



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 3, Issue 12, December 2015

Probabilistic RDF Graphs for Keyword Search

Bhansali Shubham, Madhuri Chavan, Landge Aishwarya, Lokhande Shraddha.

B.E Students, Dept. of IT, DYPIET, Pimpri. Savitribai Phule University, Pune, Maharashtra, India

ABSTRACT: In many real applications, RDF (Resource Description Framework) has been widely used to describe data in the Semantic Web in a W3C standard. RDF data may often suffer from the independency of their data sources, and exhibit errors or contrariety. Such unreliable RDF data by probabilistic RDF graphs, and study an important problem, probabilistic RDF graphs for keyword search query (namely, the pg-KWS query). To retrieve meaningful keyword search answers, propose system design the score rankings for sub graph answers specific for RDF data. The keyword searching technique over uncertain graph is introduced. The Keyword routing method is used to route the keywords to applicable source. In this Approach two methods are included. The keyword relationship graph concludes the relationship between keywords and the element mentioning them. The scoring mechanism computes the score of keywords at each level which reduces the imprecision. The result will include the sub tree of the entire graph which includes all keywords of input query having high score and it retrieves the most significant data.

KEYWORDS: Probabilistic RDF graph, Keyword search, PG-KWS Uncertain graph, Keyword routing.

I. INTRODUCTION

Keyword search has been concluded to retrieve useful data from database. Keyword search has major benefit i.e. it is easy to operate. Users do not have to understand the database schema and the query language, and can gain the knowledge rapidly how to use information retrieval. Now a days, the study of keyword search technology based on graph data has become a hot spot, and it is generally applied to the field of information retrieval. In the field of traditional graph database, the research on keyword search has already gained some achievement, but in the field of uncertain graph data, the study on keyword search has scarcely started. Especially recently, quite a lot of efforts have been put for keyword search over graphs. However, all graphs in the database are assumed to be certain or valuable, and this assumption is often not valid in real-life applications, as XML data and RDF data can be highly unreliable due to errors in the web data or data expiration. In the application of the data integration, it is needed to include such RDF data from various data sources into an incorporated database. Uncertainties or independencies often exist in this case. Like In social networks, each link between any two persons is often joined with a probability that represents the uncertainty of the link or the strength of influence a person has over another person in viral marketing. XML data having tree or graph form, uncertainties are integrated in XML documents known as probabilistic XML document.

Keyword searching in RDF data, social networks and XML data have many weighty applications. For data with XML and relational schema, specific query languages, such as SQL and XQuery, have been developed for information retrieval. In order to query such data, the user must master a complex query language and understand the underlying data schema. In relational databases, information about an object is often inflected in multiple tables due to normalization considerations, and in XML datasets, the schema are often complicated and embedded XML structures often create a lot of difficulty to express queries that are forced to traverse tree structures. Furthermore, many applications work on graph-structured data with no obvious, well-structured schema, so the option of information retrieval based on query languages is not applicable. Both XML databases and relational databases can be viewed as graphs. Specifically, XML datasets can be regarded as graphs when IDREF/ID links are taken into consideration, and a relational database can be regarded as a data graph that has triplet and keywords as nodes. In the data graph, for example, two tuples or triplets are connected by an edge if they can be

International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

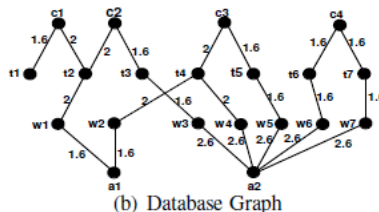
Vol. 3, Issue 12, December 2015

associated using a foreign key; a tuple or triplet and a keyword are connected if the tuple contains the keyword. Thus, traditional graph search algorithms, which extract features from graph data, and convert queries into searches over feature spaces, can be used for such data. Therefore, it is necessary to relax the strict assumption of Deterministic or certain graphs and study keyword search over uncertain graphs.

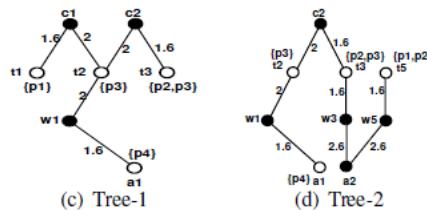
Keyword Query Analysis is the greatest possible or the ultimate goal of research on uncertain graph data management to retrieve the useful data from uncertain graph data. In Fig 1. the relational database is considered for keyword searching using graphs. The database includes author data which provides the information about author's id and author's name. Next it includes paper data which provides its id and title. The database also includes the relational data between paper and author data which includes paper id and author id. Then this, the relationship is represented among these data through a graphical structure. Whatever the input keyword query is entered, the keywords are searched in graph and routes are found out to reach keywords and display the routed sub graph in results.

Author		Paper		Paper-Author	
AID	Name	PID	Title	PID	AID
a1	Jim	t1	Keyword Search on RDBMS	t2	a1
a2	Robin	t2	Steiner Problem in DB	t4	a1
		t3	Efficient IR-Query over DB	t3	a2
		t4	Online Cluster Problems	t4	a2
		t5	Keyword Query over Web	t5	a2
		t6	Query Optimization on DB	t6	a2
		t7	Parameterized Complexity	t7	a2

(a) Database



(b) Database Graph



(c) Tree-1

(d) Tree-2

Fig 1. A Motivation Example

II. MOTIVATION AND BACKGROUND

Entity-Relationship graphs are receiving great attention for information management outside of mainstream database engines. Specifically, the Semantic-Web data model RDF is gaining popularity for applications such as biological networks, social Web2.0 applications, large-scale knowledge bases such as DB pedia or YAGO, and more generally, as a light-weight exhibition for the "Web of data". An RDF data collection consists of a set of SPO, SPO triples for short. In ER term, an SPO triple corresponds to an entity connected to the value of a named attribute or to a pair of entities connection by a named relationship. As the instance of a triple can in turn be the subject of other triples, RDF data can also be viewed



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 3, Issue 12, December 2015

as a graph where nodes correspond to entities and edges to relationships of typed nodes and typed edges (viewing attributes as relations as well). Some of the existing RDF data contain more than a billion triples.

III. LITERATURE REVIEW

[1] Clustering Large Probabilistic Graphs:

Problem of clustering probabilistic graphs is identical to the problem of clustering standard graphs, probabilistic graph clustering has numerous applications, like finding complexes in probabilistic protein-protein interaction networks and discovering groups of users in affiliation networks. The edit-distance based definition of graph clustering to probabilistic graphs. Establish a connection between objective function and correlation clustering to propose practical approximation algorithms for problem. A benefit of approach is that objective function is parameter-free. Therefore, the number of clusters is part of the output. It also develop methods for testing the statistical significance of the output clustering and study the case of noisy clustering. Using a real protein-protein interaction network and ground-truth data, methods discover the correct number of clusters and identify established protein relationships. Finally, the practicality techniques using a large social network of Yahoo! users consisting of one billion edges.

[2] Scalable Keyword Search on Large RDF Data:

Keyword search is a beneficial tool for researching large RDF datasets. Existing techniques either depend on constructing a distance matrix for pruning the search space or building upshot from the RDF graphs for query processing. Existing techniques have serious limitations in handling with realistic, large RDF data with millions of triples. Moreover, the existing summarization techniques may lead to incomplete and incorrect results. These issues can be addressed by an effective summarization algorithm to summarize the RDF data. For a given keyword query, the summaries gives significant pruning powers to exploratory keyword search and output in much better efficiency compared to previous works. Unlike existing techniques, search algorithms always return correct and complete results. In addition to this, the summaries we built can be updated efficiently and incrementally. Experiments on both large real RDF data sets and bench-mark show that techniques are scalable and efficient.

[3] Keyword Search over RDF Graphs:

Large knowledge bases consisting of entities and relationships between them have become vital sources of information, which is further more used in many applications. Most of these knowledge bases adopt the Semantic- Web data model RDF as a representation model. Querying these knowledge bases is usually done using structured queries which uses graph-pattern languages such as SPARQL. Such kind of structured queries require some expertise from users which limits the accessibility to such significant data sources for security purposes. To avoid this, keyword search must be supported. A retrieval model for keyword queries over RDF graphs. Retrieves a set of subgraphs that match the query keywords, and ranks them based on statistical language models. Retrieval model outperforms the-state-of-the-art IR and DB models for keyword search over structured data using experiments over two real-world datasets.

[4] Top-k Keyword Search Over Probabilistic XML Data:

The act of increasing of work on XML keyword query, it remains open to support keyword query over probabilistic XML data. Compared with traditional keyword search, it is far more expensive to reply a keyword query over probabilistic XML data due to the consideration of possible world semantics. The problem of studying top-k keyword search over probabilistic XML data, which is to retrieve k SLCA results with the k highest probabilities of existence. And then we propose two efficient algorithms. The first algorithm PrStack can find k SLCA results with the k highest probabilities by scanning the relevant keyword nodes only once. To further improve the efficiency, a second algorithm EagerTopK based on a set of



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 3, Issue 12, December 2015

pruning properties which can quickly prune unsatisfied SLCA candidates. Finally two algorithms and compare their performance with analysis of extensive experimental results.

[5] Searching RDF Graphs SPARQL and Keywords:

The act of increasing of knowledge-sharing communities like Wikipedia and the advances in automated information descent from Web pages enable the construction of ample knowledge bases with reality or facts about entities and their relationships. The facts can be displayed in the RDF data model, as so-called subject-property-object triples, and thus can be queried by structured query languages like SPARQL. According to theory, this allows valuable querying in the database spirit. Though, RDF data may be highly diverse and queries may return way too many results, so that ranking by informativeness measures is crucial to avoid staggering users. Furthermore, as facts are extracted from textual contexts or have community provided annotations, which can be helpful in considering keywords for formulating search requests. Ranking retrieval of RDF data with keyword-augmented structured queries is overview of ongoing and recent work. The ranking method is based on statistical language models which consists the state-of-the-art paradigm in information retrieval which develops a novel form of language models for the structured, but schema-less setting of extended SPARQL queries and RDF triples.

[6] Representing Probabilistic Relations In RDF:

Probabilistic inference will be of special significance when one needs to know how much can say with what all know given new observations. Bayesian Network is a graphical probabilistic model with which one can represent probabilistic relations intuitively and various efficient algorithms for inference are developed. Now ongoing work in its design stage which provides a vocabulary for exhibiting probabilistic knowledge in a RDF graph which is to be mapped to a Bayesian Network to do inference on it

[7] Efficient IR-style Keyword Search over Relational Databases:

Applications in which plain text coexists with structured Data are pervasive. Commercial relational database management systems generally provide querying capabilities for text attributes that incorporate state-of-the-art information retrieval (IR) relevance ranking strategies, but search functionality requires that queries specify the exact column or columns against which a given list of keywords is to be matched. This requirement can be cumbersome and inflexible from a user perspective: good answers to a keyword query might need to be assembled in perhaps unforeseen ways by joining tuples from multiple relations motivated recent research on free-form keyword search over RDBMSs. Adapt IR style document-relevance ranking strategies to the problem of processing free-form keyword queries over RDBMSs. Query model can handle queries with both AND and OR semantics, and exploits the sophisticated single-column text-search functionality often available in commercial RDBMSs. Develop query-processing strategies that build on a crucial characteristic of IR-style keyword search: only the few most relevant matches according to some definition of relevance are generally of interest. Consequently, rather than computing all matches for a keyword query, which leads to inefficient executions, techniques focus on the top-k matches for the query, for moderate values of k. A thorough experimental evaluation over real data shows the performance advantages of approach.

IV. PROPOSED SYSTEM

We propose effective pruning methods to quickly filter out false alarms. Extensive researches have been conducted to verify the effectiveness and efficiency of our proposed approaches proposed to answer keyword search queries on certain graphs. We can transform probabilistic RDF graph with uncertain vertex/edge keywords to the one of uncertain keywords in vertices only, to which we will apply our proposed approaches. We will utilize the entropy concept to propose a metric that will indicate our connection to RDF keyword search results in probabilistic RDF graphs. We will propose two pruning

International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 3, Issue 12, December 2015

strategies, score bound pruning and probabilistic pruning, which utilize score bounds or probabilistic threshold, respectively, to enable the pruning. Our proposed pruning methods via score bounds. We propose a heuristic based algorithm, which obtains PWGs of a probabilistic sub graph g with low cost. We report the experimental results of our proposed approaches for answering pg-KWS queries on both real and synthetic data. The proposed methods explored pruning technique with graph structures and matching probabilities. various probabilistic queries over uncertain data have been proposed, including probabilistic range query (PRQ) , nearest neighbor (PNN), reverse nearest neighbor (PRNN) , skyline (PSQ) , reverse skyline (PRSQ) , and similarity join (PSJ). The bidirectional search, which uses the pre-computed distance between keywords and nodes, obtains the bounds of ranking scores to enable fast pruning and retrieval.

4.1 Probabilistic-Graph Model:

Similar to deterministic graphs, probabilistic graphs may be Bidirectional and carry additional labels on the edges such as weights model assumes independence among edges it focuses on probabilistic graph which are mostly independent. It represents a probabilistic graph using tuple. The probabilistic graphs are represented with unweighted probabilistic Graph. One can think of a probabilistic graph as a generative model for deterministic graphs. A deterministic graph is generated by connecting two nodes via an edge with probability. Deterministic graphs are an instance of probabilistic graphs for which random graphs are an instance of probabilistic graphs where all edge probabilities are the same and equal. Then there are distinct graphs that can be generated .They use the term possible world to refer to each such graph.

V. ARCHITECTURAL DESIGN

The main objective of our approach is to search keyword in uncertain graph data and in addition to retrieve the relevant data for input query.

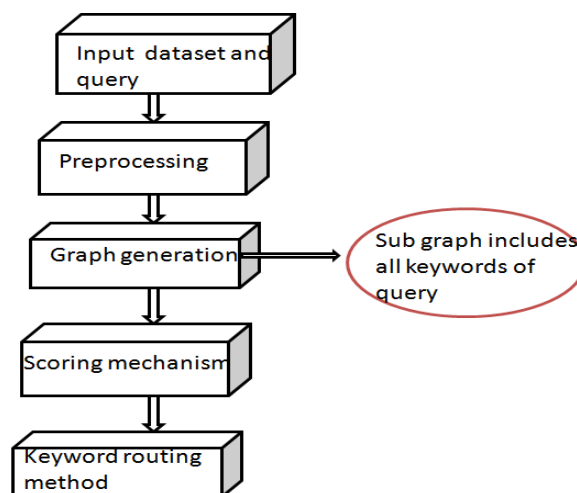


Fig 2. Architectural Design

The above mentioned modules are used in our approach to search keywords in an uncertain graph data and routes to reach the query keywords and finally shows sub tree in result which includes all keywords entered by users and in addition it shows most relevant data related to query keywords.



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 3, Issue 12, December 2015

Sr No	Author	Pros	Cons	Motivation
[1]	George Kollios	Objective function is parameter free	Usually assumed on uncertain data	Probabilistic graphs
[2]	Wangchao Le	Explores large RDF data sets	Dealing with realistic large RDF data with tens of millions of triples	Use of large RDF data sets
[3]	R. Blanco	Use of structured queries	User must have knowledge of query language for searching	Use of semantic web data model RDF
[4]	J. Li	Support keyword query over probabilistic XML data	Expensive to answer keyword query	Probabilistic XML data
[5]	S.Elbasuoni	Retrieval of RDF data with keyword augmented structured queries	User must have knowledge of query language for searching	Ranking methods for information retrieval
[6]	Y. Fukushige	Provides vocabulary for representing probabilistic knowledge in RDF graph	Needs to know how much can say with new observation	Use of vocabulary
[7]	V. Hristidis	Develop query processing strategies	Inefficient executions	Ranking strategies

Table 1. Survey Table

VI. CONCLUSION

Hence, we formulate and tackle the problem of keyword search through probabilistic RDF graphs. Besides that RDF graph creates XML file which support the structured data. Moreover, we define probabilistic graphs containing correlated adjacent edges as correlated probabilistic graphs

References

- [1] Hoang George Kollios, Michalis Potamias, and Evimaria Terzi, Clustering Large Probabilistic Graphs, IEEE vol. 25, NO. 2, February 2013.
- [2] Wangchao Le, Feifei Li, Anastasios Kementsietsidis, Songyun Duan, Scalable Keyword Search on Large RDF Data, IEEE 2013.
- [3] S. Elbasuoni and R. Blanco, Keyword search over RDF graphs, in Proc. 20th ACM Int. Conf. Inform. Knowl. Manage., 2011, pp. 237242.
- [4] J. Li, C. Liu, R. Zhou, and W. Wang, Top-k keyword search over probabilistic XML data, in Proc. IEEE 27th Int. Conf. Data Eng., 2011, pp. 673684.
- [5] S. Elbasuoni, M. Ramanath, R. Schenkel, and G. Weikum, Searching RDF graphs with SPARQL and keywords, IEEE Data Eng. Bull., vol. 33, no. 1, pp. 1624, Mar. 2010.
- [6] Y. Fukushige, Representing probabilistic relations in RDF, ISWC-URSW, pp. 106107, 2005.
- [7] V. Hristidis, L. Gravano, and Y. Papakonstantinou, Efficient IR-style keyword search over relational databases, in Proc. 29th Int. Conf. Very large data, 2003, pp. 850861.



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 3, Issue 12, December 2015

- [8] R. Cheng, D. V. Kalashnikov, and S. Prabhakar, "Evaluating probabilistic queries over imprecise data," in Proc. ACM SIGMOD Int. Conf. Manage. Data, 2003, pp. 551–562
- [9] D. Papadias, Y. Tao, G. Fu, and B. Seeger, "An optimal and progressive algorithm for skyline queries," in Proc. ACM SIGMOD Int. Conf. Manage. Data, 2003, pp. 467–478
- [10] K. Wilkinson, C. Sayers, H. Kuno, D. Reynolds, and J. Database, "Efficient RDF storage and retrieval in Jena2," in Proc. Eur. Semantic Web Databases workshop, 2003, pp. 131–150.
- [11] A. Hulgeri and C. Nakhe, "Keyword searching and browsing in databases using BANKS," in Proc. 18th Int. Conf. Data Eng., 2002, pp. 431–440
- [12] S. Borzs € onyi, D. Kossmann, and K. Stocker, "The skyline operator," in Proc. 17th Int. Conf. Data Eng., 2001, pp. 421–430.

BIOGRAPHY

Bhansali Shubham, Madhuri Chavan, Landge Aishwarya, Lokhande Shraddha would like to acknowledge and heart felt gratitude to our guide Prof. Ashwini Dhoke, DYPIET for her expert guidance and encouragement.