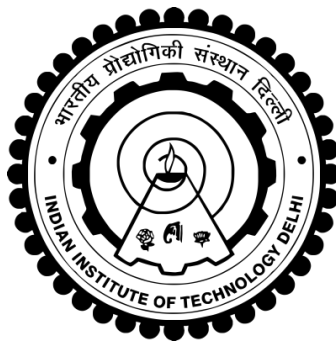


DISCOVERING RELATIONSHIPS AND PATHS IN KNOWLEDGE GRAPHS

PUNEET AGARWAL



**AMAR NATH AND SHASHI KHOSLA SCHOOL OF INFORMATION
TECHNOLOGY,
INDIAN INSTITUTE OF TECHNOLOGY DELHI
JUNE 2019**

© Indian Institute of Technology Delhi (IITD), New Delhi, 2019

DISCOVERING RELATIONSHIPS AND PATHS IN KNOWLEDGE GRAPHS

by

PUNEET AGARWAL

**AMAR NATH AND SHASHI KHOSLA SCHOOL OF INFORMATION
TECHNOLOGY**

Submitted

**in fulfilment of the requirements of the degree of Doctor of Philosophy
to the**



INDIAN INSTITUTE OF TECHNOLOGY DELHI

JUNE 2019

Certificate

This is to certify that the thesis titled **DISCOVERING RELATIONSHIPS AND PATHS IN KNOWLEDGE GRAPHS** being submitted by **Mr. PUNEET AGARWAL** for the award of **Doctor of Philosophy** in Computer Science and Engineering is a record of bona fide work carried out by him under my guidance and supervision at the Amar Nath and Shashi Khosla School of Information Technology, Indian Institute of Technology Delhi.

The work presented in this thesis has not been submitted elsewhere, either in part or full, for the award of any other degree or diploma.

Maya Ramanath
Associate Professor
Department of Computer Science and Engineering
Indian Institute of Technology Delhi
New Delhi - 110016

Gautam Shroff
Chief Scientist and Vice President
Tata Consultancy Services Ltd.
Gurgaon - 110016

Acknowledgements

Firstly, I would like to express my sincere gratitude to my advisers Prof. Maya Ramanath and Dr. Gautam Shroff for their continuous support of my Ph.D. study and related research, for their patience, motivation, and immense knowledge. Their guidance helped me at the time of research and writing of this thesis. I could not have imagined having better advisors and mentors for my Ph.D. study.

Besides my advisers, I would also like to thank the rest of my thesis committee: Prof. Amithabha Bagchi, Prof. B. Chandra, Prof. Aditeshwar Seth, and Prof. Sanjiva Prasad, for their insightful comments and encouragement, and also for the thought provoking questions which incited me to widen my research from various perspectives.

I thank Dr. Gautam Shroff, in his role as supervisor at my workplace in TCS Research, not only for motivating me to join the Ph.D. program, but also for his continuous guidance during the entire tenure. I also want to thank all the other colleagues at work in TCS Research, the insightful discussions with them have been a key source of learning for me. Most importantly, I express my gratitude to my employer Tata Consultancy Services Ltd. for allowing me to pursue Ph.D. and for the travel support related to my publications.

I also thank my fellow labmates in IIT Delhi, for the stimulating discussions, and support during entire tenure.

Last but not the least, I would like to thank my family: my parents for their love and encouragement, my wife Shilpi for managing everything at home when I was busy, and specially my daughter Navya for allowing me to invest my family time towards this endeavor. Without their immense support, it would not have been possible for me to attempt this enlightening journey.

Puneet Agarwal

Abstract

The use of knowledge graphs has proliferated into many different industrial applications in recent years, for example, in knowledge management systems, in infrastructure management systems, for the purpose of fraud detection, etc. Querying the knowledge graphs has become even more important because of such widespread industrial applications. We are interested in special types of knowledge graph queries where-in all the query keywords do not occur in the same entity, and the expected answer entity (one or more) may not contain any of the query keywords. The answer entities are connected to the entities which contain the keywords via various relationships of the knowledge graphs. As a result, standard techniques of document search are not applicable for this problem of querying the knowledge graphs, mainly because those techniques rank the documents which contain the query keywords.

In this thesis, we analyze two types of knowledge graph queries, ‘*relationship queries*’, and ‘*factoid queries*’. The first one comprises a set of keywords often indicating different entities, and the second one is a natural language sentence that mentions one or more entities. Against both of these types of queries, the answer entities may not contain the query keywords, they (answer entities) may however be connected to the entities containing the query keywords in the knowledge graph.

Relationship queries have been studied for many years, using various terminologies, e.g., keyword search, Steiner tree in a graph etc., the solutions proposed in the literature so far either don’t guarantee the optimality of answers retrieved or don’t utilize distributed parallel processing paradigm. Such an approach can be used for scaling relationship queries to large graphs having millions of nodes and edges, such graph are now publicly available in the form of ‘linked open data’(LOD). In this thesis, we present an algorithm for distributed keyword search (DKS) on large graphs, based on the graph parallel computing paradigm Pregel. We also present a proof of correctness of our algorithm. Even if terminated early, our algorithm produces approximate answers along with bounds. We describe an optimized implementation of our algorithm along with time-complexity analysis. Finally, we report experimental results on LOD data, and demonstrate efficiency of our approach on large graphs.

Further, answering natural language questions posed on a knowledge graph requires traversing an appropriate sequence of relationships starting from the mentioned entities. To answer complex queries, we often need to traverse more than two relationships. Traditional approaches traverse at most two relationships, as well as typically first retrieve

candidate sets of relationships using indexing etc., which are then compared via machine-learning. Such approaches rely on the textual labels of the relationships, rather than the structure of the knowledge graph. In this thesis, we present a novel approach KG-REP that directly predicts the embeddings of the target relationships against a natural language query, avoiding the candidate retrieval step, using a sequence to sequence neural network. Our model takes into account the knowledge graph structure via novel entity and relationship embeddings. We release a new dataset containing complex queries on a public knowledge graph that typically require traversal of as many as four relationships to answer. We also present a new benchmark result on a public dataset for this problem.

Finally, we conclude with a perspective on open problems in this domain, and highlight how the advances in this field can impact general capability of computational systems.

सार

नॉलेज ग्राफ का उपयोग हाल के वर्षों में कई अलग-अलग औद्योगिक अनुप्रयोगों में किया गया है, उदाहरण के लिए, नॉलेज प्रबंधन प्रणालियों में, बुनियादी ढांचा प्रबंधन प्रणालियों में, धोखाधड़ी का पता लगाने के उद्देश्य से, आदि। व्यापक औद्योगिक अनुप्रयोगों के कारण से नॉलेज ग्राफ को क्वेरी करना और भी महत्वपूर्ण हो गया है। हम विशेष प्रकार के नॉलेज ग्राफ प्रश्नों में रुचि रखते हैं जहां-सभी क्वेरी कीवर्ड एक ही इकाई में नहीं होते हैं, और अपेक्षित उत्तर इकाई (एक या अधिक) में कोई क्वेरी कीवर्ड नहीं हो सकता है। उत्तर इकाइयाँ उन संस्थाओं से जुड़ी होती हैं जिनमें नॉलेज ग्राफ के विभिन्न संबंधों के माध्यम से खोजशब्द होते हैं। नतीजतन, दस्तावेज़ खोज की मानक तकनीक नॉलेज ग्राफ को क्वेरी करने की इस समस्या के लिए लागू नहीं होती है, मुख्यतः क्योंकि उन तकनीकों में दस्तावेज़ होते हैं जो क्वेरी कीवर्ड होते हैं।

इस थीसिस में, हम दो प्रकार के नॉलेज ग्राफ प्रश्नों का विश्लेषण करते हैं, 'रिलेशनशिप क्वेरी', और 'फक्टोइड क्वेरी'। पहले वाले में अक्सर अलग-अलग संस्थाओं को दर्शाने वाले कीवर्ड का एक सेट होता है, और दूसरा एक प्राकृतिक भाषा का वाक्य है जिसमें एक या अधिक संस्थाओं का उल्लेख होता है। इन दोनों प्रकार के प्रश्नों के विरुद्ध, उत्तर निकाय में क्वेरी कीवर्ड नहीं हो सकते हैं, वे (उत्तर इकाइयाँ) हालांकि नॉलेज ग्राफ में क्वेरी कीवर्ड वाली संस्थाओं से जुड़े हो सकते हैं।

विभिन्न शब्दावली का उपयोग करते हुए कई वर्षों से संबंध प्रश्नों का अध्ययन किया गया है, उदाहरण के लिए, कीवर्ड खोज, ग्राफ में स्टेनर ट्री आदि, साहित्य में प्रस्तावित समाधान या तो पुनर्प्राप्त किए गए उत्तरों की इष्टतमता की गारंटी नहीं देते हैं या वितरित उपयोग नहीं करते हैं। समानांतर प्रसंस्करण प्रतिमान। इस तरह के दृष्टिकोण का उपयोग लाखों से अधिक नोड्स और किनारों वाले बड़े ग्राफ के संबंध प्रश्नों को स्केल करने के लिए किया जा सकता है, ऐसे ग्राफ अब सार्वजनिक रूप से 'लिंकड ओपन डेटा' (एलओडी) के रूप में उपलब्ध हैं। इस थीसिस में, हम रेखांकन समानांतर कंप्यूटिंग प्रतिमान के आधार पर बड़े रेखांकन पर वितरित कीवर्ड खोज (DKS) के लिए एक एल्गोरिथ्म प्रस्तुत करते हैं। हम अपने एल्गोरिथ्म की शुद्धता का प्रमाण भी प्रस्तुत करते हैं। यहां तक कि अगर जल्दी समाप्त कर दिया, हमारे एल्गोरिथ्म सीमा के साथ अनुमानित जवाब पैदा करता है। हम समय-जटिलता विश्लेषण के साथ-साथ हमारे एल्गोरिथ्म के एक अनुकूलित कार्यान्वयन का वर्णन करते हैं। अंत में, हम एलओडी डेटा पर प्रयोगात्मक परिणामों की रिपोर्ट करते हैं, और बड़े रेखांकन पर हमारे दृष्टिकोण की दक्षता प्रदर्शित करते हैं।

इसके अलावा, नॉलेज ग्राफ पर रखे गए प्राकृतिक भाषा के सवालों का जवाब देने के लिए उल्लेखित संस्थाओं से शुरू होने वाले रिश्तों का एक उचित क्रम निर्धारित करना होगा। जटिल प्रश्नों का उत्तर देने के लिए, हमें अक्सर दो से अधिक रिश्तों का पता लगाना पड़ता है। पारंपरिक दृष्टिकोण अधिकांश दो रिश्तों पर चलते हैं, साथ ही आमतौर पर सबसे पहले अनुक्रमित आदि का उपयोग करते हुए रिश्तों के उम्मीदवार सेटों को पुनः प्राप्त करते हैं, जिनकी तुलना मशीन-लर्निंग के माध्यम से की जाती है। इस तरह के दृष्टिकोण नॉलेज ग्राफ की संरचना के बजाय रिश्तों के शाब्दिक लेबल पर निर्भर करते हैं। इस थीसिस में, हम केजी-आरईपी को एक उपन्यास दृष्टिकोण प्रस्तुत करते हैं जो सीधे एक प्राकृतिक भाषा क्वेरी के खिलाफ लक्ष्य संबंधों के एम्बेडिंग की भविष्यवाणी करता है, उम्मीदवार पुनर्प्राप्ति चरण से बचते हुए, एक अनुक्रम का उपयोग करके तंत्रिका नेटवर्क का अनुक्रम करता है। हमारे मॉडल नॉलेज ग्राफ संरचना को उपन्यास इकाई और संबंध एम्बेडिंग के माध्यम से ध्यान में रखते हैं। हम एक नए डेटासेट को एक सार्वजनिक नॉलेज ग्राफ पर जटिल प्रश्नों से मुक्त करते हैं, जिसका जवाब देने के लिए आम तौर पर चार रिश्तों के ट्रावेल की आवश्यकता होती है। हम इस समस्या के लिए एक सार्वजनिक डेटासेट पर एक नया मानदंड परिणाम भी प्रस्तुत करते हैं।

अंत में, हम इस डोमेन में खुली समस्याओं पर एक परिप्रेक्ष्य के साथ निष्कर्ष निकालते हैं, और इस बात पर प्रकाश डालते हैं कि कैसे इस क्षेत्र में उन्नति कम्प्यूटेशनल सिस्टम की सामान्य क्षमता को प्रभावित कर सकती है।

Contents

List of Figures	xii
List of Tables	xv
1 Introduction	1
1.1 Organization of the thesis	5
1.1.1 Part I: Discovering Relationships in Knowledge Graphs	5
1.1.2 Part II: Discovering Paths in Knowledge Graphs	6
1.1.3 Conclusion and Appendix	6
I Discovering Relationships in Knowledge Graphs	9
2 Relationship Queries - Problem Definition	11
2.1 Motivation and Use-cases	11
2.2 Problem Description & Definitions	13
2.3 Solution Overview for Relationship Queries	14
3 Relationship Queries - Background	17
3.1 Top-K Algorithms	17
3.2 Distributed & Parallel Processing Paradigms	19
3.2.1 Map-Reduce	19
3.2.2 Pregel	20
3.3 Related Concepts	22
3.4 Summary	23
4 Relationship Queries - Related Work	25
4.1 Solution Blueprint and Organization	26
4.2 Overview of Available Solutions	27
4.2.1 Steiner Tree Solutions	28
4.2.2 Online Search in data graph	29
4.2.3 Schema Graph Based	32
4.2.4 Index Based	34

4.2.5	Distributed Parallel Approaches	37
4.3	Common Solution Aspects	38
4.3.1	Keyword Query Response	38
4.3.2	Answer Ranking Schemes	40
4.3.3	Top-K algorithms	44
4.3.4	Time Complexity and Minimum Steiner Trees	45
4.3.5	Validation	47
4.4	Summary	48
5	DKS Algorithm and Empirical Analysis	49
5.1	DKS Algorithm - Overview	50
5.1.1	Pre-processing	50
5.1.2	Algorithm overview	50
5.2	Termination conditions in DKS	51
5.2.1	Termination on finding the top- k answers	51
5.2.2	Termination on exchanging n messages	54
5.3	DKS - Detailed Description	54
5.4	Proof of Correctness	59
5.4.1	Overview and Intuition	59
5.4.2	Minimum answer weight increases	60
5.4.3	Identify Candidate Nodes	61
5.4.4	BFS Stopping Criterion	63
5.4.5	Global Vs Local TOP- k Steiner Tree	64
5.5	Time Complexity Analysis	65
5.5.1	Communication Cost	67
5.6	Experiments and Analysis	68
5.6.1	Data, Infrastructure, and Implementation	68
5.6.2	Experimental Results	69
5.6.3	Approximation, Distribution Overheads	73
5.6.4	Analysis and Discussion	76
5.7	Summary	77
II	Discovering Paths in Knowledge Graphs	79
6	Path Prediction for Natural Language Queries	81
6.1	Motivation	83
6.2	Problem Definition	84
6.3	Solution Overview	84

7	Factoid Queries: Background and Preliminaries	87
7.1	Representation Learning	87
7.2	Background on Deep Neural Networks	88
7.2.1	Recurrent Neural Network	89
7.2.2	Long Short Term Memory (LSTM)	89
7.2.3	Sequence to Sequence Model	90
7.2.4	Other Related Concepts	91
7.3	Question Answering Systems in General	91
7.4	Factoid Query Solution Blueprint	93
8	KG-REP: Related Work	95
8.1	Graph Representation	95
8.1.1	General Graph Embeddings	96
8.1.2	Knowledge Graph Structure Learning	99
8.1.3	Structure and Entity Description	101
8.2	Factoid Queries on Knowledge Graphs	103
8.2.1	Semantic Parsing based approaches	103
8.2.2	Single Fact Factoid Queries	106
8.2.3	Two Relationship Retrieval	108
8.3	Summary	111
9	KG-REP: Knowledge Graph Relationship Embedding Prediction	113
9.1	KG-REP Overview	114
9.1.1	Knowledge Graph Representation Learning	116
9.1.2	KG-REP: KG Relationship Embedding Prediction	117
9.1.3	Beam Search for Relationship Retrieval	119
9.2	Experimental Results	119
9.2.1	Description of Datasets used	119
9.2.2	Baseline Comparison - WebQSP	121
9.2.3	Experiments on QSMRQ dataset	122
9.2.4	Analysis and Discussion	122
9.3	Summary	123
10	Conclusions and Future Work	125
	Bibliography	129
	Appendix A Glossary of Terms	145
	Appendix B Knowledge Graph Embedding and Traversal	147
B.1	Vector Representation of the knowledge graph	147
B.1.1	Can node embeddings be used for graph traversal?	147

List of Figures

1.1	(a) Sample knowledge graph, showing entities e_i marked with dark circles, and relationships r_j marked with arrows; (b) Sample RDF data	2
2.1	Collusive Fraud use-case for Relationship Query	12
3.1	Shows a sample graph and how various nodes are distributed among different worker agents.	20
3.2	Part of the image taken from [100]. Shows the process of finding maximum valued node from a graph using Pregel.	21
4.1	Flow chart of abstract solution to keyword search on graph data.	26
4.2	Grouping of some of the most cited related works	28
4.3	Example of a bi-directional traversal using BANKS-2; and corresponding answer tree, with root node at u_8	31
4.4	Sample DeweyNumber Allocation, image taken from [145]	36
4.5	Answer Tree Q-Fragments, image taken from [80]	39
5.1	BFS traversal example, with answer tree shown in red color; and DKS Flowchart shown below, with respect to Pregel.	52
5.2	a) Need for Deep Message: From this example, it can be observed that using only the BFS messages, path to all keywords of the query will never become available at root nodes v_5 . b) Need to propagate deep Message: from this example it can be observed that unless we propagate the deep messages, paths to all keywords of the query will never become available, at root nodes v_8, v_9 . Here, message type is assumed to be BFS if not marked.	56
5.3	Node v_{21} receives messages $\{M1, M2, M3, M4\}$ from nodes $\{v_{17}, v_{18}, v_{19}, v_{20}\}$ respectively; Resulting local-tree of v_{21} is shown. Branches marked with dotted edges are dropped to obtain <i>filtered local tree</i>	57
5.4	For proof of Lemma 5.4.1	61
5.5	How a sorted list of the minimum partial answer weights is extracted from the partial answer weights available in order of their discovery.	62
5.6	For proof of Lemma 5.4.3	63
5.7	Shortest path to q_1 not in top- k answers	64

5.8	For <i>query set A</i> , number of keyword nodes (on a log scale) for the 20 queries for various keyword counts ($kc = \{2, 3, 4, 5, 6\}$), i.e., for all the 100 queries. Trend-line shown on exponential scale.	69
5.9	x-axis - Queries, y-axis - Execution Time taken in seconds for different values of k , on bluk-bnb knowledge graph using <i>query set A</i> . Here, Execution time = (total time taken - time taken for instantiation of the worker nodes, first time loading of the graph and serialization of the final results).	71
5.10	A) Deep Message Count for all 100 queries of <i>query set A</i> , on bluk-bnb knowledge graph, shown varying with increasing values of k	72
5.11	SPA-Ratio for queries (on x-axis) that are arranged in increasing order of keyword-count and keyword node count for bluk-bnb, when using <i>query set A</i> . Here, the SPA ratio is marked as zero (equivalent to “NC” i.e., not calculated) for those queries for which the first termination condition was satisfied and SPA ratio was not calculated. SPA ratio as more than zero and less than or equal to one, is observed for those queries for which the optimal answer was calculated in the last superstep. See Section 5.6.3 for more details.	73
5.12	Percentage of nodes explored w.r.t. $ V $, averaged for $k = \{1, 2, 5, 10\}$, on bluk-bnb knowledge graph, when using <i>query set A</i>	74
5.13	Percentage of nodes explored using <i>query set B</i> on sec-rdfabout knowledge graph, and another 20 queries on bluk-bnb knowledge graph. We can say that whole graph is not explored. In these runs uniform edge weight of all the edges was used, if the node degree was less than 1001, and infinite otherwise. For the results shown here, only the first termination condition of finding top- k results was used.	74
5.14	Total number of messages as percentage of $ E $, shown varying with respect to different values of K , for all the 100 queries on bluk-bnb dataset.	75
5.15	Parallel efficiency of DKS, for queries resulting in optimal answers. Queries on bluk-bnb dataset did not run successfully on lesser than 10 workers.	76
6.1	(a) Freebase subgraph, for a sample query taken from QSMRQ dataset. (We release this dataset. See Chapter 9, Section 9.2 for details.) (b) Freebase subgraph for a query, taken from WebQSP dataset [154].	82
7.1	Schematic diagram of LSTM Unit	90
8.1	Images taken from [57], explains the basis of formulating discrete probability distribution of a node in the graph for its outgoing edges.	97
8.2	Image taken from [154], shows the neural network architecture of used in this paper.	110

9.1	(a) Freebase subgraph, for a sample query taken from QSMRQ dataset. (We release this dataset. See Section 9.2 for details.) (b) Freebase subgraph for a query, taken from WebQSP dataset [154]	115
9.2	Deep Neural Network architecture for Sequence to Sequence Model	118
B.1	Rank of the next node for Steiner Tree traversal, among predicted nodes	148

List of Tables

4.1	Organization of the current chapter	27
4.2	Analysis of graph keyword search solution, with respect to ranking aspects adopted	42
4.3	Top-K algorithms used by various solutions	45
4.4	Time Complexity of various solutions; For Discovery [68], $ S $ = number of candidate networks, and J = average size of a candidate network. This time-complexity refers to Greedy algorithm for query execution only. For other approaches, $n = V $ (number of nodes the in the graph), $e = E $ (number of edges in the graph), T_i are the keyword nodes containing keyword q_i , and m = No of keywords.	46
5.1	Summary of notations used	52
5.2	Percentage of Time taken by DKS components, for different values of k , using <i>query set A</i> , on bluk-bnb knowledge graph	72
7.1	Sample factoid queries against a knowledge graph.	93
8.1	Graph Representation Learning	96
8.2	Knowledge Graph Question Answering	103
9.1	Summary of notations used, in general boldface variables (e.g., \mathbf{e}_i) are used for embeddings of corresponding non-boldface variable (e.g., e_i).	116
9.2	Results on WebQSP Dataset, comparing the accuracies for two situations, when using the <i>ERE embedding</i> of the mentioned entities in input query (ERE), and without using them (no-ERE), and baseline approach.	121
9.3	Results on QSMRQ Dataset, comparing the accuracies for: 1) when using the <i>ERE embedding</i> in both input and output, 2) without using them in the input queries.	123
A.1	Glossary of terms used in this thesis	145