Human-centric Computing
and Information Sciences
a SpringerOpen Journal

**RESEARCH**                                                                                    **Open Access**

# Finding relevant semantic association paths through user-specific intermediate entities

Viswanathan V[1*] and Ilango Krishnamurthi[2]

* Correspondence: visu.tgv@gmail.com
[1]Department of Computer Applications, Sri Krishna College of Engineering and Technology, Tamil nadu, India
Full list of author information is available at the end of the article
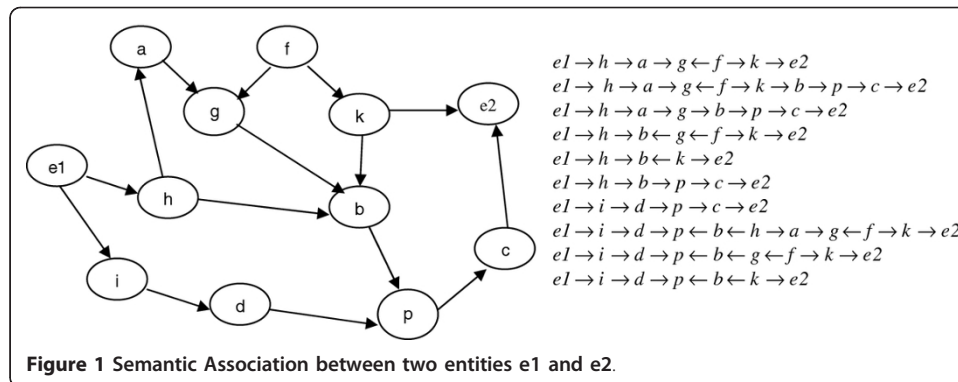
## Abstract

Semantic Associations are complex relationships between entities over metadata represented in a RDF graph. While searching for complex relationships, it is possible to find too many relationships between entities. Therefore, it is important to locate interesting and meaningful relations and rank them before presenting to the end user. In recent years e-learning systems have become very popular in all fields of higher education. In an e-learning environment, user may expect to search the semantic relationship paths between two concepts or entities. There may be numerous relationships between two entities which involve more intermediate entities. In order to filter the size of results set based on user's relevance, user may introduce one or more known intermediate entities. *In this paper, we present a Modified bidirectional Breadth-First-Search algorithm for finding paths between two entities which pass through other intermediate entities and the paths are ranked according to the users' needs.* We have evaluated our system through empirical evaluation. We have compared the execution time to discover the paths between entities for our proposed search method and existing method. According to our experiments our proposed algorithm improves search efficiently. The average correlation coefficient between the proposed system ranking and the human ranking is 0.69. It explains that our proposed system ranking is highly correlated with human ranking.

**Keywords:** Semantic Web, Semantic Association, Complex relationship, RDF, RDF Schema

## Introduction

Information retrieval over semantic metadata [1] has received a great amount of interest in both industry and academia. The semantic web contains not only resources, but also includes the heterogeneous relationships among them. In the current generation technologies of search engine, it is very difficult to find the relationships among entities. For example, how two entities are related is the most crucial question. Discovering relevant sequences of relationships between two entities answers this question. Semantic association represents a direct or indirect relationship between two entities. Different entities may be related in multiple ways. Figure 1. Illustrates a small graph of entities and the results of a query for semantic associations taking two of them as input.

There are ten different relationship paths between e1 and e2 in this graph. Here, user has to go through the entire set to find the relevant path.

$$e1 \rightarrow h \rightarrow a \rightarrow g \leftarrow f \rightarrow k \rightarrow e2$$
$$e1 \rightarrow h \rightarrow a \rightarrow g \leftarrow f \rightarrow k \rightarrow b \rightarrow p \rightarrow c \rightarrow e2$$
$$e1 \rightarrow h \rightarrow a \rightarrow g \rightarrow b \rightarrow p \rightarrow c \rightarrow e2$$
$$e1 \rightarrow h \rightarrow b \leftarrow g \leftarrow f \rightarrow k \rightarrow e2$$
$$e1 \rightarrow h \rightarrow b \leftarrow k \rightarrow e2$$
$$e1 \rightarrow h \rightarrow b \rightarrow p \rightarrow c \rightarrow e2$$
$$e1 \rightarrow i \rightarrow d \rightarrow p \rightarrow c \rightarrow e2$$
$$e1 \rightarrow i \rightarrow d \rightarrow p \leftarrow b \leftarrow h \rightarrow a \rightarrow g \leftarrow f \rightarrow k \rightarrow e2$$
$$e1 \rightarrow i \rightarrow d \rightarrow p \leftarrow b \leftarrow g \leftarrow f \rightarrow k \rightarrow e2$$
$$e1 \rightarrow i \rightarrow d \rightarrow p \leftarrow b \leftarrow k \rightarrow e2$$

**Figure 1 Semantic Association between two entities e1 and e2**.

Nowadays e-learning systems [2] have become very popular in higher education. In an e-learning environment, user may expect to know the relationships between two concepts or entities. For example, how person 'X' is related to person 'Y'. To answer this question, we need to find the semantic association paths between person 'X' and person 'Y'. Since, person 'X' may be related to person 'Y' with one or more intermediate entities, the results may be too large depending on the size of the RDF graph and most of the relationships may be irrelevant to the user. In order to filter the irrelevant paths user may introduce one more well-known intermediate entity. For example, user wants to find the relationship paths between person 'X' and person 'Y' with respect to person 'Z', he/she can introduce 'Z' as intermediate entity to find the paths. Now, we get only paths between entity 'X' and entity 'Y' that pass through the entity 'Z'. So, the size of the resultant path set can be reduced and relevancy may be improved. The challenge here is to find one or more useful paths in a quick and efficient way. A good path can bring valuable information about the relation between the two entities.

*In this paper, we propose a Modified Bidirectional Breadth-First-Search algorithm to discover the paths between two entities which pass through user-specific intermediate entities. Then the paths are ranked according to the user's need. Our proposed algorithm helps to improve the efficiency in discovery of the paths by filtering the extraneous paths.*

The paper is structured as follows: In Section 2, we present an overview of the data model and basic definitions of semantic association. In Section 3, we overview some related works in the area of semantic association. In section 4, we explain our Modified bidirectional breadth-first search algorithm. In Section 5, we explain the approach for ranking semantic associations. Experimental evaluation of the proposed approach is explained in Section 6. In Section 7, we discuss our contribution and state possible future work.

## Background

### Data model

The Resource Description Framework (RDF) is a World Wide Web Consortium (W3C) recommendation for describing Web resources. The Resource Description Framework (RDF)[3] data model provides a framework to capture the meaning of an entity by specifying how it relates to other entities. The RDF data model is a directed labeled graph, in which nodes are Resources or Literals, and edges are Properties. Properties are binary relations between entities. They are defined by 'the set of classes

that they apply to' called domain and 'class of entities from which it takes its values called range'. An RDF statement is a triple format of *Subject-Property-Object*. The subject of the statement has a property whose value is the object of the statement. An object can be either another resource or a literal value. A statement can be represented as a graph, by drawing an arc from a node representing the subject to another node representing the object, *Subject → Object*. The arc is labeled with the name of the property and resource nodes are labeled with the URIs of the resources. A special property rdf:typeOf links resources to the classes they belong to. The classes and properties are described in an RDFS (RDF Schema)[4]. RDFS provides a general schema contains classes and properties for describing domain specific vocabularies.

### Semantic associations

Semantic associations are complex relationships between resource entities [5]. Most of the useful semantic associations involve some intermediate entities and relations. It helps the user to establish the connection between different people, places and events. To describe the semantic association, we provide definitions regarding the formalization of semantic association adapted from Anyanwu et al. [6], below:

*Definition 1(Semantic Connectivity)*

Two entities $e_1$ and $e_n$ are semantically connected if there exists a sequence $e_1$, $P_1$, $e_2$, $P_2$, $e_3$, $P_{3\ldots}$ $e_{n-1}$, $P_{n-1}$, $e_n$ in an RDF graph where $e_i$ ($1 \le i \le n$) are entities and $P_j$ ($1 \le j < n$) are properties. Figure 2. Shows the semantic connectivity between $e_1$ and $e_n$.

*Definition 2 (Semantic Similarity)*

Two entities $e_1$ and $f_1$ are semantically similar if there exist two semantic paths $e_1$, $P_1$, $e_2$, $P_2$, $e_3$, $P_{3\ldots}$ $e_{n-1}$, $P_{n-1}$, $e_n$ and $f_1$, $Q_1$, $f_2$, $Q_2$, $f_3$, $Q_{3\ldots}$ $f_{n-1}$, $Q_{n-1}$, $f_n$ semantically connecting $e_1$ with $e_n$ and $f_1$ with $f_n$ respectively, and that for every pair of properties $P_i$ and $Q_i$, $1 \le i < n$, either of the following conditions holds: $P_i = Q_i$ or $P_i \subseteq Q_i$ or $Q_i \subseteq P_i$ ($\subseteq$ means rdf:subPropertyOf), then two paths originating at $e_1$ and $f_1$, respectively, are semantically similar.

*Definition 3 (Semantic Association)*

Two entities $e_x$ and $e_y$ are semantically associated if $e_x$ and $e_y$ are semantically connected or semantically similar.

### Related work

Several techniques have been proposed related to ranking of semantic associations. Some of them are summarized below:
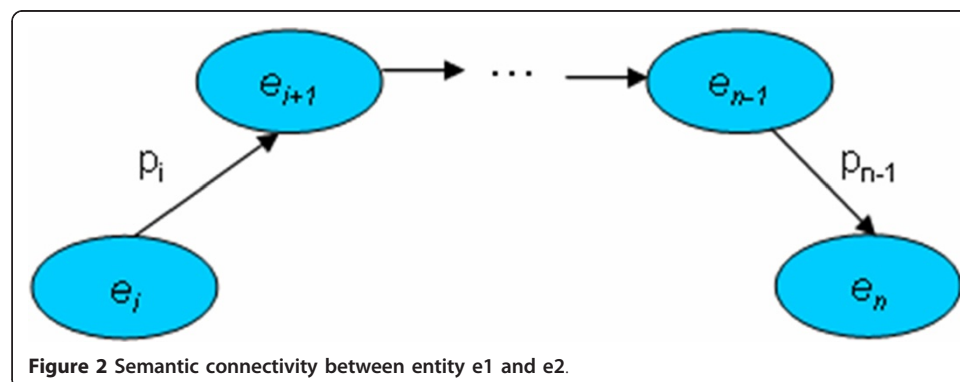


**Figure 2 Semantic connectivity between entity e1 and e2**.

Personalized semantic web search proposed by Jiang, A et al. [7] discuss a set of statistical methods to learn a user ontology from a given domain ontology and a spreading activation procedure for inferencing in the user ontology. The user ontology model with the spreading activation based inferencing procedure has been incorporated into a semantic search engine to provide personalized document retrieval services. In our method, relationships between two concepts capture in RDF data model and ranking the relationships based on the users' need with help of various metrics.

Stojanovic et al. [8] discuss an approach for semantic ranking of results which is based on estimating the distance between concepts in the query and concepts in the result, and assigning more relevance to results using closer concepts to query concepts.

For ranking the results of complex relationship searches on the semantic web, Anyanwu et al. [6] present a flexible approach called SemRank. In this method, with the help of a sliding bar user can easily vary their search mode from conventional search mode to discovery search mode.

Shahdad Shariatmadari et al. [9] present a method for finding semantic association based on the concept semantic similarity. $\rho$-operator [5] is used for discovering semantic similarities and graph similarity approach [10] is used to rank the similarity. The similarity between two paths will be calculated based on the degree of similarity of the nodes and edges using subsumption function proposed by Aleman-Meza et al. [11]. The ranking approach proposed by Anyanwu, K. et.al [5] considers 'context' based on value assignments for different ontologies.

Aleman-Meza et al. [11] discuss a frame work that uses ranking techniques to identify more interesting and more relevant semantic associations and define a ranking formula that considers Subsumption Weight $S_P$ (how much meaning a semantic association conveys depending on the places of its components in the RDF), Path Length Weight $L_P$ (that allows preference of either immediate or distant relationships), Popularity $P_P$ (number of incoming and outgoing edges), Rarity $R_P$ (rarely occurring entity) and Trust Weight $T_P$ (determining how reliable a relationship is according to its origin) and context weight, for assessing the effectiveness of the ranking scheme. In this method 'user defined weight' is assigned for each 'context regions' specified by the user and it is used to calculate the context weight. The ranking results depend on the criteria defined by the user.

Lee, M et al. [12] propose a semantic search methodology for measuring the information content of a semantic association that consists of resources and properties based on information theory and expanding the semantic network based on spreading activation. In this method, they provide search results that are connected and ordered relations between search keyword and other resources as link of relation on semantic network.

Xiao Dong et al. [13] present a prototype system called Chem2Bio2RDF Dashboard for automatic collecting semantic association within the systems chemical biology space and apply a series of ranking metrics called Quality, Specificity and Distinctiveness to select the most relevant association.

To discover semantic associations between linked data, Vidal, M et al. [14] propose an authority-flow based ranking technique that is able to assign high scores to terms that correspond to potential discoveries, and to efficiently identify these highly scored terms. They also propose an approximate solution named graph-sampling. This technique samples events in a Bayesian network that models the topology of the data

connections and it estimates ranking score that measure how important and relevant are the associations between two terms.

*In summary, while many have developed methods for ranking semantic association, most existing methods only consider relationships between two entities and there is no provision has given to the users to select the intermediate entities if they want. In our method, we have included a provision to select the intermediate entities by the discovery process that finds the semantic association paths between two entities. In order to reduce the search space we have used Modified Bidirectional BFS algorithm. This algorithm filters the irrelevant paths with respect to the user-specific intermediate entity. Then the paths are ranked based on the various statistical and semantic weights* [11]. Since Aleman-Meza et al. [11] propose a method for ranking the semantic association by considering various parameters that are influence the relevancy of the paths and also provide a flexible ranking according to the users' need, we adapt this method for ranking the semantic association paths.

### Discovering semantic association paths

To find the paths between two entities passing through intermediate entities, we have carried out the Bidirectional Breadth-First-Search algorithm with some modifications. Bidirectional Breadth-First-Search is a graph search algorithm that finds a shortest path from source node to destination node in a directed graph. Bi-directional BFS algorithm searches the paths in both forward and backward directions. When these searches meet in the middle, the algorithm stops and return the complete paths.

In our proposed Modified bidirectional BFS algorithm shown in Table 1, user has to enter the intermediate entity along with the source and destination entities. The algorithm returns the semantic association paths only if both forward and backward search meet in the middle and path pass through the user-specific intermediate entity. We have assigned different colours to identify the examined entity for both forward and backward search. When the forward search meets the entity which is already examined by the backward search, we can say that these two searches meet in the middle. When we examine the entity, we set the flag value as 1, if it meets the user-specific intermediate entity. The flag value has to reset as 0, when the new search starts. It helps to identify whether the path passes through the user-specific intermediate entity or not.

The following is Modified bidirectional breadth-first search algorithm to find the semantic association between two entities which pass through an intermediate entity.

Modified bidirectional BFS algorithm returns all the semantic association paths between source and destination entities which pass through user-specific intermediate entity. Once the paths have been discovered from the RDF, we have to rank these paths according to the user's needs in order to improve the relevancy of the results.

### Ranking the semantic association paths

Our approach defines a path rank, as a function of various intermediate weights [11]. These are described as follows:

**Context Weight $C_P$** is one of the semantic metrics which is used to determine the relevancy based on an user specific view. Consider the scenario in which someone is interested in discovering how two persons are related to each other in the

**Table 1 Modified bidirectional breadth-first search algorithm**

| Modified Bidirectional Breadth-First Search Algorithm |
| --- |
| Input: RDF Graph and user query as entity |
| Output: Semantic Association paths without ranking |
| /**Path searching for finding all paths at within k hops */ |
| 1. Enqueue the root node in the forward queue: the source Thing in the path we want to discover |
| 2. Enqueue the target node in the backward queue: the destination Thing in the path we want to discover |
| 3. Dequeue a forward node and examine it |
|     a. If the forward node is found in the backward zone//we have met in the middle: quit the search |
|         a.1 If the path passes through user specific intermediate node return the complete path |
|         a.2 quit the search |
|     b. Otherwise enqueue any successors that have not yet been examined in the forward queue and add them to the forward zone |
| 4. Dequeue a backward node and examine it |
|     a. If the backward node is found in the forward zone//we have met in the middle |
|         a.1 If the path passes through user specific intermediate node return the complete path |
|         a.2 quit the search |
|     b. Otherwise enqueue any predecessors that have not yet been examined in the backward queue and add them to the backward zone |
| 5. If the any of the two queues is empty, every node on the graph has been examined, quit the search and return "not found" |
| 6. If we are still under the depth limit, repeat from Step 3 |

domain of "Funding Company". Concepts such as "Finance" or "Financial organization" would be most relevant, whereas something like "Music Company" would be less meaningful. So it is possible to capture a user's interest through a Context Specification with help of user interface screen. Thus, using the context specified, it is possible to rank the semantic association paths according to its relevance with a user's domain of interest.

*Subsumption weigh* $S_P$: In RDF, entities that are in the lower hierarchy can be considered to be more specialized instances of those further up in the hierarchy. Thus, lower entities have more specific meaning. So, high relevance can be assigned based on subsumption.

*Path Length Weight* $L_P$: In some queries, a user may be interested in the shortest paths. This may infer a strong relationship between two entities. In some cases a user may wish to find indirect or longer paths. Hence, the user can determines which association length influence.

*Popularity* $P_P$: The number of incoming and outgoing relationships of entities called popular entities. Path contains highly popular entities may be more relevant.

*Rarity* $R_P$: Sometimes rarely occurring events are considered to be more interesting than commonly occurring ones [15,16].

*Trust Weight* $T_P$: Various entities and their relationships in semantic associations originate from different sources. Some of these sources may be more trusted than others. Thus, trust values need to be assigned to the meta-data extracted depending on its source.

We calculate Subsumption Weight $S_P$, Path Length Weight $L_P$, Popularity $P_P$, Rarity $R_P$ and Trust Weight $T_P$ [15] along with context weight $C_P$. These weights are used to

determine the path relevancy. So, all the intermediate weights are added to calculate the rank of the each path.

Overall association Rank is calculated using the criteria as

$$W_p = k_1 \times S_p + k_2 \times L_p + k_3 \times C_p + k_4 \times T_p + k_5 \times P_p + k_6 \times R_p \tag{1}$$

where $k_i (1 \leq i \leq 6)$ are preference weights and $\Sigma\ k_i = 1$. The resulting paths are ranked based on the users' domain of interest. Depending on the requirements, users can also change the preference weights to fine-tune the ranking criteria. In our experiments, we have given high weights to context component and use the other ranking components as secondary criteria.

## Experimental evaluation

For finding semantic association paths, we have used an RDF encompassing 30 classes, 50 properties and 2000 entities covering various domain such as 'Movie', 'Music', 'Finance' etc. To test the performance of our system, we have tested with 30 pairs of entities in the RDF. Semantic association paths weights have been calculated and ranked under various criteria such as favor short association or favor long association, favor popularity entities or favor unpopular entities and the favor rarity. Criteria have been selected through user interface. Semantic association paths ranking has been done by the above users through the system as well as manually.

## User interface

User interface for the system is a web based application using Servlet and Apache Tomcat. Using this interface, users can specify source, destination entities and one or more intermediate entities. The system will return ranked Semantic Associations between these entities specified by users, which are passing through intermediate entities. Snapshot of the ranked results of a sample query is given in Figure 3.
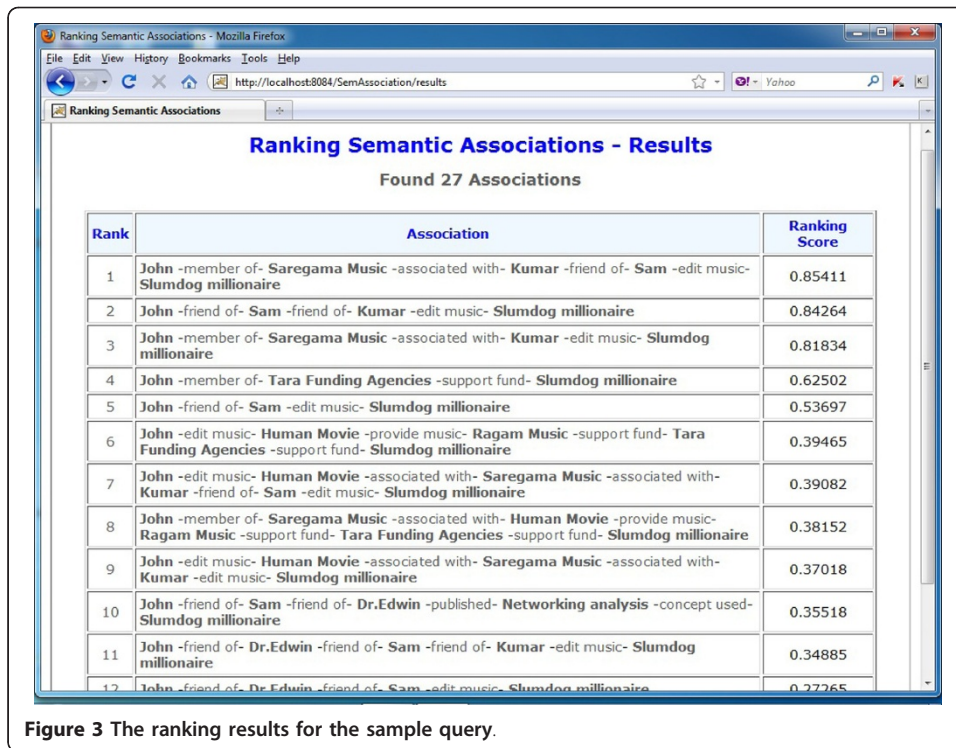
In order to demonstrate the effectiveness of the searching method, we have compared the execution time for searching the paths between two entities which are passing through user-specific intermediate entity. The experiments were carried out on a desktop PC running Windows 7 with dual Intel® Core™2, 2.40 Ghz CPU and 2 GB memory. Figure 4 shows the comparison of path searching time between bidirectional BFS and Modified bidirectional BFS algorithm.

When the number of RDF triples increases the execution time of both bidirectional BFS method and Modified bidirectional BFS method increases as well. Since our proposed algorithm reduces the searching of the paths that are irrelevant to the user, it always takes less execution time than bidirectional BFS method. According to our experiment our Modified bidirectional BFS algorithm, searches the paths efficiently.

To demonstrate the ranking scheme's effectiveness, we have used the Spearman's Foot rule [17] distance as the measure of similarity between proposed system ranking and user-human ranking. The formula is given below:
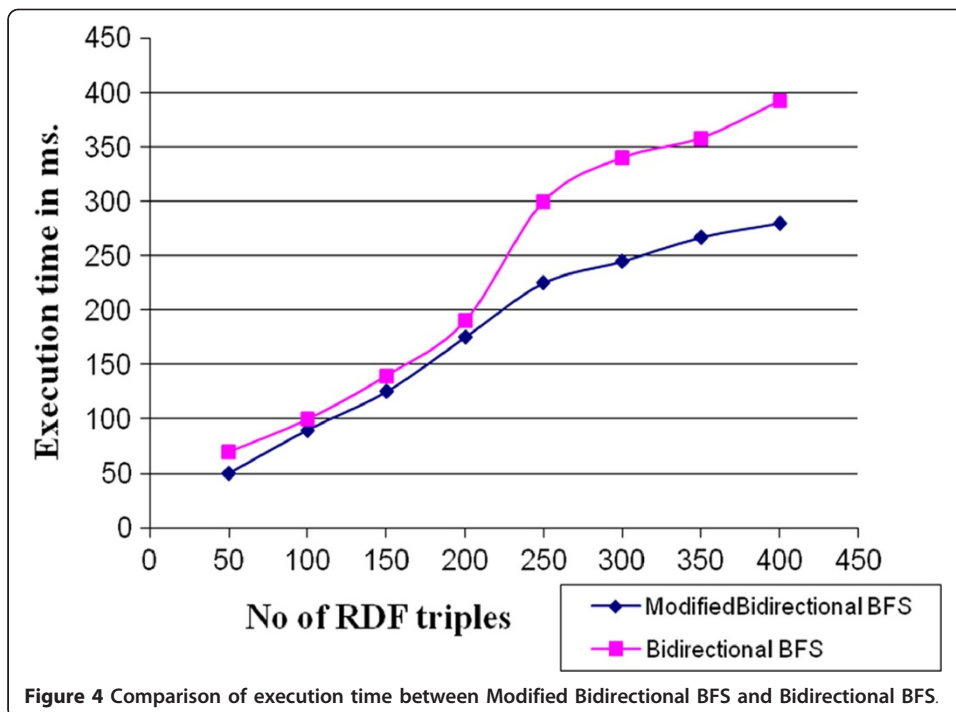
### Spearman's Foot rule distance

$$D_{(system,human)} = \sum_{i=l}^{n} \left| R_{i_{system}} - R_{i_{human}} \right| \tag{2}$$

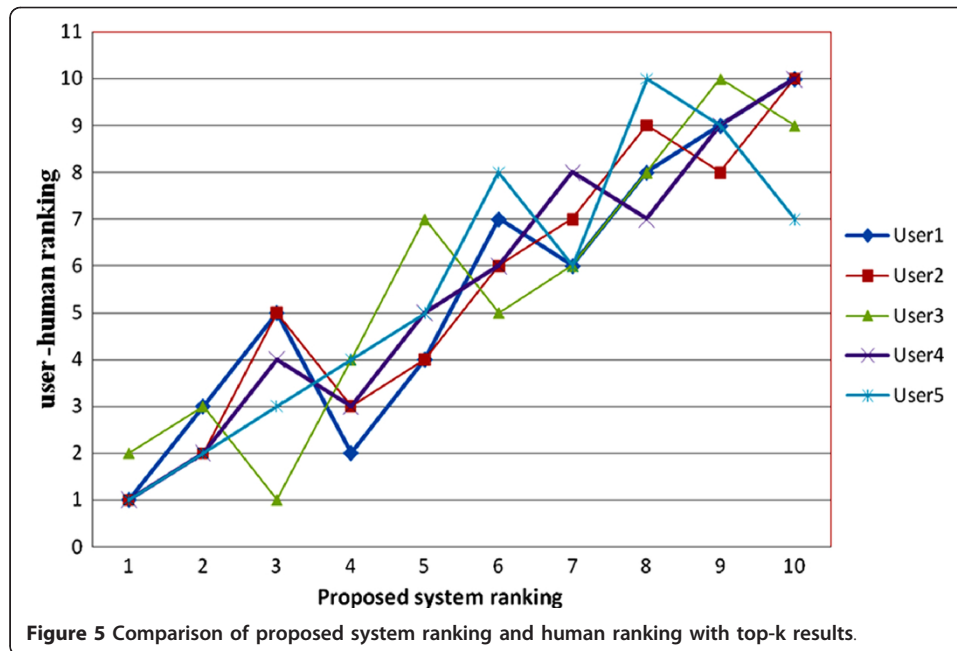**Figure 3 The ranking results for the sample query**.

$$\text{Spearman's Foot rule Coefficient} \quad C = 1 - \frac{4D}{n^2} \tag{3}$$

Figure 5 shows comparison of human ranking and proposed system ranking results between the entity sets (Entity1: 'John' and Entity2: 'Slumdog millionaire') for the 5
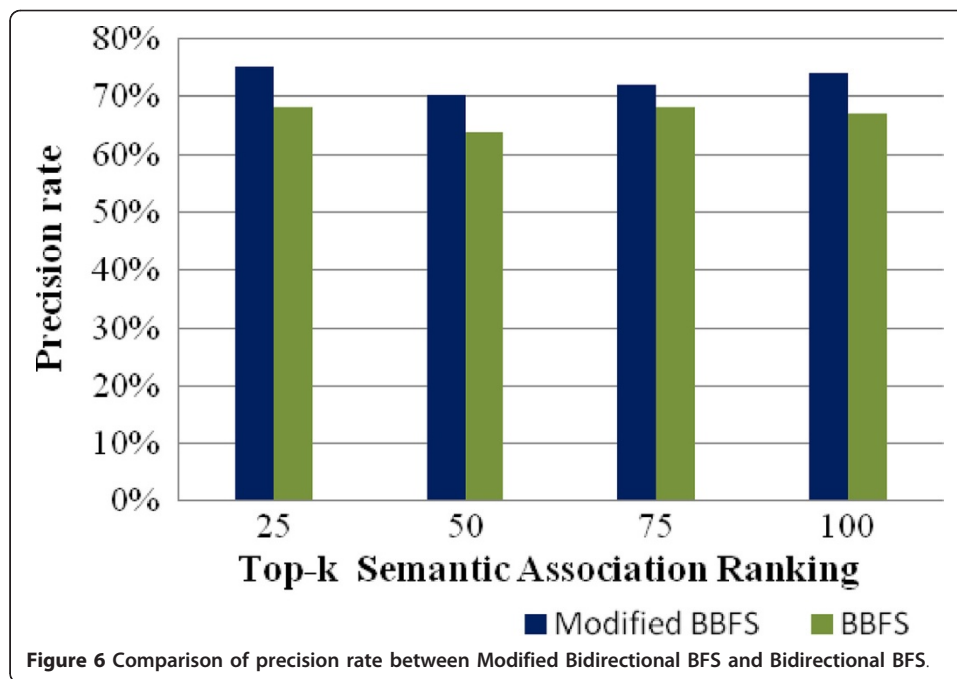


**Figure 4 Comparison of execution time between Modified Bidirectional BFS and Bidirectional BFS**.

**Figure 5 Comparison of proposed system ranking and human ranking with top-k results**.

users. Here 'Ragam music' is an intermediate entity. The x-axis represents semantic associations rank first, second and so on according to the proposed system results. The y-axis represents user-human ranking which is assigned manually by the users.

We have conducted this experiment with 50 users. In our experiment, the average correlation coefficients between the proposed system ranking and the user-human's ranking is 0.69. As the average correlation coefficient is greater than 0.5, the proposed system's ranking and user-human ranking are highly correlated. Therefore, the proposed system ranks the semantic association paths with more relevance to the user.

To measure the relevancy rate to the user, we have calculated the precision rate for top-k semantic association ranking. Precision represents the fraction of the relevant paths from top-k semantic association paths. The Figure 6 shows the comparison of average precision rate of top-k semantic associations ranking between Modified Bi-directional BFS approach and BBFS approach. Since the paths which are irrelevant to the users are filtered by the Modified Bi-directional BFS algorithm during the discovery of the paths, the precision rate has improved in ranking the semantic association paths. It explains that our proposed method provide high precision rates than existing method.

## Conclusion

A semantic association is a set of relationships between two entities in a knowledge base represented as graph paths consisting of a sequence of links. The number of relationships between entities in a knowledge base might be much greater than the number of entities. Ranking these relationship paths are required to find the relevant relationships between entities with respect to user's interest. Sometimes users' may expect the relationships between two entities with respect to other intermediate entity. In this paper, we have used Modified bidirectional BFS algorithm to discover the paths between two entities along one or more user specific intermediate entities and finally we rank the semantic association paths according to the users' needs. We compare the

**Figure 6 Comparison of precision rate between Modified Bidirectional BFS and Bidirectional BFS.**

execution time between bidirectional BFS and Modified bidirectional BFS, to search these paths. Based on the experiments, our searching method, retrieve the path efficiently. We have evaluated the correlation between our method and existing method through Spearman's correlation coefficient. The average correlation coefficient between proposed system rankings and human ranking is 0.69. It explains that our proposed system ranking is highly correlated with human ranking. Also, our method always provide higher precision rate for top-k ranking. In future, we will improve our method for discovering semantic association to rank the result in more accurate way according to the specific users' interest. For this purpose we plan to generate the semantic web usage ontology from web usage information of each user and it may be used to get personalized semantic associations ranking.

### Endnotes
[1]Resource Description Framework

**Author details**
[1]Department of Computer Applications, Sri Krishna College of Engineering and Technology, Tamil nadu, India
[2]Department of Computer Science and Engineering, Sri Krishna College of Engineering and Technology, Tamil nadu, India

**Authors' contributions**
The described approach was developed through discussions collectively by the both authors. VV has implemented the graphical user interface for the testing environment. Both VV and IK involved in the manuscript drafting and revising. All authors read and approved the final manuscripts.

**Competing interests**
The authors declare that they have no competing interests.

**References**
1. Berners-Lee T, Hendler J, Lassila O (2001) The Semantic Web: a new form of web content that is meaningful to computers will unleash a revolution of new possibilities. Sci Am 285(5):34–43. doi:10.1038/scientificamerican1101-34.

2.  Klasnja-Milicevic A, Vesin B, Ivanovic M, Budimac Z (2011) E-Learning personalization based on hybrid recommendation strategy and learning style identification. Comput Educ 56:885–899. doi:10.1016/j.compedu.2010.11.001.
3.  Lassila O, Swick RR (1999) Resource Description Frame work(RDF) Model and syntax specification, W3C Recommendation. http://www.w3.org/TR/PR-rdf-syntax/
4.  Brickley D, Guha RV (2000) Resource Description Framework (RDF) Schema Specification 1.0, W3C Candidate Recommendation. http://www.w3.org/TR/rdf-schema/
5.  Anyanwu K, Sheth A (2003) ρ -Queries: Enabling Querying for Semantic Associations on the Semantic web. Proc of the 12th International World Wide Web Conference. ACM Press  pp 690–699
6.  Anyanwu K, Maduko A, Sheth A (2005) SemRank: Ranking Complex Relationship Search Results on the Semantic Web. Proc of the 14th International World Wide Web Conference. ACM Press  pp 117–127
7.  Jiang X, Tan A (2009) Learning and inferencing in user ontology for personalized Semantic Web search. Inform Sci 179:2794–2808. doi:10.1016/j.ins.2009.04.005.
8.  Stojanovic N, Mädche A, Staab S, Studer R (2001) Sure Y, SEAL: a framework for Developing SEmantic PortALs. Proceedings of K-CAP 155–162
9.  Shariatmadari S, Mamat A, Ibrahim H, Mustapha N (2008) SwSim: Discovering semantic similarity association in semantic web. Proc of the International Symposium on IT Sim 1–4
10. Sokolsky O, Kannan S, Lee I (2006) Simulation - Based Graph Similarity. Proc of 12th International Conference on Tools and Algorithms for the Construction and Analysis of Systems. Springer-verlag LNCS 3920:426–440
11. Aleman-Meza B, Halaschek C, Arpinar IB, Sheth A (2005) Ranking Complex Relationships on the Semantic Web. IEEE Internet computing 9(3):37–44. doi:10.1109/MIC.2005.63.
12. Lee M, Kim W (2009) Semantic Association Search and Rank Method based on Spreading Activation for the Semantic Web. Proc of IEEE International Conference on Industrial Engineering and Engineering Management 1523–1527
13. Dong X, Ding Y, Wang H, Chen B, Wild DJ (2010) Ranking semantic associations in systems chemical biology space, 19th International World Wide Web Conference. FWCS, Raleigh
14. Vidal M, Rashid L, Ibabez L, Rivera J, Rodrogiez H, Ruckhaus E (2010) A Ranking-Based Approach to Discover Semantic Association Between Linked Data. The 2nd International Workshop on Inductive Reasoning and Machine Learning for the Semantic Web 611:18–29
15. Anderson R, Khattak A (1998) The Use of Information Retrieval Techniques for Intrusion Detection. Proc 1st International Workshop Recent Advances in Intrusion Detection http://www.cl.cam.ac.uk/ftp/users/rja14/raid.ps.gz
16. Lin S, Chalupsky H (2003) Unsupervised Link Discovery in Multi-Relational Data via Rarity Analysis. Proc 3rd IEEE Int'l Conf. Data mining. IEEE CS Press  pp 171–178
17. Diaconis P, Graham R (1977) Spearman's Foot rule as a Measure of Disarray. J Royal Statistical Soc, Series B 39(no.2):262–268