

Optimizing SPARQL Queries using Shape Statistics

Kashif Rabbani

Aalborg University, Denmark
kashifrabbani@cs.aau.dk

Matteo Lissandrini

Aalborg University, Denmark
matteo@cs.aau.dk

Katja Hose

Aalborg University, Denmark
khose@cs.aau.dk

ABSTRACT

With the growing popularity of storing data in native RDF, we witness more and more diverse use cases with complex SPARQL queries. As a consequence, query optimization – and in particular cardinality estimation and join ordering – becomes even more crucial. Classical methods exploit global statistics covering the entire RDF graph as a whole, which naturally fails to correctly capture correlations that are very common in RDF datasets, which then leads to erroneous cardinality estimations and suboptimal query execution plans. The alternative of trying to capture correlations in a fine-granular manner, on the other hand, results in very costly preprocessing steps to create these statistics. Hence, in this paper we propose *shapes statistics*, which extend the recent SHACL standard with statistic information to capture the correlation between classes and properties. Our extensive experiments on synthetic and real data show that shapes statistics can be generated and managed with only little overhead without disadvantages in query runtime while leading to noticeable improvements in cardinality estimation.

1 INTRODUCTION

Driven by diverse movements, such as Linked Open Government Data, Open Street Map, DBpedia [3], and YAGO [21], more and more data is being published in RDF [7] capturing a multitude of diverse information. Along with the growing popularity, increasingly complex queries formulated in SPARQL [6] are being executed over such data to answer business and research questions. Query logs of the public DBpedia SPARQL endpoint, for instance, contain SPARQL queries with up to 10 joins [4] and analytic queries in the biomedical field can involve more than 50 joins per query [9]. Therefore, the need for high-performance SPARQL query processing is now more pressing than ever.

Existing approaches for query optimization in RDF stores often adapt techniques from relational databases modeling an RDF dataset as a single large table with three columns [5, 16] (one column for each of the components of an RDF triple: subject, predicate, and object). Nevertheless, accurate cardinality estimation is at the heart of any query optimizer that does not rely on heuristics but instead uses a cost model to find the best query execution plan for a given query. Cardinality estimation then relies on the availability of statistics describing the characteristics of the data to estimate the sizes of intermediate results produced while query execution. However, general statistics typically result in highly imprecise estimations since they are mostly gathered on the RDF graph as a whole, in contrast to the relational case where it is possible to create such statistics with higher precision since data is separated into multiple tables [15]. Furthermore, assuming independence when joining parts of SPARQL queries (triple patterns) leads to erroneous estimations [9] as co-occurrences of certain predicates are highly correlated [19].

Hence, exploiting more fine-grained statistics capturing correlations among RDF triples leads to more accurate join cardinality estimations [19]. However, creating such statistics comes at the price of a very time and resource-intensive preprocessing step. On the other hand, the alternative of online, query-dependent, sampling [20] results in overheads during query optimization. Instead, what we propose in this paper is to better exploit the information that is often provided along with an RDF dataset: SHACL (Shapes Constraint Language) [14] constraints, which is a recent standard for validating RDF datasets that are becoming more and more popular. SHACL defines so-called shapes describing the relationships between entities of a specific class, their properties, and their connections to other classes of entities. Although they are currently only used for validation purposes, we show in this paper that by slightly extending them with basic statistics, they can also be exploited for join cardinality estimation.

In summary, this paper makes the following contributions. First, we extend the SHACL definition to capture statistical information to replace the need for creating complex (and expensive) statistics over RDF datasets. To the best of our knowledge, this is the first proposal of this kind. Second, we introduce an algorithm to enhance SHACL shapes with statistical information and to exploit these statistics for join cardinality estimation and query optimization. Third, we study the impact of our approach using both synthetic (LUBM [10], WatDiv [2]) and real (YAGO-4 [21]) datasets, demonstrating that shapes statistics can provide higher precision for query optimization with only a little overhead.

This paper is structured as follows. While Sections 2 and 3 discuss related work and introduce preliminaries, Section 4 formally defines the problem. Section 5 then describes our proposed extension of the SHACL standards, and Section 6 presents techniques to exploit the additional information for cardinality estimation and query optimization. Section 7 discusses the results of our extensive experimental study, and Section 8 concludes the paper with an outlook to future work.

2 RELATED WORK

Cardinality estimation has been studied extensively in the context of relational databases [20]. For SPARQL queries, existing techniques adapt relational approaches [13, 24] and focus mostly on specific type of queries [19]. Usually, these approaches construct different kinds of single or multidimensional synopses over databases that can be used to estimate cardinalities [23]. While algorithms designed to generate synopsis for unlabelled graphs are not applicable here (as the edges in RDF graphs are labeled), consequently approaches to generate RDF summaries either produce very large summaries [23], have very high computational complexities, or they are unable to preserve the RDF schema while constructing the summaries [23]. Therefore, the most promising approaches aim at using statistics computed directly from edge label frequencies. In particular, RDF-3X proposes a histogram-based technique for cardinality estimation based on edge label frequencies. This technique was later extended by exploiting the statistical information of Characteristic Sets [19], which compute frequencies of sets of predicates sharing the same

subject to estimate the cardinalities. This approach shows high performance for star-shaped queries while it suffers from significant underestimation due to the independence assumption in the general case [20]. This approach was extended as Characteristic Pairs [18] to overcome this limitation, but it could only support multi-chain star queries. Moreover, extracting Characteristic Sets from large heterogeneous graphs is computationally expensive. SumRDF [23] is another cardinality estimation approach based on a graph summarization. It fails to handle large queries due to a prohibitive computation cost, and it is costly to construct such summaries over large RDF graphs [20].

A recent benchmark, G-CARE [20], analyzed the performance of existing cardinality estimation techniques for subgraph matching. This analysis revealed that the techniques based on sampling and designed for online aggregation outperform the cardinality estimation techniques for RDF graphs. *This calls for a more in-depth study on how to perform cardinality estimation for SPARQL query optimization appropriately.*

In a recent work, Shape Expressions (ShEx) [22] have been used to reorder triple patterns to enable SPARQL query optimization [1], i.e., it estimates an order of execution for the triple patterns based on some heuristic inference on which triples are more selective. For instance, if a shape definition says that every instructor has one or more courses, but every course has exactly one instructor, it infers that the cardinality of courses is at least the same as the cardinality of instructors and probably larger. Hence, this optimization procedure is not based on actual data.

Therefore, contrary to existing works, we aim at exploiting fine-grained statistics based on shapes to produce more precise cardinality estimations for query planning. This will allow us to overcome the limitations of existing methods that only use the global-statistics [11]. To this end, instead of creating large expensive summaries and characteristic sets over the RDF graphs to estimate the cardinalities, we exploit SHACL shapes constraints (which are as expressive as ShEx [22]) and annotate the *Node* and *Property Shapes* with the statistics of the input RDF graph. Compared to other solutions, it requires a lightweight preprocessing and retains the structure of original RDF and SHACL shapes graphs. Moreover, *this allow us to study more closely the effect of more fine-grained statistics, and more accurate cardinality estimation for the task of SPARQL query optimization.*

3 PRELIMINARIES

RDF Graphs: RDF graphs model entities and their relationships in the form of triples consisting of SPO \langle subjects, predicates, objects \rangle . We present a simplified example of an RDF graph based on the LUBM [10] dataset in Figure 1, where oval and rectangular shapes represent IRIs and literal nodes, respectively. An RDF graph is formally defined as:

Definition 3.1 (RDF Graph). Given pairwise disjoint sets of IRIs \mathcal{I} , blank nodes \mathcal{B} , and literals \mathcal{L} , an RDF Graph \mathcal{G} is a finite set of RDF triples $\mathcal{T} \subseteq \mathcal{I} \times \mathcal{P} \times (\mathcal{I} \cup \mathcal{B} \cup \mathcal{L})$.

SPARQL: SPARQL [6] is a standard query language for RDF. A SPARQL query consists of a finite set of triple patterns (known as basic graph pattern, BGP) and some conditions that have to be met in order for data to be selected and returned from an RDF graph. Each SPO position in a triple pattern can be concrete (i.e., bound) or a variable (i.e., unbound). The variable names in a SPARQL query are prefixed by a ‘?’ symbol, e.g., ?X. To answer a BGP, we require a *mapping between variables to values in an RDF graph*, all the resulting triples existing in the RDF graph

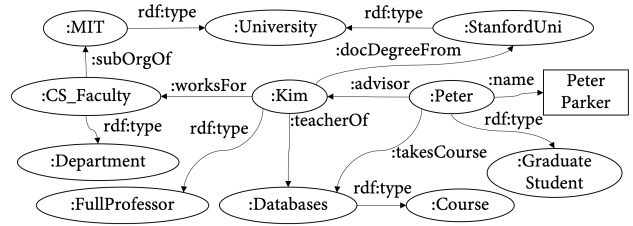


Figure 1: An RDF Graph

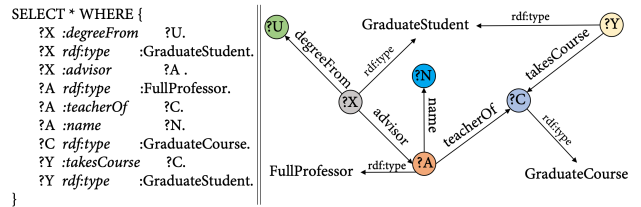


Figure 2: Query & and its Graph &

obtained by replacing the variables with values are answers to the BGP. Figure 2 shows an example SPARQL query (&) and its query graph & on the graph of Figure 1. A BGP is defined as:

Definition 3.2 (BGP). Given a set of IRIs \mathcal{I} , literals \mathcal{L} , and variables \mathcal{V} , a BGP is defined as $\mathcal{G} = \{ (s, p, o) \mid s \in \mathcal{I} \cup \mathcal{V}, p \in \mathcal{P}, o \in \mathcal{I} \cup \mathcal{L} \}$, whose elements are called triple patterns.

Shapes Graphs: Several schema languages have been proposed for RDF in the past, where the most common are RDF Schema (RDFS¹) and OWL [17]. RDFS is primarily used to infer implicit facts, and OWL is an extension of RDF and RDFS to represent ontologies. The declarative Shapes Constraint Language (SHACL) [14] became a W3C standard recently. SHACL schema provides high-level information about the structure and contents of an RDF graph. It allows to define and validate structural constraints over RDF graphs. SHACL models the data in two components: the *data graph* and the *shape graph*. The *data graph* contains the actual data to be validated, while the *shape graph* contains the constraints against which resources in the *data graph* are validated. These constraints are modeled as node and property shapes, which consist of attributes encoding the constraints. The node shapes constraints are applicable on nodes that are instances of a specific type in the *data graph* while the property shapes constraints are applicable to predicates associated with nodes of specific types. We define a SHACL shapes graph as follows:

Definition 3.3 (SHACL Shapes Graph). A SHACL shapes graph \mathcal{S} is an RDF graph describing a set of node shapes \mathcal{N} and a set of property shapes \mathcal{P} , such that $\text{COA64C} : (\mathcal{V} \rightarrow \mathcal{I})$ and $\text{COA64C} : \mathcal{P} \rightarrow \mathcal{V}$ are injective functions mapping each node shape \mathcal{N} (and each property shape \mathcal{P}) to the IRI of a target class and a target predicate in \mathcal{I} respectively, and $q : (\mathcal{V} \rightarrow \mathcal{V})$ is a surjective function assigning to each node shape \mathcal{N} a subset \mathcal{P} of property shapes.

For example in Figure 3, node shape constraints are applicable on node ub: GraduateStudent and its property shapes constraints are applicable on predicates like takesCourse, and advisor. This information is declared with attributes *sh:targetClass* for node shapes and *sh:path* for property shapes. Note that the attributes in the dark shaded boxes are part of our extension of the SHACL definition, explained in Section 5.

¹https://www.w3.org/TR/rdf-schema/

5 EXTENDING SHACL WITH STATISTICS

To compute more accurate join cardinality estimations (Problem 1), we capture the correlations between RDF triples by extending SHACL's node and property shapes with fine-grained statistics of the RDF graph. We denote these statistics as *shapes statistics*. These include the total triple count ($sh:count$), minimum ($sh:minCount$) and maximum ($sh:maxCount$) number of triples for each instance, and the number of distinct objects for property instantiations ($sh:distinctCount$). The attributes shown in the dark shaded boxes in Figure 3 are the annotated statistical attributes of their respective node and property shapes. These statistics are computed by executing analytical SPARQL queries over the RDF graph. For instance, to compute the number of instances of `GraduateStudent` in the dataset, i.e., the value of attribute $sh:count$ of node shape `GraduateStudent`, the annotator issues the SPARQL query: `SELECT COUNT(*) WHERE {?x a ub:GraduateStudent}`. Along with shapes statistics, we also define global statistics by extending VOID statistics with more precise statistics of RDF properties, i.e., the distinct subject count (DSC) and distinct object count (DOC) of each property of the RDF graph.

Figure 3: SHACL Shapes Graph

The Shapes Expression (ShEx) language also serves a similar purpose as SHACL to validate RDF graphs. Nonetheless, the two formulations diverge mostly at the syntactic level [14], and our approach can be extended to work using ShEx or other constraints languages as well without the loss of generality.

4 PROBLEM FORMULATION

Given an input query Q , a query optimizer has the goal to find a query plan expected to answer in the minimum amount of time [15]. Constructing a SPARQL query plan includes finding a join ordering between triple patterns of its BGPs. In this paper, we focus on the join ordering of BGPs defined as follows:

Definition 4.1 (Join Ordering) Given a set of triple patterns $\{t_1, \dots, t_n\}$, the join order O for BGPs is defined as a total ordering of T so that for every $t_i, t_j \in T$, either $t_i \prec t_j$ or $t_j \prec t_i$.

To find an optimal plan, a query optimizer needs to explore the search space of semantically equivalent join orders and choose the optimal (cheapest) plan according to some cost function. It is crucial to accurately estimate the join cardinality between triple patterns of a given query to construct a query plan with an efficient join ordering [9]. In line with the related work [20], we neglect other cost factors and focus on join cardinality as the most dominant cost factor to find a join ordering. We formally define the problem of estimating join cardinalities as follows:

Problem 1 (Join Cardinality Estimation). Given a set of triple patterns $\{t_1, \dots, t_n\}$, apply a cardinality estimation function ρ such that for every pair of triple patterns $t_i, t_j \in T$, $\rho(t_i, t_j) \in \mathbb{N}$.

We extend the above estimation problem also to the case of joining a triple pattern with the intermediate results of prior join operations, e.g., to estimate the total cardinality $\rho(t_i, \rho(t_j, t_k))$. Then, given such estimates, an optimal query plan minimizes the total number of operations to compute, i.e., the execution cost $\sum_{t_i \prec t_j} \rho(t_i, t_j)$. In practice, this total join cost is obtained by summing up the intermediate cardinalities of each join operation in their respective join order. Hence, we formalize the problem of join order optimization as follows:

Problem 2 (Join Order Optimization). Given a set of triple patterns $\{t_1, \dots, t_n\}$ and a join cardinality estimation function ρ , find the join order O obtained as $\arg \min_O \sum_{t_i \prec t_j} \rho(t_i, t_j)$.

6 QUERY PLANNING

In this section, we present our approach to exploit global and shapes statistics to obtain more accurate join cardinality estimates (Problem 1). These estimates, in turn, are used for join order optimization (Problem 2).

6.1 Cardinality Estimation of Triple Patterns

A SPARQL query contains joins between multiple triple patterns. Hence, the first step is to estimate how many triples match every triple pattern individually. We exploit the statistical information contained in the extended SHACL shapes graph (Section 5) to obtain this estimate. Hence, for each triple pattern, we obtain their corresponding node or property shapes using the values of the $sh:targetClass$ and $sh:path$ attributes.

First, all triples of the type $\langle ?x, a, [Class] \rangle$ (i.e., instances with $rdf:type [Class]$) are mapped to the node shape having that class as the value of the attribute $sh:targetClass$. Then, triples having variable $?x$ as a subject are also assigned to that node shape. The triple predicate determines instead its corresponding candidate property shapes, i.e., those with a matching value for $sh:path$. For example, given triple $\langle ?x, rdfs:type, ub:GraduateStudent \rangle$ and $\langle ?x, ub:name, ?n \rangle$, the subject $?x$ is assigned to node shape: `GraduateStudent` while the predicate $rdfs:type$ matches shape: `name` (Figure 3, top left and top right).

Once the candidate shapes for all the triple patterns are identified, their statistical information combined with the distinct subject and object count (DSC & DOC) from global statistics are used in combination with the formulas shown in Table 1 to compute their expected cardinality. These formulas, inspired by a previous work [11], cover all possible types of triple patterns. The term 2 in the formulas denotes the count of t in the RDF graph; 2_{C_1} denotes the count of all triples and 2_{1942C_1} the count of all objects. Similarly, 2 represents the count of having \dots . This can be used, for instance, to derive that there are 85K triples matching $\langle ?x, rdfs:type, ub:FullProfessor \rangle$ (Table 2a).

While both global and shapes statistics can be used to estimate the cardinality of triple patterns using these formulas, they can lead to different estimated cardinalities. When the query does not contain any type-defined triple, only global statistics are used.

²Vocabulary of Interlinked Datasets: <https://www.w3.org/TR/void/>

Triple Pattern	Cardinality	Triple Pattern	Cardinality
?s ?p obj	$\frac{2_{CA8?;4B}}{2_{1942CB}}$?s ?p ?o	$\frac{2_{CA8?;4B}}{2_{38BC(D19)}}$
subj ?p obj	$\frac{2_{CA8?;4B}}{2_{38BC(D19)}}$	subj ?p ?o	$\frac{2_{CA8?;4B}}{2_{7A43}}$
?s pred obj	$\frac{2_{7A43}}{2_{7A43;19}}$?s pred ?o	$\frac{2_{7A43}}{2_{38BC(D19)}}$
subj pred obj	$\frac{2_{7A43}}{2_{38BC(D19)}}$	sub pred ?o	$\frac{2_{7A43}}{2_{A35C-74}}$
?s rdf:type obj	$\frac{2_{4=C8C84BC-74;19}}{1 > A0}$?s rdf:type ?o	$\frac{2_{A35C-74}}{2_{A35C-74d01}}$
subj rdf:type obj	$1 > A0$	subj rdf:type ?o	$\frac{2_{A35C-74d01}}{2_{A35C-74d01}}$

Table 1: Cardinality estimation of triple patterns

6.2 Cardinality Estimation of Joins

The join operation is performed on a common variable between two triple patterns. We consider three possible types of joins between two triple patterns based on the position of the common variable, namely: Subject-Subject (SS), Subject-Object (SO), and Object-Object (OO). If there is no common variable between two triple patterns, the join will result in a Cartesian product. Inspired by related work [8], we estimate the SS, SO, and OO join cardinalities using the formulas stated in Equations 1, 2, and 3. Note that $\$$ and $\$$ in the formulas represent the distinct subject and object count of triple patterns respectively.

$$\frac{2_{0A3?C?}}{<0G^1} Z_{(C?)} = \frac{2_{0A3?} \cdot 2_{0A3?}}{<0G^1 \cdot (\$ \cdot \$)} \quad (1)$$

$$\frac{2_{0A3?C?}}{<0G^1} Z_{(\$ C?)} = \frac{2_{0A3?} \cdot 2_{0A3?}}{<0G^1 \cdot (\$ \cdot \$)} \quad (2)$$

$$\frac{2_{0A3?C?}}{<0G^1} Z_{(\$ \$ C?)} = \frac{2_{0A3?} \cdot 2_{0A3?}}{<0G^1 \cdot (\$ \cdot \$)} \quad (3)$$

6.3 Join Ordering

Given an RDF graph G , its shapes statistics graph B^0 , and global statistics graph B^g , we propose an algorithm to compute the join ordering for an input query Q (Algorithm 1). In the first step, the triple patterns of Q are sorted in ascending order of their estimated cardinalities using only global statistics. The algorithm starts with the triple pattern having the least cardinality and then estimates its join cardinality with the rest of the triple patterns using the formulas from Section 6.2. The algorithm iterates over all the triple patterns and chooses a triple pattern with the least estimated join cardinality (size of intermediate result) given the triple already selected. This produces a first join ordering based on global statistics. In the second step, shapes statistics are taken into account, and both the estimated cardinalities and the join ordering proposed in the first step are revised using these shapes specific fine-grained statistics. The algorithm also computes the cost of each join ordering by adding the estimated join cardinalities in each iteration. Its complexity is cubic to the number of triple patterns in the query, i.e., $O(n^3)$.

Given our example query Q , and the cardinalities of its triple patterns $\{C_1, C_2, C_3, C_4, C_5, C_6\}$ estimated with both global and shape statistics, Tables 2a and 2b show the join ordering computed only using global statistics (O_{GB}) and via shapes statistics (O_{SB}), respectively. There is a significant difference between the estimated and true join cardinalities and their final total cost. The estimated join cardinalities for O_{GB} are much closer to the true cardinalities of the query than the estimates for O_{GB} with two exceptions for C_3 and C_4 where shapes statistics largely overestimate their cardinalities due to skewed distribution of data.

Algorithm 1 Join Ordering

```

Input:  $G, B^0, B^g$ 
Output: Join order  $O$  of  $Q$ 
1:  $Q = \{C_1, \dots, C_n\}$ 
2:  $2 > BC_0 = 0; 20A3 = 0; @D4D4 @D4D4^8=34G$ 
3:  $C?B = 64C \% B^0;$ 
4:  $C?B = 64C = 3830C4(0?4B, B, 6B);$ 
5:  $C?B = 2 > ?DC4 0A38=0; 8C?B^0;$ 
6:  $B > A?B2 \cdot C?B 0A38=0; 8C?B^0;$ 
7:  $? ?033?C?B; A^033 ;; (C?B C?B);$ 
8:  $2 > BC C?B 20A38=0; 8C?B^0;$ 
9:  $@D4D4^03G?B^8=34G;$ 
10: for  $C? C? C? B$  do
11:  $8=34G = C?^8=34G 2 > BC_20; = 2 > BC$ 
12:  $@D4D4 @D4D4$ 
13: while  $!@D4D4^8 < ?C$ -do
14:  $C? = @D4D4^8; ^10;$ 
15: for  $C? 2 A$  do
16:  $2 = 0;$ 
17: if  $tp_0 Z$   $tp_1$  then
18:  $c = ^1C? \cdot C?^0;$ 
19: else  $2 = 2^1C? \cdot C?^0;$ 
20: if  $2 \checkmark 2 > BC_20;$  then
21:  $2 > BC_20; = 2; 8=34G = C?^8=34G 20A3= 2;$ 
22:  $2 > BC_2 > BC_20;$ 
23:  $@D4D4^038=34G; ?^033?C?B 64C8=34G;$ 
24:  $A^0A < > EC?B 64C8=34G;$ 
25:  $O @D4D4^8? > ;$ 

```

7 EXPERIMENTAL EVALUATION

We investigated the performance of query plans proposed using our algorithm (with global and shapes statistics) compared to the plans proposed by two state-of-the-art query engines (Apache Jena ARQ and GraphDB) as well as two state-of-the-art RDF cardinality estimation approaches (Characteristic Sets [19] and SumRDF [23]). All experiments are performed on a single machine with Ubuntu 18.04, having 16 cores and 256GB RAM. Datasets: We used LUBM [10], WatDiv [2], and YAGO-4 [21] to study various query plans on different datasets and sizes (Table 3). In particular, we used LUBM-500, two variants of WatDiv datasets (WATDIV-S (Small) with ~108.9 M triples and WATDIV-L (Large) with 1 billion triples), and for YAGO-4 we used the subset containing instances that have an English Wikipedia article. Implementation: Nowadays, constraints languages are having widespread application to validate RDF graphs [25]. We assume the availability of SHACL shapes graph with the dataset and provide a Shapes Annotator to extend it with statistics of the graph. For cases where they are not present, the SHACL Query Engine is commonly used to generate shapes graphs and we also use it in our case (e.g., for YAGO-4). All shapes are then extended with the required statistics using our Shapes Annotator implemented in Java. The SHACL shapes graph for LUBM, for instance, is 45 KB, and the size of extended shapes is 68 KB. The time required to extend the SHACL shapes depends on the number of its nodes and property shapes. The process of extending LUBM shapes graph took 16 minutes, WATDIV-S took 8 minutes, and for YAGO-4 (which consists of 8888 nodes and 80831 property shapes) it took 62 minutes. We implemented our join ordering algorithm in Java using Jena. The source code is available on our website.

³<https://jena.apache.org/documentation/>
⁴<https://graphdb.ontotext.com>
⁵<https://pypi.org/project/shaclgen/>
⁶<https://relweb.cs.aau.dk/rdfshapes/>

