

## Research Article

# pSPARQL: A Querying Language for Probabilistic RDF Data

**Hong Fang** 

*College of Arts and Sciences, Shanghai Polytechnic University, Shanghai 201209, China*

Correspondence should be addressed to Hong Fang; [fanghong@sspu.edu.cn](mailto:fanghong@sspu.edu.cn)

Received 19 December 2018; Accepted 19 February 2019; Published 26 March 2019

Guest Editor: Jianxin Li

Copyright © 2019 Hong Fang. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

More and more linked data (taken as knowledge) can be automatically generated from nonstructured data such as text and image via learning, which are often uncertain in practice. On the other hand, most of the existing approaches to processing linked data are mainly designed for certain data. It becomes more and more important to process uncertain linked data in theoretical aspect. In this paper, we present a querying language framework for probabilistic RDF data (an important uncertain linked data), where each triple has a probability, called pSRARQL, built on SPARQL, recommended by W3C as a querying language for RDF databases. pSPARQL can support the full SPARQL and satisfies some important properties such as well-definedness, uniqueness, and some equivalences. Finally, we illustrate that pSPARQL is feasible in expressing practical queries in a real world.

## 1. Introduction

Resource Description Framework (RDF) [1] is the standard data model in the Semantic Web. In our real world, RDF data (as a knowledge base) possibly contains some uncertainty data due to the diversity of data sources, where RDF data are automatically extracted from different sources, such as YAGO [2]. For instance, some RDF data is generated from raw data via knowledge extraction and machine learning [3]. Indeed, uncertainty is generally a basic feature of data [4–6]. However, RDF model itself provides little support for uncertain data [7]. SPARQL [8], as a querying language for RDF data officially recommended by W3C [9], is unable to process uncertain data [4].

There are many approaches to processing probabilistic RDF [10]. Reference [4] proposes a probabilistic model for SQL over relational data. Reference [5] presents a Bayesian network to represent probabilistic relations in RDF. Reference [11] develops a framework for evaluating SPARQL conjunctive queries (i.e., basic graph patterns, BGP) on RDF probabilistic databases. Reference [12] proposes answering SPARQL queries with RDFS reasoning on probabilistic models that encode statistical relationships among correlated triples, where the proposed probability models are based on either probability distribution function or a disjunctive normal form probability problem. Reference [13] presents

effective pruning mechanisms, as well as structural and probabilistic pruning for query answering of SPARQL conjunctive queries (i.e., BGP) over probabilistic RDF data graphs. Reference [14] presents a RESCAL-based approach to query processing in relational data via factorization. Reference [14] presents a heuristic algorithm for query answering of SPARQL conjunctive queries (i.e., BGP) over incomplete and uncertain RDF. Reference [15] presents a framework for SPARQL query answering over probabilistic databases by extending the rich semantics offered by ontologies with probabilistic information. Reference [16] presents a probabilistic knowledge base system, ARCHIMEDESONE, for query answering with inference by scaling up the knowledge expansion and statistical inference algorithms. Reference [17] proposes a probabilistic automata-based framework of query evaluation in the presence of uncertainty efficiently.

Although those approaches can query probabilistic RDF, most of them mainly process SPARQL conjunctive queries, that is, BGP queries. However, those existing probability models have little support for expressive operators (for instance, neither [13] nor [16] discusses OPTIONAL query for RDF) such as OPTIONAL, which is the least conventional operator of SPARQL [18], and DIFF, a difference operator in SPARQL 1.1 [19], which brings more expressivity [20].

In this paper, we present an extended querying language (called pSPARQL: *probabilistic SPARQL*) for probabilistic

RDF databases with support of the full SPARQL fragment. We show that the semantics of pSPARQL can satisfy some important properties such as well-definedness, uniqueness, and some equivalences. Compared with the previous poster in ISWC 2016 [21], in this paper, we present a totally new probabilistic representation model and prove that the newly proposed model can preserve some important properties such as uniqueness and distributive law of equivalence.

The remainder of this paper is structured as follows: the next section recalls RDF and SPARQL. Section 3 introduces the syntax and semantics of pSPARQL and Section 4 discusses some important properties. Finally, we summarize our work in the last section.

## 2. RDF and SPARQL

In this section, we briefly recall the syntax and semantics of SPARQL. For more readings, please refer to the core SPARQL formalization in [22].

*2.1. RDF Graphs.* Let  $I$  and  $L$  be infinite sets of IRIs and literals, respectively, with  $I \cap L = \emptyset$ . Let  $U = I \cup L$ . A triple  $(s, p, o) \in U \times I \times U$  is called an *RDF triple*. An *RDF graph* is a finite set of RDF triples.

*2.2. Syntax of SPARQL.* Let  $V$  be a set of variables. SPARQL patterns are inductively defined as follows:

- (i) Any triple from  $(U \cup V) \times (I \cup V) \times (U \cup V)$  is a pattern (called a *triple pattern*).
- (ii) If  $P_1$  and  $P_2$  are patterns, then so are the following:  $P_1$  UNION  $P_2$ ,  $P_1$  AND  $P_2$ ,  $P_1$  DIFF  $P_2$ , and  $P_1$  OPT  $P_2$ .
- (iii) If  $P$  is a pattern and  $F$  is a constraint (defined next), then  $P$  FILTER  $F$  is a pattern; we call  $F$  the *filter*, which is a Boolean combination of *atomic constraints*, one of the three following forms:  $\text{bound}(?x)$  (*bound*),  $?x = ?y$  (*equality*), and  $?x = c$  (*constant equality*), for  $?x, ?y \in V$  and  $c \in U$ .

*2.3. Semantics of SPARQL.* Now, given a graph  $G$  and a pattern  $P$ , we define the semantics of  $P$  on  $G$ , denoted by  $\llbracket P \rrbracket_G$ , as a set of mappings (i.e., partial functions from  $V$  to  $U$ , in the following manner, where we use  $\text{dom}(\mu)$  to denote the domain of  $\mu$ )

- (i)  $\llbracket (u, v, w) \rrbracket_G := \{\mu : \{u, v, w\} \cap V \rightarrow U \mid (\mu(u), \mu(v), \mu(w)) \in G\}$ .
- (ii)  $\llbracket P_1 \text{ UNION } P_2 \rrbracket_G := \llbracket P_1 \rrbracket_G \cup \llbracket P_2 \rrbracket_G$ .
- (iii)  $\llbracket P_1 \text{ AND } P_2 \rrbracket_G = \{\mu_1 \cup \mu_2 \mid \mu_1 \in \llbracket P_1 \rrbracket_G \text{ and } \mu_2 \in \llbracket P_2 \rrbracket_G \text{ and } \mu_1 \sim \mu_2\}$ .

Here, two mappings  $\mu_1$  and  $\mu_2$  are called *compatible*, denoted by  $\mu_1 \sim \mu_2$ , if for any

$$?x \in \text{dom}(\mu_1) \cap \text{dom}(\mu_2) = \emptyset, \mu_1(?x) = \mu_2(?x).$$

- (iv)  $\llbracket P_1 \text{ DIFF } P_2 \rrbracket_G = \{\mu_1 \in \llbracket P_1 \rrbracket_G \mid \neg \exists \mu_2 \in \llbracket P_2 \rrbracket_G, \mu_1 \sim \mu_2\}$ .
- (v)  $\llbracket P_1 \text{ OPT } P_2 \rrbracket_G = \llbracket (P_1 \text{ AND } P_2) \text{ UNION } ((P_1 \text{ DIFF } P_2)) \rrbracket_G$ .

$$(vi) \llbracket P_1 \text{ FILTER } P_2 \rrbracket_G := \{\mu \in \llbracket P_1 \rrbracket_G \mid \mu(F) = true\}.$$

Here, for any mapping  $\mu$  and filter  $F$ , the evaluation of  $F$  on  $\mu$ , denoted by  $\mu(F)$ , is defined in terms of a three-valued logic with truth values *true*, *false*, and *error*. Recall that

$F$  is a Boolean combination of atomic constraints.

For a bound constraint  $\text{bound}(?x)$ , we define

$$\mu(\text{bound}(?x)) = \begin{cases} true & \text{if } ?x \in \text{dom}(\mu); \\ false & \text{otherwise.} \end{cases} \quad (1)$$

For an equality constraint  $?x = ?y$ , we define

$$\begin{aligned} & \mu(?x = ?y) \\ &= \begin{cases} true & \text{if } ?x, ?y \in \text{dom}(\mu) \text{ and } \mu(?x) = \mu(?y); \\ false & \text{if } ?x, ?y \in \text{dom}(\mu) \text{ and } \mu(?x) \neq \mu(?y); \\ error & \text{otherwise.} \end{cases} \end{aligned} \quad (2)$$

Thus, when  $?x$  and  $?y$  do not both belong to  $\text{dom}(\mu)$ , the equality constraint evaluates to *error*. Similarly, for a constant-equality constraint  $?x = c$ , we define

$$\begin{aligned} & \mu(?x = c) \\ &= \begin{cases} true & \text{if } ?x \in \text{dom}(\mu) \text{ and } \mu(?x) = c; \\ false & \text{if } ?x \in \text{dom}(\mu) \text{ and } \mu(?x) \neq c; \\ error & \text{otherwise.} \end{cases} \end{aligned} \quad (3)$$

A Boolean combination is evaluated using the truth tables given in Table 1.

## 3. Probabilistic RDF and pSPARQL

In this section, we present probabilistic RDF and introduce the syntax and semantics of pSPARQL.

*3.1. Probabilistic RDF.* A *probabilistic RDF*  $R$  is a pair  $(G, \rho)$  where  $G$  is an RDF graph and  $\rho$  is a total function from  $G \rightarrow (0, 1]$ . Intuitively speaking,  $\rho$  is a probability function mapping each triple to a probability.

For instance, let  $R_{\text{John}} = (G, \rho)$  be a probabilistic RDF with  $G = \{t_1, t_2, t_3\}$  and  $\rho$  is a function from  $G \rightarrow [0, 1]$  defined in Table 2.

Note that we assign a triple to a probability so that we could take triples as atoms in our scenario analogously treated in [13, 16]. This treatment is not direct to characterize the probability of subjects/objects in triples.

*3.2. pSPARQL: A Probabilistic SPARQL.* In this section, we introduce a probabilistic SPARQL (for short, pSPARQL).

*The Syntax of pSPARQL.* The syntax of pSPARQL is slightly different from the syntax of SPARQL [22] in filters, where we

newly introduce a fixed variable  $?p$  to express constraints of probability.

A probabilistic atomic filter is one of the four following forms:  $?p \geq s$ ,  $?p \leq s$ ,  $?p = s$  and  $?p \neq s$ , where  $s \in (0, 1]$ . The filter of pSPARQL is a Boolean combination of atomic filters and probabilistic atomic filters.

All patterns are called probabilistic patterns in pSPARQL.

*The Semantics of pSPARQL.* The semantics of probabilistic patterns are defined in terms of sets of pairs of the form  $(\mu, p)$  (called a solution (with probability), denoted by  $\tau$ ), where  $\mu$  is a solution of probabilistic patterns and  $p \in (0, 1]$ . Note that we only consider pairs of form  $(\mu, p)$  where  $p > 0$ .

Now, given a probabilistic RDF graph  $R = (G, \rho)$  and a probabilistic pattern  $P$ , we define the semantics of  $P$  on  $R$ , denoted by  $\llbracket P \rrbracket_R$ , as a set of solutions with probability, in the following manner:

$$\begin{aligned} \llbracket (u, v, w) \rrbracket_R &:= \{(\mu, p) : \{u, v, w\} \cap V \\ &\longrightarrow U \mid (\mu(u), \mu(v), \mu(w)) \in G \text{ and } p \\ &= \rho(\mu(u), \mu(v), \mu(w))\}. \\ \llbracket P_1 \text{ UNION } P_2 \rrbracket_R &:= \{(\mu, p) \mid (\mu, p_1) \\ &\in \llbracket P_1 \rrbracket_R \text{ or } (\mu, p_2) \in \llbracket P_2 \rrbracket_R \text{ and } p \\ &= \max\{p_1, p_2\}\}. \end{aligned} \quad (4)$$

By default, we set  $p = \max\{p\}$ .

$$\begin{aligned} \llbracket P_1 \text{ AND } P_2 \rrbracket_R &:= \{\mu_1 \cup \mu_2, p \mid (\mu_i, p_i) \\ &\in \llbracket P_i \rrbracket_R \ (i = 1, 2) \\ \text{and } p &= \max_{\mu_1 \cup \mu_2 = \mu_i \cup \mu_j} \min\{p_i, p_j \mid (\mu_i, p_i) \\ &\in \llbracket P_1 \rrbracket_R \text{ and } (\mu_j, p_j) \in \llbracket P_2 \rrbracket_R \text{ and } \mu_i \sim \mu_j\} \\ \llbracket P_1 \text{ DIFF } P_2 \rrbracket_R &:= \{(\mu_1, p_1) \in \llbracket P_1 \rrbracket_R \mid \neg \exists (\mu_2, p_2) \\ &\in \llbracket P_2 \rrbracket_R \text{ s.t. } \mu_1 \sim \mu_2\}. \end{aligned} \quad (5)$$

$$\begin{aligned} \llbracket P_1 \text{ OPT } P_2 \rrbracket_R & \\ &= \llbracket (P_1 \text{ AND } P_2) \text{ UNION } (P_1 \text{ DIFF } P_2) \rrbracket_R. \\ \llbracket P_1 \text{ FILTER } F \rrbracket_G &:= \{\tau = (\mu, p) \in \llbracket P_1 \rrbracket_G \mid \tau(F) \\ &= \text{true}\}. \end{aligned}$$

(i) For a nonprobabilistic filter  $F$ ,  $\tau(F) = \mu(F)$ .

(ii) For a Boolean combination  $F$ ,  $\tau(F) = \mu(F)$ .

(iii) For a probabilistic filter  $?p \geq s$ , we define

$$\tau(?p \geq s) = \begin{cases} \text{true} & \text{if } p \geq s; \\ \text{false} & \text{otherwise.} \end{cases} \quad (6)$$

TABLE 1: Truth tables for the three-valued semantics.

(a)			
$p$	$q$	$p \wedge q$	$p \vee q$
true	true	true	true
true	false	false	true
true	error	error	true
false	true	false	true
false	false	false	false
false	error	false	error
error	true	error	true
error	false	false	error
error	error	error	error

(b)	
$p$	$\neg p$
true	false
false	true
error	error

(iv) For a probabilistic filter  $?p \leq s$ , we define

$$\tau(?p \leq s) = \begin{cases} \text{true} & \text{if } p \leq s; \\ \text{false} & \text{otherwise.} \end{cases} \quad (7)$$

(v) For a probabilistic filter  $?p = s$ , we define

$$\tau(?p = s) = \begin{cases} \text{true} & \text{if } p = s; \\ \text{false} & \text{otherwise.} \end{cases} \quad (8)$$

(vi) For a probabilistic filter  $?p \neq s$ , we define

$$\tau(?p \neq s) = \begin{cases} \text{true} & \text{if } p \neq s; \\ \text{false} & \text{otherwise.} \end{cases} \quad (9)$$

*Example 1.* Given a pattern  $P = (?x, \text{sufferFrom}, ?y)\text{FILTER } ?p \geq 0.5$  (i.e., we query those persons who have suffered from some illness with probability over 0.5), we can compute that  $\llbracket P \rrbracket_{R_{\text{john}}} = \{(\mu, p)\}$  where  $\mu = \{?x \rightarrow \text{John}, ?y \rightarrow \text{Schizophrenia}\}$  and  $p = 0.84$ . However, let  $\tau = (\{?x \rightarrow \text{John}, ?y \rightarrow \text{Schizophrenia}\}, 0.32)$ ,  $\tau \in \llbracket P \rrbracket_{R_{\text{john}}}$  since  $\tau(?p \geq 0.5) = \text{false}$ .

*Example 2.* Given a pattern  $P = (?x, \text{sufferFrom}, ?y)\text{ AND } (?x, \text{Treatedby}, ?z)$  (i.e., we query those who have suffered from some illness and have been treated), we can compute that  $\llbracket P \rrbracket_{R_{\text{john}}} = \{\tau_1, \tau_2\}$ , where

- (i)  $\tau_1 = (\{?x \rightarrow \text{John}, ?y \rightarrow \text{Schizophrenia}, ?z \rightarrow \text{Psychiatrist}\}, 0.32)$ ;
- (ii)  $\tau_2 = (\{?x \rightarrow \text{John}, ?y \rightarrow \text{MentalDisorder}, ?z \rightarrow \text{Psychiatrist}\}, 0.84)$ .

TABLE 2

No	Triple: ( $t$ )	Probability: ( $\rho(t)$ )
$t_1$	(John, sufferedFrom, Schizophrenia)	0.32
$t_2$	(John, sufferedFrom, MentalDisorder)	0.84
$t_3$	(John, Treatedby, Psychiatrist)	0.95

*Example 3.* Given a pattern  $P = (?x, sufferFrom, Schizophrenia) \cup (?x, Treatedby, Psychiatrist)$  (i.e., we query those who have suffered from schizophrenia or those who are treated by psychiatrists); we can compute that  $\llbracket P \rrbracket_{R, \text{john}} = \{(?x \rightarrow \text{John}, 0.95)\}$ .

Note that  $?p$  is slightly different from variables where the value of  $?p$  is variable via probability computation, while the value of other variables is fixed. Moreover, we disallow the comparison of probability in filters.

#### 4. Well-Definedness, Uniqueness, and Equivalence of pSPARQL

In this section, we discuss some important properties of pSPARQL.

Firstly, we introduce a property called well-definedness, which can ensure that the semantics of pSPARQL are well defined.

**Proposition 4** (well-definedness). *For any pSPARQL pattern  $P$ , for any probabilistic RDF  $R$ , for any solution  $(\mu, p) \in \llbracket P \rrbracket_R$ , we have  $p \in (0, 1]$ .*

*Proof.* By induction on the structure of  $P$ ,

if  $P$  is a triple pattern  $(u, v, w)$ , then  $p = ((\mu(u), \mu(v), \mu(w))) \in (0, 1]$ ;

if  $P$  is of the form  $P_1 \cup P_2$ , then let us discuss the three cases:

(i) if  $(\mu, p) \in \llbracket P_1 \rrbracket_R$  but  $(\mu, p) \notin \llbracket P_2 \rrbracket_R$ , then  $p \in (0, 1]$ ;

(ii) if  $(\mu, p) \in \llbracket P_2 \rrbracket_R$  but  $(\mu, p) \notin \llbracket P_1 \rrbracket_R$ , then  $p \in (0, 1]$ ;

(iii) if  $(\mu, p) \in \llbracket P_1 \rrbracket_R$  and  $(\mu, p) \in \llbracket P_2 \rrbracket_R$ , then  $p = \max\{p_1, p_2\} \in (0, 1]$  by induction.

If  $P$  is of the form  $P_1 \text{ AND } P_2$ , then this claim holds by induction, since there exists some solution  $(\mu_1, p_1) \in \llbracket P_1 \rrbracket_R$  and some solution  $(\mu_2, p_2) \in \llbracket P_2 \rrbracket_R$  with  $\mu_1 \sim \mu_2$  such that

$$\mu = \mu_1 \cup \mu_2$$

$$\text{and } p \geq \min\{p_1, p_2\}$$

$$\text{and } p = \max_{\mu_1 \cup \mu_2 = \mu} \min\{p_i, p_j \mid (\mu_i, p_i) \in \llbracket P_1 \rrbracket_R \text{ and } (\mu_j, p_j) \in \llbracket P_2 \rrbracket_R \text{ and } \mu_i \sim \mu_j\} \in (0, 1]. \quad (10)$$

If  $P$  is of the form  $P_1 \text{ DIFF } P_2$  or  $P_1 \text{ FILTER } F$ , then this claim holds by induction, since

$$\llbracket P \rrbracket_R \subseteq \llbracket P_1 \rrbracket_R. \quad (11)$$

Finally, if  $P$  is of the form  $P_1 \text{ OPT } P_2$ , then this claim holds by the cases of  $P_1 \text{ AND } P_2$ ,  $P_1 \text{ UNION } P_2$ , and  $P_1 \text{ DIFF } P_2$  by induction.  $\square$

**Proposition 5** (uniqueness). *For any pSPARQL pattern  $P$ , for any probabilistic RDF  $R$ , for any two solutions  $(\mu, p), (\mu', p') \in \llbracket P \rrbracket_R$ , if  $\mu = \mu'$ , then  $p = p'$ .*

*Proof.* By induction on the structure of  $P$ , we have the following.

If  $P$  is a triple pattern  $(u, v, w)$ , then this claim directly holds by definition, since  $p_1 = p_2 = \rho(\mu(u), \mu(v), \mu(w))$ .

If  $P$  is of the form  $P_1 \cup P_2$ , then let us discuss the three cases:

(i) If  $(\mu, p) \in \llbracket P_1 \rrbracket_R$  but  $(\mu, p) \notin \llbracket P_2 \rrbracket_R$ , then this claim holds by induction.

(ii) If  $(\mu, p) \in \llbracket P_2 \rrbracket_R$  but  $(\mu, p) \notin \llbracket P_1 \rrbracket_R$ , then this claim holds by induction.

(iii) If there exist  $(\mu, p_1) \in \llbracket P_1 \rrbracket_R$  and  $(\mu, p_2) \in \llbracket P_2 \rrbracket_R$ , then this claim holds by induction, since  $p = \max\{p_1, p_2\}$ .

If  $P$  is of the form  $P_1 \text{ AND } P_2$ , then this claim holds by induction, since there exists some solution  $(\mu_1, p_1) \in \llbracket P_1 \rrbracket_R$  and some solution  $(\mu_2, p_2) \in \llbracket P_2 \rrbracket_R$  with  $\mu_1 \sim \mu_2$  such that  $\mu = \mu_1 \cup \mu_2$  and  $p = \min\{p_i, p_j \mid (\mu_i, p_i) \in \llbracket P_1 \rrbracket_R \text{ and } (\mu_j, p_j) \in \llbracket P_2 \rrbracket_R \text{ and } \mu_i \sim \mu_j\}$ . Therefore,  $p$  is unique.

If  $P$  is of the form  $P_1 \text{ DIFF } P_2$  or  $P_1 \text{ FILTER } F$ , then this claim holds by induction, since

$$\llbracket P \rrbracket_R \subseteq \llbracket P_1 \rrbracket_R. \quad (12)$$

Finally, we discuss the equivalence of patterns in pSPARQL. Let  $P$  and  $Q$  be two patterns in pSPARQL. We say that  $P$  is *equivalent* to  $Q$ , denoted by  $P \equiv_p Q$ , if for any probabilistic RDF  $R$ ,  $\llbracket P \rrbracket_R = \llbracket Q \rrbracket_R$ .  $\square$

Next, we show that pSPARQL satisfies the distributive law of equivalence, which is proven to be important in SPARQL.

**Proposition 6** (distributive law). *Let  $P_1, P_2$ , and  $P_3$  be three patterns in pSPARQL and let  $F$  be a filter. The following holds:*

(1)  $(P_1 \cup P_2) \text{ FILTER } F \equiv_p (P_1 \text{ FILTER } F) \cup (P_2 \text{ FILTER } F)$ ;

(2)  $P_1 \text{ AND } (P_2 \cup P_3) \equiv_p (P_1 \text{ AND } P_2) \cup (P_1 \text{ AND } P_3)$ ;

- (3)  $(P_1 \text{ UNION } P_2) \text{ AND } P_3 \equiv_p (P_1 \text{ AND } P_3) \text{ UNION } (P_2 \text{ AND } P_3)$ ;
- (4)  $(P_1 \text{ UNION } P_2) \text{ DIFF } P_3 \equiv_p (P_1 \text{ DIFF } P_3) \text{ UNION } (P_2 \text{ DIFF } P_3)$ ;
- (5)  $(P_1 \text{ UNION } P_2) \text{ OPT } P_3 \equiv_p (P_1 \text{ OPT } P_3) \text{ UNION } (P_2 \text{ OPT } P_3)$ .

*Proof (sketch).* The first claim directly holds by the definition.

Now, we show the second item. Let  $R$  be a probabilistic RDF of the form  $(G, \rho)$ . If  $(\mu, p) \in \llbracket P_1 \text{ AND } (P_2 \text{ UNION } P_3) \rrbracket_R$ , then there must exist some solution  $(\mu, p') \in \llbracket (P_1 \text{ AND } P_2) \text{ UNION } (P_1 \text{ AND } P_3) \rrbracket_R$ . By Proposition 5, we can conclude that  $p = p'$ . Then  $\llbracket P_1 \text{ AND } (P_2 \text{ UNION } P_3) \rrbracket_R \subseteq \llbracket (P_1 \text{ AND } P_2) \cup (P_1 \text{ AND } P_3) \rrbracket_R$ .

On the other hand,  $(\mu, p) \in \llbracket (P_1 \text{ AND } P_2) \cup (P_1 \text{ AND } P_3) \rrbracket_R$ ; then there must exist some solution  $(\mu, p') \in \llbracket P_1 \text{ AND } (P_2 \text{ UNION } P_3) \rrbracket_R$ . By Proposition 5, we can conclude that  $p = p'$ . Then  $\llbracket (P_1 \text{ AND } P_2) \cup (P_1 \text{ AND } P_3) \rrbracket_R \subseteq \llbracket P_1 \text{ AND } (P_2 \text{ UNION } P_3) \rrbracket_R$ .

Analogously, we can prove the third item and the fourth item.

Finally, we could prove the fifth item by using the third item and the fourth item.  $\square$

## 5. A Practical Example

In this section, we illustrate the application of pSPARQL in a real world via a practical example, where a probabilistic RDF is introduced in [11] shown in Figure 1.

Consider the following four queries ( $Q_1, Q_2, Q_3, Q_4$ ) in pSPARQL.

(1)  $Q_1$ : What causes fatigue associated with some illness over 0.65 probability?

$Q_1$  is formally expressed in pSPARQL as follows:

SELECT  $?x$  ((*Fatigue, CauseOf, ?x*) AND (( $?x, \text{AssociatedWith}, ?z$ ) FILTER  $?p > 0.65$ )).

The solution of  $Q_1$  is as follows:

$$\llbracket Q_1 \rrbracket_{R_{rsv}} = \{(\{?x \rightarrow \text{Flu}\}, 0.6)\}. \quad (13)$$

Note that  $\llbracket (?x, \text{AssociatedWith}, ?z) \rrbracket_{R_{rsv}} = \{(\{?x \rightarrow \text{Flu}, ?z \rightarrow \text{Cough}\}, 0.7), (\{?x \rightarrow \text{Pneumonia}, ?z \rightarrow \text{Bronchitis}\}, 0.6)\}$ . Thus  $\llbracket (?x, \text{AssociatedWith}, ?z) \text{ FILTER } ?p > 0.65 \rrbracket_{R_{rsv}} = \{(\{?x \rightarrow \text{Flu}, ?z \rightarrow \text{Cough}\}, 0.7)\}$ . Then  $\llbracket Q_1 \rrbracket_{R_{rsv}} = \{(\{?x \rightarrow \text{Flu}\}, 0.6)\}$ , since  $\llbracket (\text{Fatigue, CauseOf}, ?x) \rrbracket_{R_{rsv}} = \{(\{?x \rightarrow \text{Flu}\}, 0.6), (\{?x \rightarrow \text{Pneumonia}\}, 0.6)\}$ .

(2)  $Q_2$ : What are associated with cough over 0.7 probability?

$Q_2$  is formally expressed in pSPARQL as follows:

SELECT  $?x$  (( $?x, \text{AssociatedWith}, ?y$ ) FILTER  $?y = \text{'Cough'}$   $\wedge ?p > 0.7$ ).

The solution of  $Q_2$  is as follows:

$$\begin{aligned} \llbracket Q_2 \rrbracket_{R_{rsv}} \\ = \{(\{?x \rightarrow \text{Bronchitis}, ?y \rightarrow \text{Cough}\}, 0.8)\}. \end{aligned} \quad (14)$$

TABLE 3

$?x$	$?y$	$?p$
<i>Bronchitis</i>	<i>Cough</i>	0.8
<i>Bronchitis</i>	<i>RSV</i>	0.6
<i>RSV</i>	<i>Cough</i>	0.7
<i>Flu</i>	<i>Cough</i>	0.7
<i>Pneumonia</i>	<i>Bronchitis</i>	0.6

Note that  $\llbracket (?x, \text{AssociatedWith}, ?y) \rrbracket_{R_{rsv}}$  is shown in Table 3.

Thus  $\llbracket (?x, \text{AssociatedWith}, ?y) \text{ FILTER } ?y = \text{'Cough'} \rrbracket_{R_{rsv}} = \{(\{?x \rightarrow \text{Bronchitis}, ?y \rightarrow \text{Cough}\}, 0.8), (\{?x \rightarrow \text{RSV}, ?y \rightarrow \text{Cough}\}, 0.7), (\{?x \rightarrow \text{Flu}, ?y \rightarrow \text{Cough}\}, 0.7)\}$ .

Then  $\llbracket Q_2 \rrbracket_{R_{rsv}} = \{(\{?x \rightarrow \text{Bronchitis}, ?y \rightarrow \text{Cough}\}, 0.8)\}$ .

(3)  $Q_3$ : What is probability of bronchitis associated with cough directly or indirectly?

$Q_3$  is formally expressed as follows:

SELECT  $?x$  (( $?x, \text{AssociatedWith}, ?y$ ) UNION (( $?x, \text{AssociatedWith}, ?z$ ) AND ( $?z, \text{AssociatedWith}, ?y$ )) FILTER  $?x = \text{'Bronchitis'} \wedge ?y = \text{'Cough'}$ ).

The solution of  $Q_3$  is as follows:

$$\begin{aligned} \llbracket Q_3 \rrbracket_{R_{rsv}} \\ = \{(\{?x \rightarrow \text{Bronchitis}, ?y \rightarrow \text{Cough}\}, 0.8)\}. \end{aligned} \quad (15)$$

Note that  $\llbracket ((?x, \text{AssociatedWith}, ?y)) \rrbracket_{R_{rsv}}$  is shown in Table 4.

Note that  $\llbracket (?x, \text{AssociatedWith}, ?z) \text{ AND } (?z, \text{AssociatedWith}, ?y) \rrbracket_{R_{rsv}}$  is shown in Table 5.

Thus  $\llbracket (?x, \text{AssociatedWith}, ?y) \text{ UNION } ((?x, \text{AssociatedWith}, ?z) \text{ AND } (?z, \text{AssociatedWith}, ?y)) \rrbracket_{R_{rsv}}$  is shown in Table 6.

Then

$$\begin{aligned} \llbracket Q_3 \rrbracket_{R_{rsv}} \\ = \{(\{?x \rightarrow \text{Bronchitis}, ?y \rightarrow \text{Cough}\}, 0.8)\}. \end{aligned} \quad (16)$$

(4)  $Q_4$ : What are associated with cough excluding bronchitis?

$Q_4$  is expressed in pSPARQL as follows:

SELECT  $?x$  ((( $?x, \text{AssociatedWith}, ?y$ ) FILTER  $?y = \text{'Cough'}$ ) DIFF (( $?x, \text{AssociatedWith}, ?y$ ) FILTER  $?x = \text{'Bronchitis'} \wedge ?y = \text{'Cough'}$ )).

The solution of  $Q_4$  is as follows:

$$\llbracket Q_4 \rrbracket_{R_{rsv}} = \{(\{?x \rightarrow \text{Flu}, ?y \rightarrow \text{Cough}\}, 0.7)\}. \quad (17)$$

Note that  $\llbracket (?x, \text{AssociatedWith}, ?y) \text{ FILTER } ?y = \text{'Cough'} \rrbracket_{R_{rsv}}$  is shown in Table 7.

Note that  $\llbracket (?x, \text{AssociatedWith}, ?y) \text{ FILTER } ?x = \text{'Bronchitis'} \wedge ?y = \text{'Cough'} \rrbracket_{R_{rsv}}$  is shown in Table 8.

Then  $\llbracket Q_4 \rrbracket_{R_{rsv}} = \{(\{?x \rightarrow \text{Flu}, ?y \rightarrow \text{Cough}\}, 0.7)\}$ .

In short, we could express many interesting queries with respect to probabilistic RDF via pSPARQL, which are useful

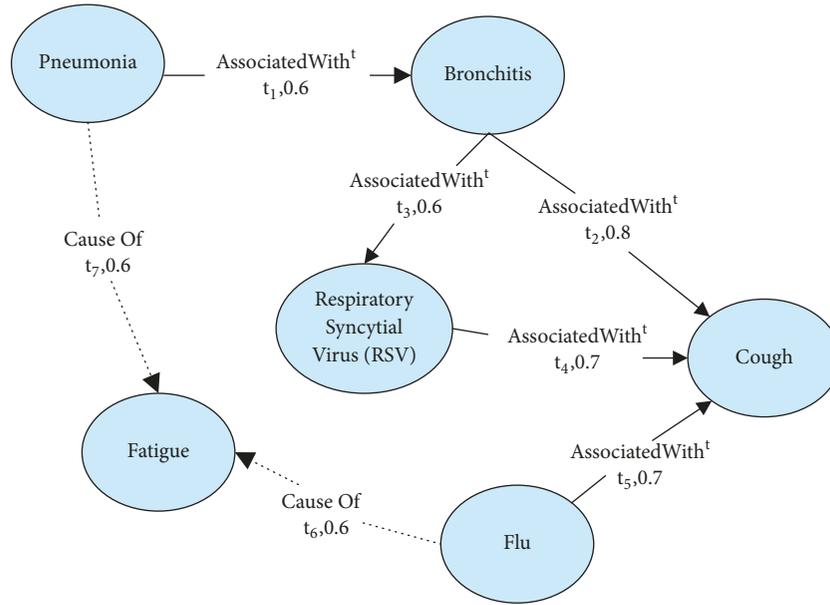
FIGURE 1: A virus RDF graph  $R_{rsv}$  [11].

TABLE 4

$?x$	$?y$	$?p$
<i>Bronchitis</i>	<i>Cough</i>	0.8
<i>Bronchitis</i>	<i>RSV</i>	0.6
<i>RSV</i>	<i>Cough</i>	0.7
<i>Flu</i>	<i>Cough</i>	0.7
<i>Pneumonia</i>	<i>Bronchitis</i>	0.6

TABLE 5

$?x$	$?y$	$?z$	$?p$
<i>Bronchitis</i>	<i>Cough</i>	<i>RSV</i>	0.6
<i>Pneumonia</i>	<i>Cough</i>	<i>Bronchitis</i>	0.6
<i>Pneumonia</i>	<i>RSV</i>	<i>Bronchitis</i>	0.6

TABLE 6

$?x$	$?y$	$?z$	$?p$
<i>Bronchitis</i>	<i>Cough</i>		0.8
<i>Bronchitis</i>	<i>RSV</i>		0.6
<i>RSV</i>	<i>Cough</i>		0.7
<i>Flu</i>	<i>Cough</i>		0.7
<i>Pneumonia</i>	<i>Bronchitis</i>		0.6
<i>Bronchitis</i>	<i>Cough</i>	<i>RSV</i>	0.6
<i>Pneumonia</i>	<i>Cough</i>	<i>Bronchitis</i>	0.6
<i>Pneumonia</i>	<i>RSV</i>	<i>Bronchitis</i>	0.6

in a practical world. Compared with SPARQL, where we obtain only connection via SPARQL querying, we could quantize the connection via pSPARQL, so that we could obtain more specific solutions.

## 6. Conclusions

In this paper, we extended SPARQL to support querying over probabilistic RDF. In the future, we will discuss some

TABLE 7

$?x$	$?y$	$?p$
<i>Bronchitis</i>	<i>Cough</i>	0.8
<i>RSV</i>	<i>Cough</i>	0.7
<i>Flu</i>	<i>Cough</i>	0.7

TABLE 8

$?x$	$?y$	$?p$
<i>Bronchitis</i>	<i>Cough</i>	0.8

foundational properties of pSPARQL and implement it in a prototype to provide the full SPARQL query answering services for probabilistic RDF. As a future work, we are interested in presenting probabilistic semantics of RDF graphs in a unified framework, where many applications could be supported.

## Data Availability

No data were used to support this study.

## Disclosure

An earlier version of this work was presented at “International Conference on Big Scientific Data Management 2018.”

## Conflicts of Interest

The author declares that there are no conflicts of interest regarding the publication of this paper.

## Acknowledgments

This work is supported by the program of the key discipline “Applied Mathematics” of Shanghai Polytechnic University (XXKPY1604).

## References

- [1] “RDF primer, W3C Recommendation, February 2004”.
- [2] F. M. Suchanek, G. Kasneci, and G. Weikum, “Yago: a core of semantic knowledge,” in *Proceedings of the 16th International World Wide Web Conference (WWW '07)*, pp. 697–706, Alberta, Canada, May 2007.
- [3] U. Fayyad, G. Piatetsky-Shapiro, and P. Smyth, “From data mining to knowledge discovery in databases,” *AI Magazine*, vol. 17, no. 3, pp. 37–53, 1996.
- [4] N. Dalvi and D. Suciu, “Efficient query evaluation on probabilistic databases,” in *Proceedings of the VLDB'04*, pp. 864–875, 2004.
- [5] Y. Fukushige, “Representing probabilistic relations in RDF in,” in *Proceedings of the ISWC-URSW'05*, pp. 106–107, 2005.
- [6] D. Suciu, *Probabilistic Databases, Encyclopedia of Database Systems*, Springer, 2009.
- [7] O. Udrea, V. Subrahmanian, and Z. Majkic, “Probabilistic RDF,” in *Proceedings of the 2006 IEEE International Conference on Information Reuse & Integration*, pp. 172–177, Waikoloa Village, HI, USA, September 2006.
- [8] “SPARQL query language for RDF, W3C Recommendation, January 2008”.
- [9] P. T. Wood, “Query languages for graph databases,” *SIGMOD Record*, vol. 41, no. 1, pp. 50–60, 2012.
- [10] A. Khan and L. Chen, “On uncertain graphs modeling and queries,” in *Proceedings of the PVLDB Endowment*, vol. 8, pp. 2042–2043, 2015.
- [11] H. Huang and C. Liu, “Query evaluation on probabilistic RDF databases,” in *Proceedings of the WISE'09*, pp. 307–320, 2009.
- [12] C. Szeto, E. Hung, and Y. Deng, “SPARQL query answering with RDFS reasoning on correlated probabilistic data,” in *Proceedings of the WAIM'11*, pp. 56–67, 2011.
- [13] X. Lian and L. Chen, “Efficient query answering in probabilistic RDF graphs,” in *Proceedings of the the 2011 international conference*, p. 157, Athens, Greece, June 2011.
- [14] D. Krompaß, M. Nickel, and V. Tresp, “Querying factorized probabilistic triple databases,” in *Proceedings of the ISWC'14*, pp. 114–129, 2014.
- [15] J. Schoenfish, “Querying probabilistic ontologies with SPARQL,” in *Proceedings of the KI'14*, pp. 2245–2256, 2014.
- [16] X. Zhou, Y. Chen, and D. Z. Wang, “ArchimedesOne: Query processing over probabilistic knowledge bases,” *Proceedings of the VLDB Endowment*, vol. 9, no. 13, pp. 1461–1464, 2016.
- [17] T. Andronikos, A. Singh, K. Giannakis, and S. Sioutas, “Computing probabilistic queries in the presence of uncertainty via probabilistic automata,” in *Proceedings of the ALGO CLOUD'17*, pp. 106–120, 2017.
- [18] X. Zhang and J. Van den Bussche, “On the primitivity of operators in SPARQL,” *Information Processing Letters*, vol. 114, no. 9, pp. 480–485, 2014.
- [19] “SPARQL 1.1 query language, W3C Recommendation, March 2013”.
- [20] X. Zhang, J. Van den Bussche, K. Wang, and Z. Wang, “On the satisfiability problem of patterns in SPARQL 1.1,” in *Proceedings of the AAAI'18*, pp. 2054–2061, 2018.
- [21] H. Fang and X. Zhang, “pSPARQL: a querying language for probabilistic RDF (extended abstract),” in *Proceedings of the ISWC'16, Posters*, 2016.
- [22] J. Pérez, M. Arenas, and C. Gutierrez, “Semantics and complexity of SPARQL,” *ACM Transactions on Database Systems (TODS)*, vol. 34, no. 3, pp. 1–45, 2009.

