommender system. The proposed approach is compared with Matrix Factorization based approach presented in (Berjani, 2011) and the proposed

approach shows an improvement in average RMSE of (approx.) 4.63%. The recommendation accuracy is 94% (approx.).

The third problem is to propose, the semantic user similarity measure for recommender system. The proposed approach for user similarity is based on semantic stay trajectories using Earth Mover's Distance (EMD). The proposed approach performs better than the popular user similarity measures such as Pearson's correlation, Jacard and Dice. It shows an average percentage improvement of 10.7% in RMSE and 5.73% in MAE as compared to the above mentioned approaches.

Finally, the problem of Information enrichment of recommender system is addressed. This technique used the concept of location based crowdsourcing with the aim of improving the quality of the recommendation. Here, fuzzy inference system is used to convert the contextual information obtained from the crowd to an appropriate rank of the recommended spot. The experimental evaluation demonstrates the better satisfaction level of users with enriched information.

सार

अनुशंसात्मक प्रणाली (Recommender System) में शोध पिछले दो दशकों में परिपक्व हो गया है। हाल ही में शोध में बेतार संचार उपकरणों और स्थान अधिग्रहण तकनीकों जैसे वैश्विक स्थिति प्रणाली(GPS) की उन्नति के कारण स्थान आधारित सेवाओं(LBS) पर अधिक ध्यान केंद्रित किया गया है। इस युग में मोबाइल और हाथ में धारण करने योग्य उपकरण संचार और सूचना प्राप्ति के प्राथमिक तरीके हैं । मोबाइल/चल उपकरणों के क्षेत्र में अनुसंधान के मुद्दे पारम्परिक वेब आधारित अनुशंसात्मक तकनीक की तुलना में अधिक चुनौतीपूर्ण हैं । इन कारणों में छोटे/लघु स्क्रीन वाले उपकरण और गतिशील उपयोगकर्ता के पास वास्तविक समय में जानकारी को प्रसार कर पाना सम्मिलित है ।

प्रस्तावित अनुसंधान कार्य का उद्देश्य व्यक्तिगत स्थान आधारित अनुशंसात्मक प्रणाली हेतु विभिन्न तकनीकों का प्रारूपण और विकसित करना है । इस प्रकार की प्रणाली मोबाइल उपकरण की प्रयोज्यता को उसकी विकास यात्रा के दौरान कई गुणा तक बड़ा देगी । पर्यटन दुनिया के सबसे वृहत्तम उद्योगों में से एक है। हम पर्यटन के क्षेत्र में उपयोगकर्ता को दिलचस्प पर्यटन स्थलों की जानकारी देने और पर्यटन से संबंधित विवरण से उन्हें समृद्ध करने के लिए शोध को केंद्रित करते हैं ।

शोध कार्य मुख्य तौर पर स्थान आधारित पर्यटन स्थल अनुशंसात्मक प्रणाली के क्षेत्र में अधोलिखित चुनौतियों को सम्बोधित करता है ।

सबसे पहले, यह एक भौगोलिक क्षेत्र में नए लोकप्रिय पर्यटन स्थलों की खोज और पर्यटन स्थलों की अर्थिक टिप्पणियों की स्वतः खोज की समस्या को संबोधित करता है। स्थानीय शब्दार्थों को domain ontology का उपयोग करते हुए एक केंद्रित पद्धति से वेब को क्रॉल करके इकट्ठा किया जाता है। हमारा प्रयोग बताता है कि (Zheng Y. a.-Y., 2009) में प्रस्तुत पदानुक्रमित ग्राफ आधारित दृष्टिकोण से और (Khetarpaul, 2011) में प्रस्तुत रिलेशनल बीजगणित पद्धति से प्राप्त 20 से 23% स्थान पर्यटन के लिए दिलचस्प नहीं है ।

दूसरी बात , उपयोगकर्ता की प्राथमिकताओं की स्वतः खोज और व्यक्तिगत स्थान आधारित अनुशंसा की स्थापना की समस्या पर चर्चा की गई है। विगत यात्रा विवरण और आनुवंशिक विधि (algorithm) पर आधारित तकनीक का उपयोग पर्यटन स्थल अनुशंसा करने वाली प्रणाली की निजीकरण समस्या के समाधान के लिए किया जाता है।

प्रस्तावित पद्धति की तुलना मैट्रिक्स फैक्टरिजेशन आधारित पद्धति (Berjani, 2011) से की गई है और प्रस्तावित विधि में औसत RMSE (लगभग) 4.63% का सुधार दर्शाता है। अनुशासनात्मक परिशुद्धता लगभग 94% है।

तीसरी बात, अनुशासनात्मक तकनीक हेतु सार्थक उपयोगकर्ता समानता उपाय प्रस्तावित करना है।

उपयोगकर्ता सादृश्य हेतु प्रस्तावित प्रणाली सिमेंटिक ट्राजेक्टोरिज(semantic trajectories) और Earth Mover's Distance (EMD) तकनीक का उपयोग करती है। प्रस्तावित पद्धति लोकप्रिय उपयोगकर्ता समानता उपायों जैसे Pearson's Correlation Coefficient method, Jacard Method तथा Dice Method आदि से बेहतर प्रदर्शन कर रही है। उपर्युक्त विधियों की तुलना में प्रस्तावित विधि औसत रूट मीन स्क्वायर एरर (RMSE) में 10.7% और औसत मीन स्क्वायर एरर (MAE) में 5.73% की औसत प्रतिशत का सुधार दर्शाती है।

और अततोगत्वा अनुशासनात्मक प्रणाली की सुचना संवर्धन (Information Enrichment) की समस्या को समबोधित किया जाता है। इस तकनीक ने अनुशासा की गुणवत्ता में सुधार लाने के उद्देश्य से स्थान आधारित क्रोवडसोर्सिंग (crowdsourcing) की अवधारणा का प्रयोग किया। यहां एक स्थल पर उपलब्ध भीड़ से प्राप्त प्रासंगिक पैरामीटर फजी निष्कर्ष प्रणाली(Fuzzy inference system) द्वारा उस पर्यटन स्थल के उपयुक्त स्तर(rank) प्राप्त करने के लिए उपयोग किया जाता है। प्रायोगिक मूल्यांकन सूचना समृद्धीकरण (information enrichment) वाले अनुशासनात्मक प्रणाली के उपयोगकर्ताओं के बेहतर संतुष्टि स्तर को दर्शाते हैं।

Table of Contents

<i>Certificate</i>	<i>i</i>
<i>Acknowledgements</i>	<i>iii</i>
<i>Abstract</i>	<i>v</i>
<i>संक्षेप</i>	<i>vii</i>
<i>Table of Contents</i>	<i>ix</i>
<i>List of Figures</i>	<i>xiii</i>
<i>List of Tables</i>	<i>xv</i>
<i>Acronyms</i>	<i>xvi</i>
Chapter 1 : Introduction	1
1.1. Location Based Recommender System Basics	1
1.2. Traditional Recommender Systems	3
1.2.1. Collaborative Filtering Recommender Systems (CFRS).....	4
1.2.2. Content Based Recommender Systems (CBRS)	5
1.2.3. Knowledge Based Recommender Systems (KBRS)	5
1.2.4. Hybrid Recommender Systems (HRS).....	6
1.3. Context Aware Recommender System (CARS)	6
1.4. Location Based Recommender Systems	7
1.4.1. Research Challenges in Location Based Recommender Systems	8
1.4.2. Broad Research Challenges Addressed.....	11
1.4.3. Related Work	11
1.5. Motivation	14
1.6. Research Problem Statement	15
1.7. Proposed Research Work Preliminaries	16
1.7.1. GPS Log and Related Terms	16
1.7.2. Dataset Used.....	18
1.7.3. Point of Interests (POI)	19
1.8. Thesis Outline	20
1.8.1. Mining Popular Tourist Spots	20

1.8.2. Implicit User Preference Discovery	21
1.8.3. Evolving user's preference for unvisited spot	22
1.8.4. Semantic User Similarities.....	23
1.8.5. Information Enrichment.....	23
Chapter 2 : Location Semantic based Mining of Popular Tourist Spots from GPS	
Trajectories	25
2.1. Introduction	25
2.2. Related Work	27
2.3. Problem Definition and Challenges	29
2.4. Proposed Solution.....	32
2.4.1. Stay Point Generation	33
2.4.2. Clustering of Stay Points into Popular Spots	37
2.4.3. Semantic annotation of spots.	41
2.5. Implementation and Evaluation of the Proposed Approach	50
2.6. Conclusion	57
Chapter 3 : Evolving Rating of Unvisited Spots Using Genetic Algorithm	59
3.1. Introduction	59
3.2. Related Work	61
3.3. Problem Definition and Challenges	66
3.4. Proposed Methodology.....	67
3.4.1. Popular Spot Mining.....	69
3.4.2. Implicit Preference Discovery	69
3.4.3. Unknown Rating Predictor	72
3.4.4. Recommendation Engine	78
3.4.5. Mobile User	79
3.5. Implementations and Results.....	80
3.5.1. Experimental Settings and Results	80
3.5.2. Validation of Predicted Values	91
3.6. Conclusion	95
Chapter 4 : Semantic Based User Similarity Using Past Travel Histories	97
4.1. Introduction	97

4.2. Background and Related Work.....	98
4.2.1. Semantic Trajectories	99
4.2.2. User Similarity.....	100
4.2.3. Earth Mover Distance	102
4.3. Problem Definition.....	103
4.4. Proposed System Design	104
4.4.1. Stay Point Generation.....	104
4.4.2. Semantic annotation of stay points.....	105
4.4.3. Earth Mover Distance for Semantic Stay Sequences.....	106
4.5. Implementation and Results	111
4.5.1. Scenario-1 (RMSE)	111
4.5.2. Scenario-2(MAE)	113
4.5.3. Scenario-3 (Loss of Information)	113
4.6. Evaluation	115
4.6.1. Ground Truth and Evaluation	115
4.7. Conclusion	115
<i>Chapter 5 : Information Enrichment for Tourist Spot Recommender System Using Location Aware Crowdsourcing</i>	<i>117</i>
5.1. Introduction	117
5.2. Background and Related Work.....	119
5.2.1. Crowdsourcing.....	119
5.2.2. Information Enrichment of Recommender Systems	122
5.2.3. Fuzzy Inference System	123
5.3. Problem Definition and Challenges	124
5.4. Proposed Solutions	124
5.4.1. Basic Information Enrichment System.....	125
5.4.2. Fuzzy Information Enrichment System	130
5.5. Implementation and Evaluation of the Proposed Approach	135
5.5.1. Implementation	135
5.5.2. Evaluation	141
5.6. Content Verification Services	143
5.7. Conclusion	145

Chapter 6 : Conclusion and Future Directions.....	147
References.....	151
Appendix A: Sample Code & Experimental Results [Chapter 2]	169
Some Results of Proposed Approach.....	169
Sample Code for Crawler	176
Sample Code for Computing Stay Points of All Users	179
Appendix B: Sample Code & Experimental Results [Chapter 3]	183
Matrix Factorization for RS	183
Sample MF code.....	183
Some Results of Proposed GA Based Approach.....	184
Part of GA code	185
Appendix C: Sample Code & Experimental Results [Chapter 4]	189
Some Results of Proposed EMD based Approach.....	189
Appendix D: Sample Code & Experimental Results [Chapter 5]	191
Sample Code for Crowd Platform	191
PhP code for registration	191
PhP code for collecting responses from crowd.....	192
Part of the FIS used for experimentation.....	194
Publications out of Thesis	197
Brief Resume	199

List of Figures

Figure 1.1 Location Based System Architecture	1
Figure 1.2 Personalized Recommendations overview.....	3
Figure 1.3 CARS Overview	7
Figure 1.4 General Architecture of LBRS.....	8
Figure 1.5 GPS logs	16
Figure 1.6 GPS Trajectory	17
Figure 1.7 Architecture of pattern analysis	18
Figure 1.8 Overall design of proposed work.....	20
Figure 2.1 GPS trajectory and stay points of a user U.	31
Figure 2.2 Proposed Solution Design.....	33
Figure 2.3 Sample trajectory dataset of GeoLife dataset.....	34
Figure 2.4 Preprocessed trajectories of a user	34
Figure 2.5 Some of the stay points generated by multiple users are plotted on the Map	37
Figure 2.6 Example of core, border and noise point	39
Figure 2.7 Results of DBSCAN for $\epsilon = 200$ meters and $\text{Min_points} = 6$	40
Figure 2.8 Design of semantic annotation module	41
Figure 2.9 Sample Results of reverse geocoding using OpenStreetMap	43
Figure 2.10 Working of a basic crawler	44
Figure 2.11 Working of Ontology based focused crawler.....	45
Figure 2.12 Retrieval of Frequent Terms using Tagcrowd (Tagcrowd, 2012).....	46
Figure.2.13 Part Ontology for tourism domain.	48
Figure 2.14 Frequency distribution of Spots	53
Figure 2.15 Ontology Based Focused Crawler Interface	54
Figure 3.1 Proposed Solution Design.....	68
Figure 3.2 User-spot Matrix	70
Figure 3.3 Part of User Spot Matrix	70
Figure 3.4 Example Chromosome.....	73
Figure 3.5 Example of block uniform crossover	77
Figure 3.6 Example of mutation operator	77
Figure 3.7 Plot of percentage improvement of RMSE for DS-1 to DS-5.....	84
Figure 3.8 Comparisons of RMSE over five runs	85
Figure 3.9 RMSE comparisons for different runs	85

Figure 3.10 Comparison of results for GA and MF based approaches over varying number of users	86
Figure 3.11 Comparison of MAE for GA based and MF based approach.....	88
Figure 3.12 Average RSME of GA based and MF based approach for sparse data	91
Figure 3.13 Example of reverse experiments	92
Figure 3.14 Percentage error for each data samples when experiments are conducted in reverse direction....	92
Figure 3.15 % error for varying number of users for experiments are conducted in reverse direction.....	93
Figure 4.1 User trajectories.....	99
Figure 4.2 Semantic trajectory.....	100
Figure 4.3 Proposed System Design.....	105
Figure 4.4 Semantic stay sequences	106
Figure 4.5 Example Flow graph	109
Figure 4.6 Average percentage improvement of proposed approach.....	112
Figure 4.7 % improvement in Average MAE of proposed method	113
Figure 4.8 Comparison of stay based and trajectory based approach	114
Figure 5.1 Crowdsourcing Architecture (Howe, 2006)	120
Figure 5.2 Proposed Information Enrichment System Design.....	125
Figure 5.3 Sample Questionnaire Form.....	127
Figure 5.4 Proposed Fuzzy Information Enrichment System Design.....	131
Figure 5.5 Membership function of input variable “weather”	132
Figure 5.6 Membership function of input variable “traffic”	133
Figure 5.7 Membership function of input variable “crowdedness”	133
Figure 5.8 Membership function of input variable “security”	133
Figure 5.9 Membership function of input variable “Personalized Rank”	134
Figure 5.10 Membership function of output variable “Rating”	134
Figure 5.11 Sample Fuzzy if-then Rules	134
Figure 5.12 Modules of proposed prototype systems	136
Figure 5.13 Navigation flow of proposed prototype system.....	137
Figure 5.14 some of the tourist spots considered in experiments	138
Figure 5.15 Consolidate feedback for proposed information enrichment system.....	142
Figure 5.16 TSRS Vs TSRS with IES	142
Figure 5.17 Content Uploading GUI	144
Figure 5.18 Chatting	144
Figure 5.19 Audio Calling	145

List of Tables

Table 1.1 Data point after pre-processing	18
Table 2.1 Results obtained for different values of time thresholds.	36
Table 2.2 Term Weight in Tourism Ontology	47
Table 2.3 Example tags associated with some of the spots	50
Table 2.4 Categories of POI in HERE Map Beijing POI dataset	51
Table 2.5 Top twenty spots based on user frequency.....	52
Table 2.6 Results after semantic annotation of spots from tourism perspective	55
Table 2.7 Here Maps Ratings of Top 10 Relevant tourist spots.....	56
Table 3.1 Average RMSE for each dataset and varying population size	81
Table 3.2 Average RMSE for GA and MF methods	82
Table 3.3 Experiments over all five dataset using GA and MF based approaches.....	83
Table 3.4 Confusion Matrix	89
Table 3.5 Comparison of Avg. RMSE for MovieLens Dataset.....	94
Table 4.1 Tag similarity example.....	109
Table 4.2 Average RMSE of different similarity measures	112
Table 5.1 Sample Result Generated by TSRS.....	138
Table 5.2 Additional information collected by crowdsourcing.....	139
Table 5.3 Integrated Fuzzy Information Enrichment	140

Acronyms

AGPS	Assisted Global Positioning System
API	Application Program Interface
CARS	Context Aware Recommender System
CBRS	Content Based Recommender System
CF	Collaborative Filtering
CFRS	Collaborative Filtering Recommender Systems
DBSCAN	Density Based Spatial Clustering of Applications with Noise
DS	Dataset
EMD	Earth Mover's Distance
GA	Genetic Algorithm
GIS	Geographic Information System
GPS	Global Positioning System
GSM	Global System for Mobile Communication
HGSM	Hierarchical Graph Based Similarity Measurement
HRS	Hybrid Recommender Systems
IES	Information Enrichment System
KBRS	Knowledge Based Recommender Systems
LBRS	Location Based Recommender Systems
LBS	Location Based Services
LBSN	Location Based Social Network
LCARS	location and context aware recommender system

MAE	Mean Absolute Error
MF	Matrix Factorization
PHP	Hypertext Preprocessor
POI	Point of Interest
RFID	Radio-frequency identification
RMSE	Root Mean Squared Error
RS	Recommender systems
SDK	Software Development Kit
TSRS	Tourist Spot Recommender Systems
URL	Uniform Resource Locator
VVIP	Very Very Important Person
WWW	World Wide Web
XML	Extensible Markup Language
PND	Personal Navigation Device