



Learning Commonalities in SPARQL

Sara El Hassad François Goasdoué Hélène Jaudoin

IRISA, Univ. Rennes 1, Lannion, France

ISWC 2017 - 21 - 26 October 2017

Introduction

Least general generalization (lgg)

- ▶ Machine Learning in the early 70's by Gordon Plotkin
- ▶ Knowledge representation domain in the early 90's
- ▶ Recently in semantic web

Introduction

Least general generalization (lgg)

- ▶ Machine Learning in the early 70's by Gordon Plotkin
- ▶ Knowledge representation domain in the early 90's
- ▶ Recently in semantic web

Applications of lgg

- ▶ Query optimization: identify candidate views, or potential query sharing
- ▶ Query approximation: a set of queries by a single query
- ▶ Social context: recommending users asking for enough relates things

Introduction

Least general generalization (lgg)

- ▶ Machine Learning in the early 70's by Gordon Plotkin
- ▶ Knowledge representation domain in the early 90's
- ▶ Recently in semantic web

Applications of lgg

- ▶ Query optimization: identify candidate views, or potential query sharing
- ▶ Query approximation: a set of queries by a single query
- ▶ Social context: recommending users asking for enough relates things

Goal

To study the problem in the *entire* conjunctive fragment of SPARQL setting.

Outline

Introduction

Preliminaries

Finding commonalities between SPARQL conjunctive queries

Experiments

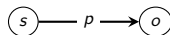
Related work

Conclusion

RDF graphs

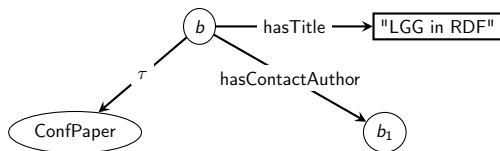
- ▶ Specification of RDF graphs with triples:

$$(s, p, o) \in (\mathcal{U} \cup \mathcal{B}) \times \mathcal{U} \times (\mathcal{U} \cup \mathcal{L} \cup \mathcal{B})$$



- ▶ Built-in property URIs to state RDF statements

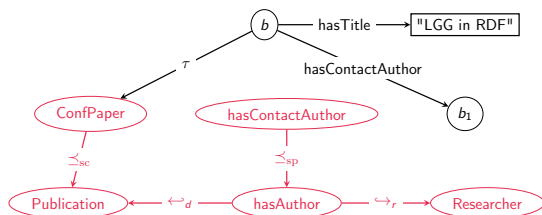
RDF statement	Triple
Class assertion	$(s, \text{rdf:type}, o)$
Property assertion	(s, p, o) with $p \neq \text{rdf:type}$



Adding ontological knowledge to RDF graphs

- ▶ Built-in property URIs to state RDF Schema statements, i.e., ontological constraints.

RDFS statement	Triple
Subclass	(s, \preceq_{sc}, o)
Subproperty	(s, \preceq_{sp}, o)
Domain typing	$(s, \leftrightarrow_d, o)$
Range typing	$(s, \hookrightarrow_r, o)$



Deriving the implicit triples

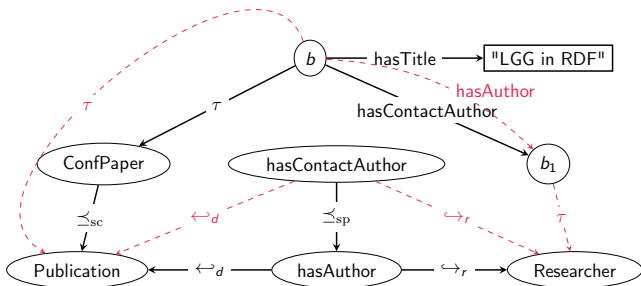


Figure: RDF graph \mathcal{G}

How to derive implicit triples of an RDF graph ?

Sample set of entailment rules

Rule [W3C-RDFS, 2014]	Entailment rule
rdfs2	$(p, \xleftrightarrow{d}, o), (s_1, p, o_1) \rightarrow (s_1, \tau, o)$
rdfs3	$(p, \xleftrightarrow{r}, o), (s_1, p, o_1) \rightarrow (o_1, \tau, o)$
rdfs5	$(p_1, \preceq_{sp}, p_2), (p_2, \preceq_{sp}, p_3) \rightarrow (p_1, \preceq_{sp}, p_3)$
rdfs7	$(p_1, \preceq_{sp}, p_2), (s, p_1, o) \rightarrow (s, p_2, o)$
rdfs9	$(s, \preceq_{sc}, o), (s_1, \tau, s) \rightarrow (s_1, \tau, o)$
rdfs11	$(s, \preceq_{sc}, o), (o, \preceq_{sc}, o_1) \rightarrow (s, \preceq_{sc}, o_1)$
ext1	$(p, \xleftrightarrow{d}, o), (o, \preceq_{sc}, o_1) \rightarrow (p, \xleftrightarrow{d}, o_1)$
ext2	$(p, \xleftrightarrow{r}, o), (o, \preceq_{sc}, o_1) \rightarrow (p, \xleftrightarrow{r}, o_1)$
ext3	$(p, \preceq_{sp}, p_1), (p_1, \xleftrightarrow{d}, o) \rightarrow (p, \xleftrightarrow{d}, o)$
ext4	$(p, \preceq_{sp}, p_1), (p_1, \xleftrightarrow{r}, o) \rightarrow (p, \xleftrightarrow{r}, o)$

Table: Sample RDF entailment rules \mathcal{R}

Semantics of RDF graphs

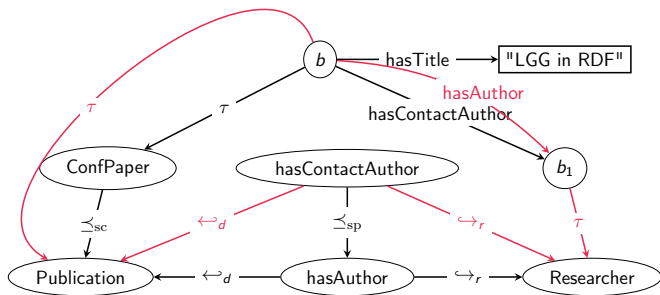


Figure: Saturated RDF graph \mathcal{G}^∞

Basic graph pattern queries (BGPQ)

- ▶ BGPQ : conjunctive fragment of SPARQL queries, is the counterpart of the select-project-join queries for databases
- ▶ $(s, p, o) \in (\mathcal{V} \cup \mathcal{U}) \times (\mathcal{V} \cup \mathcal{U}) \times (\mathcal{V} \cup \mathcal{U} \cup \mathcal{L})$

Basic graph pattern queries (BGPQ)

- ▶ BGPQ : conjunctive fragment of SPARQL queries, is the counterpart of the select-project-join queries for databases
- ▶ $(s, p, o) \in (\mathcal{V} \cup \mathcal{U}) \times (\mathcal{V} \cup \mathcal{U}) \times (\mathcal{V} \cup \mathcal{U} \cup \mathcal{L})$

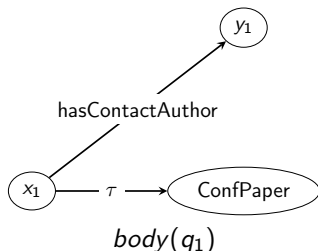
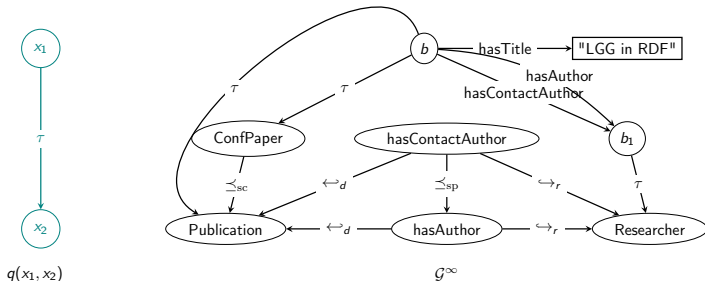


Figure: Sample BGPQ $q_1(x_1)$

Entailing and answering queries

Query entailment

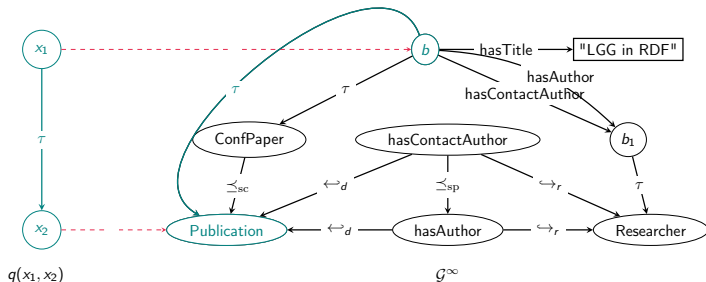
$$\mathcal{G} \models_{\mathcal{R}} q \iff \mathcal{G}^{\infty} \models_{\mathcal{R}} q$$



Entailing and answering queries

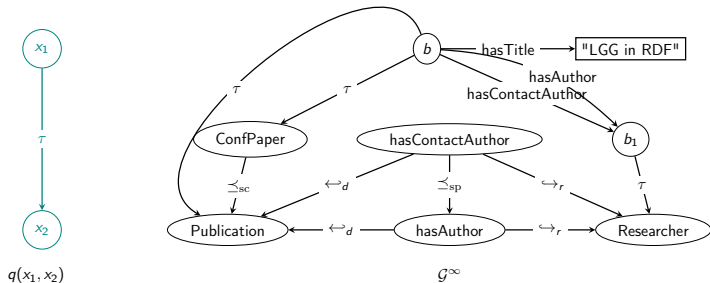
Query entailment

$$\mathcal{G} \models_{\mathcal{R}} q \iff \mathcal{G}^{\infty} \models_{\mathcal{R}} q$$



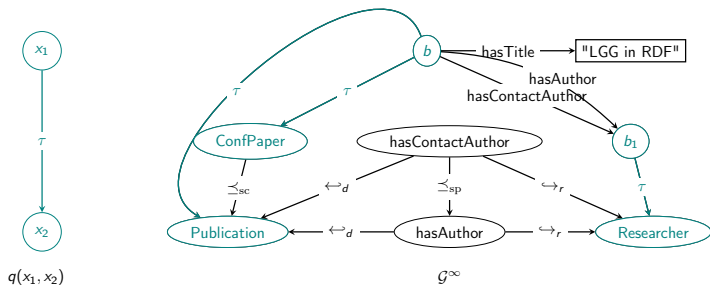
Entailing and answering queries

Query answering



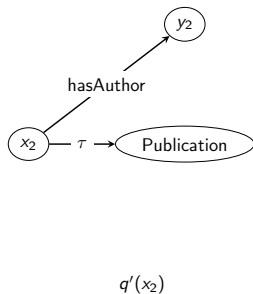
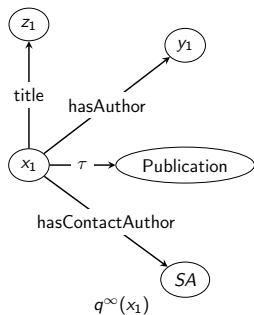
Entailing and answering queries

Query answering



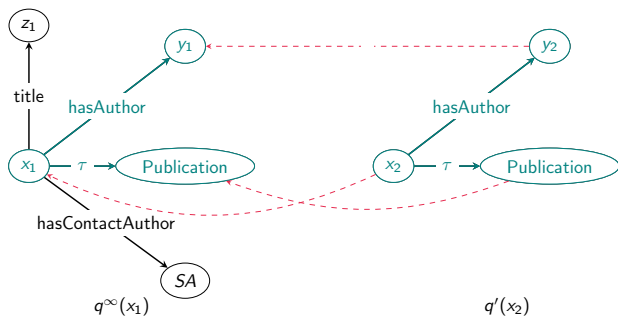
Entailing between BGPQs

$$q \models_{\mathcal{R}} q' \iff q^{\infty} \models q'$$



Entailing between BGPQs

$$q \models_{\mathcal{R}} q' \iff q^{\infty} \models q'$$



Outline

Introduction

Preliminaries

Finding commonalities between SPARQL conjunctive queries

Experiments

Related work

Conclusion

Towards defining lgg in SPARQL conjunctive fragment

A *least general generalization* (lgg) of n descriptions d_1, \dots, d_n is a most specific description d generalizing every $d_{1 \leq i \leq n}$ for some generalization/specialization relation between descriptions (G.Plotkin).

lgg in our SPARQL setting

- ▶ descriptions are BGP Queries
- ▶ relation generalization/specialization is entailment between queries

Defining the lgg of queries

lgg of BGPQs

Let q_1, \dots, q_n be BGPQs with the same arity and \mathcal{R} a set of RDF entailment rules.

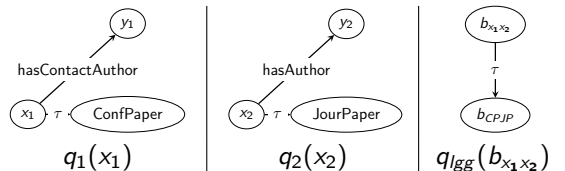
- ▶ A *generalization* of q_1, \dots, q_n is a BGPQ q_g such that $q_i \models_{\mathcal{R}} q_g$ for $1 \leq i \leq n$.
- ▶ A *least general generalization* of q_1, \dots, q_n is a generalization q_{lgg} of q_1, \dots, q_n such that for any other generalization q_g of q_1, \dots, q_n : $q_{\text{lgg}} \models_{\mathcal{R}} q_g$.

Defining the lgg of queries

lgg of BGPQs

Let q_1, \dots, q_n be BGPQs with the same arity and \mathcal{R} a set of RDF entailment rules.

- ▶ A *generalization* of q_1, \dots, q_n is a BGPQ q_g such that $q_i \models_{\mathcal{R}} q_g$ for $1 \leq i \leq n$.
- ▶ A *least general generalization* of q_1, \dots, q_n is a generalization q_{lgg} of q_1, \dots, q_n such that for any other generalization q_g of q_1, \dots, q_n : $q_{\text{lgg}} \models_{\mathcal{R}} q_g$.

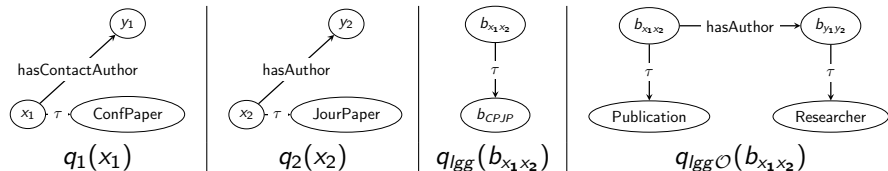


Defining the lgg of queries

lgg of BGPQs

Let q_1, \dots, q_n be BGPQs with the same arity and \mathcal{R} a set of RDF entailment rules.

- ▶ A *generalization* of q_1, \dots, q_n is a BGPQ q_g such that $q_i \models_{\mathcal{R}} q_g$ for $1 \leq i \leq n$.
- ▶ A *least general generalization* of q_1, \dots, q_n is a generalization q_{lgg} of q_1, \dots, q_n such that for any other generalization q_g of q_1, \dots, q_n : $q_{\text{lgg}} \models_{\mathcal{R}} q_g$.



Entailment relation between BGPQs w.r.t. background knowledge

Entailment between BGPQs w.r.t. \mathcal{R}, \mathcal{O}

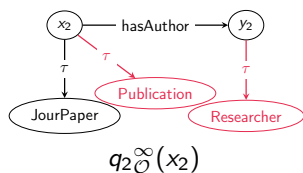
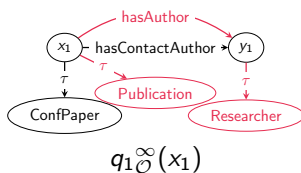
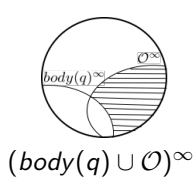
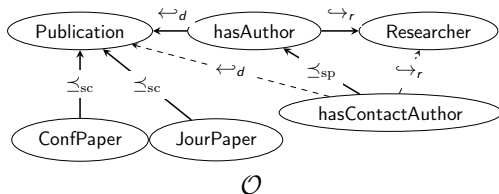
Given a set \mathcal{R} of RDF entailment rules, a set \mathcal{O} of RDFS statements, and two BGPQs q_1 and q_2 with the same arity, q_1 *entails* q_2 w.r.t. \mathcal{O} , denoted $q_1 \models_{\mathcal{R}, \mathcal{O}} q_2$, iff $q_1^{\infty}_{\mathcal{O}} \models q_2$ holds.

Well-founded relation : $q_1 \models_{\mathcal{R}, \mathcal{O}} q_2$

- ▶ **Query entailment:** if $\mathcal{G} \models_{\mathcal{R}} q_1$ holds then $\mathcal{G} \models_{\mathcal{R}} q_2$ holds,
- ▶ **Query answering:** $q_1(\mathcal{G}) \subseteq q_2(\mathcal{G})$ holds.

Saturation of queries

BGPQ saturation w.r.t. RDFS constraints



Defining the 1gg of queries w.r.t. background knowledge

Definition (1gg of BGPQs w.r.t. RDFS constraints)

Let \mathcal{R} be a set of RDF entailment rules, \mathcal{O} a set of RDFS statements, and q_1, \dots, q_n n BGPQs with the same arity.

- ▶ A *generalization* of q_1, \dots, q_n w.r.t. \mathcal{O} is a BGPQ q_g such that $q_i \models_{\mathcal{R}, \mathcal{O}} q_g$ for $1 \leq i \leq n$.
- ▶ A *least general generalization* of q_1, \dots, q_n w.r.t. \mathcal{O} is a generalization q_{1gg} of q_1, \dots, q_n w.r.t. \mathcal{O} such that for any other generalization q_g of q_1, \dots, q_n w.r.t. \mathcal{O} : $q_{1gg} \models_{\mathcal{R}, \mathcal{O}} q_g$.

Theorem

An 1gg of BGPQs w.r.t. RDFS statements may not exist for some set of RDF entailment rules; when it exists, it is unique up to entailment ($\models_{\mathcal{R}, \mathcal{O}}$).

Defining the lgg of queries w.r.t. background knowledge

Definition (lgg of BGPQs w.r.t. RDFS constraints)

Let \mathcal{R} be a set of RDF entailment rules, \mathcal{O} a set of RDFS statements, and q_1, \dots, q_n n BGPQs with the same arity.

- ▶ A *generalization* of q_1, \dots, q_n w.r.t. \mathcal{O} is a BGPQ q_g such that $q_i \models_{\mathcal{R}, \mathcal{O}} q_g$ for $1 \leq i \leq n$.
- ▶ A *least general generalization* of q_1, \dots, q_n w.r.t. \mathcal{O} is a generalization q_{lgg} of q_1, \dots, q_n w.r.t. \mathcal{O} such that for any other generalization q_g of q_1, \dots, q_n w.r.t. \mathcal{O} : $q_{\text{lgg}} \models_{\mathcal{R}, \mathcal{O}} q_g$.

Result : lgg of n BGPQ queries vs lgg of two BGPQ queries

$$\ell_3(q_1, q_2, q_3) \equiv_{\mathcal{R}, \mathcal{O}} \ell_2(\ell_2(q_1, q_2), q_3)$$

...

$$\begin{aligned} \ell_n(q_1, \dots, q_n) &\equiv_{\mathcal{R}, \mathcal{O}} \ell_2(\ell_{n-1}(q_1, \dots, q_{n-1}), q_n) \\ &\equiv_{\mathcal{R}, \mathcal{O}} \ell_2(\ell_2(\dots \ell_2(\ell_2(q_1, q_2), q_3) \dots), q_{n-1}), q_n) \end{aligned}$$

Defining the lgg of queries w.r.t. background knowledge

Definition (lgg of BGPQs w.r.t. RDFS constraints)

Let \mathcal{R} be a set of RDF entailment rules, \mathcal{O} a set of RDFS statements, and q_1, \dots, q_n n BGPQs with the same arity.

- ▶ A *generalization* of q_1, \dots, q_n w.r.t. \mathcal{O} is a BGPQ q_g such that $q_i \models_{\mathcal{R}, \mathcal{O}} q_g$ for $1 \leq i \leq n$.
- ▶ A *least general generalization* of q_1, \dots, q_n w.r.t. \mathcal{O} is a generalization q_{lgg} of q_1, \dots, q_n w.r.t. \mathcal{O} such that for any other generalization q_g of q_1, \dots, q_n w.r.t. \mathcal{O} : $q_{\text{lgg}} \models_{\mathcal{R}, \mathcal{O}} q_g$.

Result : lgg of n BGPQ queries vs lgg of two BGPQ queries

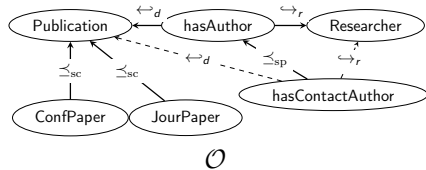
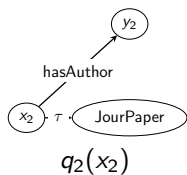
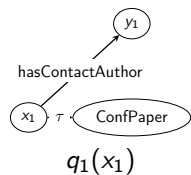
$$\ell_3(q_1, q_2, q_3) \equiv_{\mathcal{R}, \mathcal{O}} \ell_2(\ell_2(q_1, q_2), q_3)$$

...

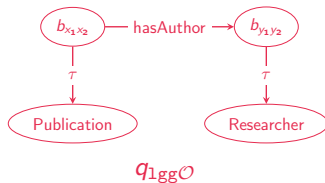
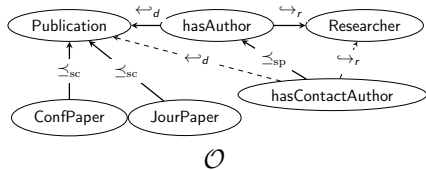
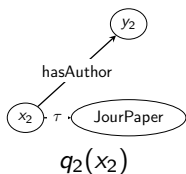
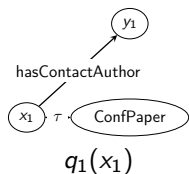
$$\begin{aligned} \ell_n(q_1, \dots, q_n) &\equiv_{\mathcal{R}, \mathcal{O}} \ell_2(\ell_{n-1}(q_1, \dots, q_{n-1}), q_n) \\ &\equiv_{\mathcal{R}, \mathcal{O}} \ell_2(\ell_2(\dots \ell_2(\ell_2(q_1, q_2), q_3) \dots), q_{n-1}), q_n) \end{aligned}$$

We focus on computing lgg of two BGPQ queries

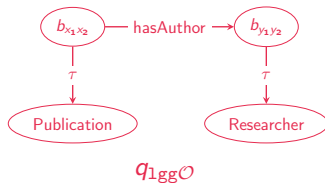
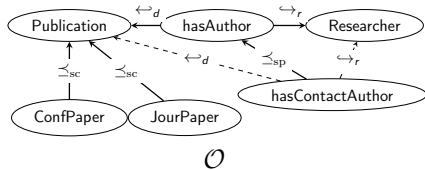
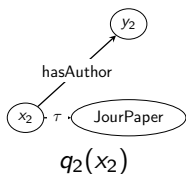
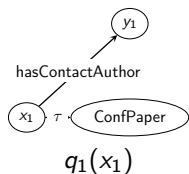
Defining the lgg of queries



Defining the lgg of queries



Defining the lgg of queries



How to compute this query ?

The cover of SPARQL queries

Definition (Cover query)

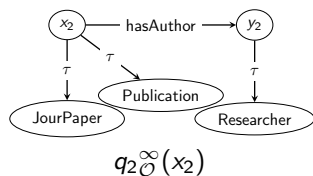
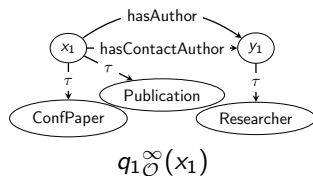
Let q_1, q_2 be two BGPQs with the same arity n .

If there exists the BGPQ q such that

- ▶ $head(q_1) = q(x_1^1, \dots, x_1^n)$ and $head(q_2) = q(x_2^1, \dots, x_2^n)$ iff
 $head(q) = q(v_{x_1^1 x_2^1}, \dots, v_{x_1^n x_2^n})$
- ▶ $(t_1, t_2, t_3) \in body(q_1)$ and $(t_4, t_5, t_6) \in body(q_2)$ iff
 $(t_7, t_8, t_9) \in body(q)$ with, for $1 \leq i \leq 3$, $t_{i+6} = t_i$ if $t_i = t_{i+3}$ and
 $t_i \in \mathcal{U} \cup \mathcal{L}$, otherwise t_{i+6} is the variable $v_{t_i t_{i+3}}$

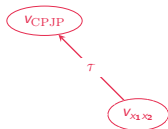
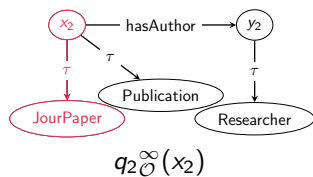
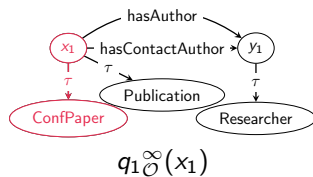
then q is the *cover query* of q_1, q_2 .

The cover of SPARQL queries



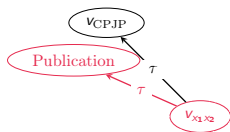
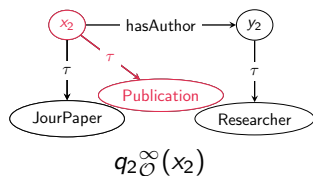
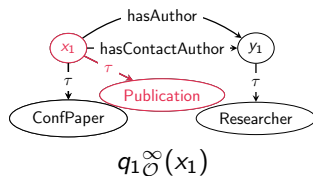
$q(V_{x_1x_2})$

The cover of SPARQL queries



$q(V_{x_1x_2})$

The cover of SPARQL queries



$q(V_{x_1x_2})$

Cover graph vs lgg

Theorem

Given a set \mathcal{R} of RDF entailment rules, a set \mathcal{O} of RDFS statements and two BGPQs q_1, q_2 with the same arity,

1. the cover query q of $q_1^{\infty}_{\mathcal{O}}, q_2^{\infty}_{\mathcal{O}}$ exists iff an lgg of q_1, q_2 w.r.t. \mathcal{O} exists;
2. the cover query q of $q_1^{\infty}_{\mathcal{O}}, q_2^{\infty}_{\mathcal{O}}$ is an lgg of q_1, q_2 w.r.t. \mathcal{O} .

Corollary

A cover query-based lgg of two BGPQs q_1 and q_2 is computed in $O(|body(q_1^{\infty}_{\mathcal{O}})| \times |body(q_2^{\infty}_{\mathcal{O}})|)$ and its size is $|body(q_1^{\infty}_{\mathcal{O}})| \times |body(q_2^{\infty}_{\mathcal{O}})|$.

Outline

Introduction

Preliminaries

Finding commonalities between SPARQL conjunctive queries

Experiments

Related work

Conclusion

lgg of DBpedia queries

$$q_{1\text{gg}}^{\mathcal{O}_{\text{DBpedia}}} \models q_{1\text{gg}}$$

lgg of:	Q_1Q_2	Q_1Q_3	Q_1Q_4	Q_2Q_3	Q_4Q_5	Q_5Q_6	Q_5Q_7	Q_7Q_8
Time to compute $q_{1\text{gg}}$	3	3	5	4	4	5	6	5
$ q_{1\text{gg}}(\mathcal{G}_{\text{DBpedia}}) $	477,455	34,747,102	34,901,117	34,747,102	1,977	1,221	35	70
Time to compute $q_{1\text{gg}}^{\mathcal{O}_{\text{DBpedia}}}$	13	14	14	15	15	14	17	18
$ q_{1\text{gg}}^{\mathcal{O}_{\text{DBpedia}}}(\mathcal{G}_{\text{DBpedia}}) $	10,637	7,874,768	456,690	4,537,824	1,701	780	34	36
Gain in precision	97.77	77.33	98.69	86.94	13.96	36.11	2.85	48.57

Table: Characteristics of cover query-based lgg of test queries, w/ or w/o using the DBpedia RDFS constraints.

lgg3 of :	$Q_1Q_2Q_3$	$Q_1Q_2Q_4$	$Q_1Q_3Q_4$	$Q_2Q_3Q_4$	$Q_4Q_7Q_8$	$Q_5Q_7Q_8$	$Q_6Q_7Q_8$
Time to compute $q_{1\text{gg}}$	5	4	5	6	10	11	12
$ q_{1\text{gg}}(\mathcal{G}_{\text{DBpedia}}) $	34,747,102	34,901,117	34,901,117	34,901,117	70	1,977	4,969
Time to compute $q_{1\text{gg}}^{\mathcal{O}_{\text{DBpedia}}}$	19	20	20	24	27	27	33
$ q_{1\text{gg}}^{\mathcal{O}_{\text{DBpedia}}}(\mathcal{G}_{\text{DBpedia}}) $	7,874,768	615,339	7,874,779	4,537,824	36	1,701	335
Gain in precision	77.33	98.23	77.43	86.99	48.57	13.96	93.25

Table: Characteristics of cover query-based lgg of 3 test queries, w/ or w/o using the DBpedia RDFS constraints; times are in ms.

Outline

Introduction

Preliminaries

Finding commonalities between SPARQL conjunctive queries

Experiments

Related work

Conclusion

Structural approaches

- ▶ RDF
 - ▶ Rooted graphs, ignore RDF entailment :
 - [Colucci et al., 2016].
- ▶ SPARQL : tree queries
 - [Lehmann and Bühmann, 2011].
- ▶ Description Logics
 - [Zarriß and Turhan, 2013].
 - [Baader et al., 1999].

Approaches independent of the structure

- ▶ RDF
 - [Hassad et al., 2017].
 - [Petrova et al., 2017].
- ▶ Conceptual Graphs
 - [Chein and Mugnier, 2009].
- ▶ First Order Clauses
 - [Nienhuys-Cheng and de Wolf, 1996].
 - [Plotkin, 1970].

Conclusion

- ▶ We revisited the problem of computing a least general generalization of general BGPQs w.r.t. background knowledge.
- ▶ We defined **new** entailment relationship between BGPQs w.r.t. background knowledge.
- ▶ We studied the added-value of considering background knowledge when learning lggs.

Perspective:

- ▶ Heuristics in order to compute 1gg without redundant triples.

Thank you !



Questions?

References I

- [Baader et al., 1999] Baader, F., Küsters, R., and Molitor, R. (1999).
Computing least common subsumers in description logics with existential restrictions.
In *IJCAI*.
- [Chein and Mugnier, 2009] Chein, M. and Mugnier, M. (2009).
Graph-based Knowledge Representation - Computational Foundations of Conceptual Graphs.
Springer.
- [Colucci et al., 2016] Colucci, S., Donini, F., Giannini, S., and Sciascio, E. D. (2016).
Defining and computing least common subsumers in RDF.
J. Web Semantics, 39(0).
- [Hassad et al., 2017] Hassad, S. E., Goasdoué, F., and Jaudoin, H. (2017).
Learning commonalities in RDF.
In *The 14th Extended Semantic Web Conference, ESWC 2017, Portorož, Slovenia, May 28 - June 1, 2017, Proceedings, Part I*, pages 502–517.
- [Lehmann and Bühmann, 2011] Lehmann, J. and Bühmann, L. (2011).
Autosparql: Let users query your knowledge base.
In *ESWC*.
- [Nienhuys-Cheng and de Wolf, 1996] Nienhuys-Cheng, S. and de Wolf, R. (1996).
Least generalizations and greatest specializations of sets of clauses.
J. Artif. Intell. Res.
- [Petrova et al., 2017] Petrova, A., Sherkhonov, E., Grau, B. C., and Horrocks, I. (2017).
Entity comparison in RDF graphs.
In *International Semantic Web Conference (ISWC)*. Springer.
- [Plotkin, 1970] Plotkin, G. D. (1970).
A note on inductive generalization.
Machine Intelligence, 5.
- [W3C-RDFS, 2014] W3C-RDFS (2014).
RDF 1.1 semantics.
<https://www.w3.org/TR/rdf11-mt/>.

References II

[Zarriß and Turhan, 2013] Zarriß, B. and Turhan, A. (2013).
Most specific generalizations w.r.t. general EL-TBoxes.
In *IJCAI*.