

Learning SHACL Shapes from Knowledge Graphs

Pouya Ghiasnezhad Omran^{a,*}, Kerry Taylor^a, Sergio Rodríguez Méndez^a, and Armin Haller^a

^a*School of Computing, The Australian National University, ACT, Australia*

E-mails: P.G.Omran@anu.edu.au, Kerry.Taylor@anu.edu.au, Sergio.RodriguezMendez@anu.edu.au, Armin.Haller@anu.edu.au

Abstract. Knowledge Graphs (KGs) have proliferated on the Web since the introduction of knowledge panels to Google search in 2012. KGs are large data-first graph databases with weak inference rules and weakly-constraining data schemes. SHACL, the Shapes Constraint Language, is a W3C recommendation for expressing constraints on graph data as shapes. SHACL shapes serve to validate a KG, to underpin manual KG editing tasks, and to offer insight into KG structure. Often in practice, large KGs have no available shape constraints and so cannot obtain these benefits for ongoing maintenance and extension.

We introduce Inverse Open Path (IOP) rules, a predicate logic formalism which presents specific shapes in the form of paths over connected entities that are present in a KG. IOP rules express simple shape patterns that can be augmented with minimum cardinality constraints and also used as a building block for more complex shapes, such as trees and other rule patterns. We define formal quality measures for IOP rules and propose a novel method to learn high-quality rules from KGs. We show how to build high-quality tree shapes from the IOP rules. Our learning method, SHACLEARNER, is adapted from a state-of-the-art embedding-based open path rule learner (OPRL).

We evaluate SHACLEARNER on some real-world massive KGs, including YAGO2s (4M facts), DBpedia 3.8 (11M facts), and Wikidata (8M facts). The experiments show that our SHACLEARNER can effectively learn informative and intuitive shapes from massive KGs. The shapes are diverse in structural features such as depth and width, and also in quality measures that indicate confidence and generality.

Keywords: SHACL Shape Learning, Shapes Constraint Language, Knowledge Graph, Inverse Open Path Rule

1. Introduction

While public knowledge graphs (KGs) became popular with the development of DBpedia [1] and Yago [2] more than a decade ago, interest in enterprise knowledge graphs [3] has taken off since the inclusion of knowledge panels on the Google Search engine in 2012, driven by its internal knowledge graph. Although these KGs are massive and diverse, they are typically incomplete. Regardless of the method that is used to build a KG (e.g., collaboratively vs individually, manually vs automatically), it will be incomplete because of the evolving nature of human knowledge, cultural bias [4] and resource constraints. Consider Wikidata [5], for example, where there is more

complete information for some types of entities (e.g., pop stars), while less for others (e.g., opera singers). Even for the same type of entity, for example, computer scientists, there are different depths of detail (for similarly accomplished scientists) depending on their country of residence.

However, the power of KGs comes from their data-first approach, enabling contributors to extend a KG in a relatively arbitrary manner. By contrast, a relational database typically employs not-null and other schema-based constraints that require some attributes to be instantiated in a defined way at all times. Large KGs are typically populated by automatic and semi-automatic methods using non-structured sources such as Wikipedia that are prone to errors of omission (i.e., incompleteness) and commission (i.e., falsity). Both kinds of errors can be highlighted for correction

*Corresponding author. E-mail: P.G.Omran@anu.edu.au.

by a careful application of schema constraints. However, such constraints are commonly unavailable and, if available, uncertain and frequently violated in a KG for valid reasons, arising from the intended data-first approach of KG applications.

SHACL [6] was formally recommended by the W3C in 2017 to express such constraints on a KG as *shapes*. For example, SHACL can be used to express that a person (in a specific use case) needs to have a name, birth date, and place of birth, and that these attributes have particular types: a string; a date; and a location. The shapes are used to guide the population of a KG, although they are not necessarily enforced. Typically, SHACL shapes are manually-specified. However, as for multidimensional relational database schemes [7], shapes could, in principle, be inferred from KG data. As frequent patterns, the shapes characterise a KG and can be used for subsequent data cleaning or ongoing data entry. There is scant previous research on this topic [8, 9].

While basic SHACL [6] and its advanced features [10] allows the modelling of diverse shapes including rules and constraints, most of these shapes are previously well known when expressed by alternative formalisms, including closed rules [11], trees, existential rules [12], and graph functional dependencies [13]. We claim that the common underlying form of all these shapes is the *path*, over which additional constraints induce alternative shapes. For example, in DBpedia we see the following path, $\langle dbo : type \rangle _ \langle db : Song \rangle (x) \rightarrow \langle dbo : album \rangle (x, y) \wedge \langle dbo : recordLabel \rangle (y, z)$. which expresses that if an entity x is a song, then x is in an album y which has a record label z .

Since the satisfaction of a less-constrained shape is a necessary condition for satisfaction of a more complex shape (but not a sufficient condition), in this paper we focus on learning paths, the least constrained shape for our purposes. In addition, we learn cardinality constraints that can express, for example, that a song has (hypothetically) at least 2 producers. We also investigate the process of constructing one kind of more complex shape, that is a *tree*, out of paths. For example, we discover a tree about an entity which has song as its type as shown in Fig. 1. In a KG context, the tree suggests that if we have an entity of type song in the KG, then we would expect to have the associated facts too.

In this paper, we present a system, SHACLEARNER, that mines shapes from KG data. For this purpose we propose a predicate calculus formalism in which rules

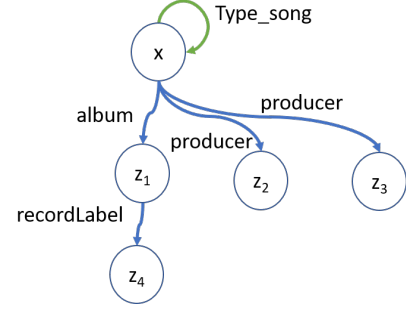


Fig. 1. A tree shape for the Song concept from DBpedia.

have one body atom and a chain of conjunctive atoms in the head with a specific variable binding pattern. Since these rules are an inverse version of *open path rules* [14], we call them *inverse open path (IOP) rules*. To learn IOP rules we adapt an embedding-based open path rule learner, OPRL [14]. We define quality measures to express the validity of IOP rules in a KG. SHACLEARNER uses the mined IOP rules to subsequently discover more complex tree shapes. Each IOP rule or tree is a SHACL shape, in the sense that it can be syntactically rewritten in SHACL. Our mined shapes are augmented with a novel numerical confidence measure to express the strength of evidence in the KG for each shape.

In summary, the main contributions of this paper are:

- We introduce a new formalism called Inverse Open Path rules, that serves as a building block for more complex shapes such as trees, together with cardinality constraints and quality measurements;
- We extend the Open Path rule learning method, OPRL [14, 15], to learn IOP (**Inverse Open Path**) rules annotated with cardinality constraints, while introducing unary predicates that can act as class or type constraints; and
- We propose a method to aggregate IOP rules to produce tree shapes.

This paper is organised as follows. After presenting some foundations in Section 2, we describe our SHACL learning method in Section 3, including the formalism of IOP rules, the embedding-based method that discovers IOP rules from a KG, and the method for aggregating IOP rules into trees. In Section 4, we present related work. We discuss results of an experimental evaluation in Section 5 before we conclude in Section 6.

2. Preliminaries

The presentation of closed path rules and open path rules in this section is adapted and extended from [14].

An *entity* e is an identifier for an object such as a place or a person. A *fact* (also known as a *link*) is an RDF triple (e, P, e') , written here as $P(e, e')$, meaning that the *subject* entity e is related to an *object* entity e' via the binary *predicate* (also known as a *property*), P . In addition, we allow *unary* predicates of the form $P(e)$, also equivalently written here as the binary fact $P(e, e)$. We model unary predicates as self-loops to have a unary predicate act as the label of a link in the graph, as shown in Fig. 1, just as for binary predicates. Unary predicates may, but are not limited to, represent class assertions expressed in an RDF triple as $(e, \text{rdf:type}, P)$ where P is a class or a datatype. A *knowledge graph* (KG) is a pair $K = (E, F)$, where E is a set of entities and F is a set of facts and all the entities occurring in F also occur in E .

2.1. Closed-Path Rules

KG rule learning systems employ various rule languages to express rules. RLvLR [16] and SCALEKB [17] use so-called *closed path* (CP) rules. Each consists of a *head* at the front of the implication arrow and a *body* at the tail. We say the rule is *about* the predicate of the head. The rule forms a closed path, or single unbroken loop of links between the variables. It has the following general form.

$$P_t(x, y) \leftarrow P_1(x, z_1) \wedge P_2(z_1, z_2) \wedge \dots \wedge P_n(z_{n-1}, y). \quad (1)$$

We interpret these kinds of rules with universal quantification of all variables at the outside, and so we can infer a fact that instantiates the head of the rule by finding an instantiation of the body of the rule in the KG. For example, from the rule $\text{citizenOf}(x, y) \leftarrow \text{livesIn}(x, z) \wedge \text{locatedIn}(z, y)$ and the facts in the KG: $\text{livesIn}(\text{Mary}, \text{Canberra})$ and $\text{locatedIn}(\text{Canberra}, \text{Australia})$, we can infer and assert the new fact: $\text{citizenOf}(\text{Mary}, \text{Australia})$.

Rules are considered of more use if they generalise well, that is, they explain many facts. To quantify this idea we recall measures *support*, *head coverage* and *standard confidence* that are used in some major approaches to rule learning including [17] and [11].

Definition 1 (satisfies, support). *Let r be a CP rule of the form (1). A pair of entities (e, e') satisfies the body of r , denoted $\text{body}_r(e, e')$, if there exist entities e_1, \dots, e_{n-1} in the KG such that all of $\{P_1(e, e_1), P_2(e_1, e_2), \dots, P_n(e_{n-1}, e')\}$ are facts in the KG. Further (e, e') satisfies the head of r , denoted $P_t(e, e')$, if $P_t(e, e')$ is a fact in the KG. Then the support of r counts the rule instances for which the body and the head are both satisfied in the KG.*

$$\text{supp}(r) = |\{(e, e') : \text{body}_r(e, e') \text{ and } P_t(e, e')\}|$$

Standard confidence (SC) and *head coverage* (HC) offer standardised measures for comparing rule quality. SC describes how frequently the rule is true, i.e., of the number of entity pairs that satisfy the body in the KG, what proportion of the inferred head instances are satisfied? It is closely related to *confidence* widely used in association rule mining [18]. HC measures the explanatory power of the rule, i.e., what proportion of the facts satisfying the head of the rule could be inferred by satisfying the rule body? It is closely related to *cover* which is widely used for rule learning in inductive logic programming [19]. A non-recursive rule that has both 100% SC and HC is redundant with respect to the KG, and every KG fact that is an instance of the rule head is redundant with respect to the rule.

Definition 2 (standard confidence, head coverage). *Let r, e, e', body_r be as given in definition 1. Then standard confidence is*

$$SC(r) = \frac{\text{supp}(r)}{|\{(e, e') : \text{body}_r(e, e')\}|}$$

and head coverage is

$$HC(r) = \frac{\text{supp}(r)}{|\{(e, e') : P_t(e, e')\}|}$$

2.2. Open-Path Rules: Rules with Open Variables

Unlike earlier work in rule mining for KG completion, OPRL for *active* knowledge graph completion [14] defines *open path* (OP) rules of the form:

$$P_1(z_0, z_1) \wedge P_2(z_1, z_2) \wedge \dots \wedge P_n(z_{n-1}, y) \rightarrow \exists x P_t(x, z_0). \quad (2)$$

Here, P_i is a predicate in the KG, each of $\{x, z_i, y\}$ are entity variables, and all free variables are universally quantified at the outside.

We call a variable *closed* if it occurs in at least two distinct predicate terms, such as z_0 here, and otherwise it is *open*. If all variables of a rule are closed then the rule is closed. An OP rule has two open variables, y and x . [14] uses OP rules since they imply the existence of a fact, like $spouse(x, y) \rightarrow \exists z child(x, z)$ [11]. Unlike CP rules, OP rules do not necessarily form a loop, but a straightforward variable unification transforms an OP rule to a CP rule, and every entity-instantiation of a CP rule is also an entity-instantiation of the related OP rule (but not vice-versa).

To assess the quality of OP rules, we use *open path standard confidence (OPSC)* and *open path head coverage (OPHC)* which are derived in [14] from the closed path forms (Definition 2).

Definition 3 (open path: OPsupp, OPSC, OPHC). *Let r be an OP rule of the form (2). Then a pair of entities (e, e') satisfies the body of r , denoted $body_r(e, e')$, if there exist entities e_1, \dots, e_{n-1} in the KG such that $P_1(e, e_1), P_2(e_1, e_2), \dots, P_n(e_{n-1}, e')$ are facts in the KG. Also (e', e) satisfies the head of r , denoted $P_t(e', e)$, if $P_t(e', e)$ is a fact in the KG. The open path support, open path standard confidence, and open path head coverage of r are given respectively by*

$$OPsupp(r) = |\{e : \exists e', e'' \text{ s.t. } body_r(e, e') \text{ and } P_t(e'', e)\}|$$

$$OPSC(r) = \frac{OPsupp(r)}{|\{e : \exists e' \text{ s.t. } body_r(e, e')\}|}$$

$$OPHC(r) = \frac{OPsupp(r)}{|\{e : \exists e' \text{ s.t. } P_t(e', e)\}|}$$

2.3. SHACL Shapes

A KG is a schema-free database and does not need to be augmented with schema information natively. However, many KGs are augmented with type information that can be used to understand and validate data and can also be very helpful for inference processes on the KG. In 2017 the Shapes Constraint Language (SHACL) [6] was introduced as a W3C recommendation to define schema information for KGs stored as an RDF graph. SHACL defines constraints for graphs as *shapes*. KGs can then be validated against a set of shapes.

Shapes can serve two main purposes: validating the quality of a KG and characterising the frequent patterns in a KG. In Fig. 2, we illustrate an example of

a shape from Wikidata¹, where x and z_i s are variables that are instantiated by entities. Although the shape is originally expressed in ShEx [20], we translate it to SHACL in Fig. 2.

SHACL, together with SHACL advanced features [10] is extensive. Here we focus on the core of SHACL in which *node* shapes constrain a *target* predicate (e.g., the unary predicate *human* in Fig. 2), with *property* shapes expressing constraints over facts related to the target predicate. We particularly focus on property shapes which act to constrain an argument of the target predicate. In Fig. 2 the shape expresses that each entity x which satisfies $human(x)$ should satisfy the following paths: (1) $citizenOf(x, z_1) \wedge country(z_1)$, (2) $father(x, z_2) \wedge human(z_2)$, and (3) $nativeLanguage(x, z_3) \wedge language(z_3)$.

```

humanShape
  a sh:NodeShape;
  sh:targetClass class:human ;
  sh:property [
    sh:path property:citizenOf ;
    sh:class class:country ;
    sh:nodeKind sh:IRI ;
  ];
  sh:property [
    sh:path property:father ;
    sh:class class:human ;
    sh:nodeKind sh:IRI ;
  ];
  sh:property [
    sh:path property:nativeLanguage ;
    sh:class class:language ;
    sh:nodeKind sh:IRI ;
  ];

```

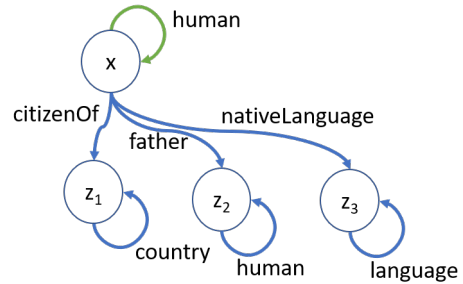


Fig. 2. An example of a SHACL shape from Wikidata.

Various formalisms with corresponding shapes have been proposed to express diverse kinds of patterns exhibited in KGs, such as k-cliques, Closed rules

¹<https://www.wikidata.org/wiki/EntitySchema:E10>

(CR) [11] (that include closed path rules), Functional Graph Dependency (FGD) [13], and trees [12]. CRs are used for inferring new facts. FGDs define constraints like the type of entities in the domain and range of predicates, or the number of entities to which an entity can be related by a specific predicate. These constraints are deployed to make the KG consistent. Regardless of differences amongst these formalisms, they share a key feature in their syntax. The building block for expressing all of these shape constraints is a sequence of predicates.

We focus on such path shapes for our shape learning system. A *path* is a sequence of predicates connected by closed intermediate variables but terminating with open variables at both ends. Although shapes in the form of a path are less constrained than some more complex shapes, they are a more general template for more complex shapes like closed rules or trees, which are also paths (with further restrictions). We will define *Inverse Open Path* rules induced from paths that have a straightforward interpretation as shapes, and also propose a method to mine such rules from a KG. To demonstrate the potential for these kind of shapes to serve as building blocks for more complex shapes, we then propose a method that builds trees out of mined rules, and briefly discuss the application of such trees to KG completion.

3. SHACL learning

3.1. Rules with Open Variables or Uncertain Shapes

We observe that the converse of OP rules, that we call *inverse open path rules* (IOP), correspond to SHACL shapes. For example, the shape of Fig. 2 can be derived from the following three IOP rules:

$$\text{human}(x, x) \rightarrow \text{citizenOf}(x, z_1) \wedge \text{country}(z_1, z_1).$$

$$\text{human}(x, x) \rightarrow \text{father}(x, z_2) \wedge \text{human}(z_2, z_2).$$

$$\text{human}(x, x) \rightarrow \text{nativeLanguage}(x, z_3) \wedge \text{language}(z_3, z_3).$$

The general form of an IOP rule is given by

$$P'_i(x, z_0) \rightarrow \exists(z_1, \dots, z_{n-1}, y) P'_1(z_0, z_1) \wedge P'_2(z_1, z_2) \wedge \dots \wedge P'_n(z_{n-1}, y). \quad (3)$$

where each P'_i is either a predicate in the KG or its reverse with the subject and object bindings swapped,

and free variables are universally quantified at the outside. We often omit the quantifiers when writing IOP rules. In an IOP rule the body of the rule is P_i and its head is the sequence of predicates, $P_1 \wedge P_2 \wedge \dots \wedge P_n$. Hence we instantiate the atomic body to predict an instance of the head. IOP rules are not Horn. The pattern of existential variables in the head and universal variables in the body has been investigated in the literature as *existential rules* [12].

To assess the quality of IOP rules we follow the style of quality measures for OP rules [14].

Definition 4 (inverse open path: IOPsupp, IOPSC, IOPHC). *Let r be an IOP rule of the form (3). Then a pair of entities (e, e') satisfies the head of r , denoted $\text{head}_r(e, e')$, if there exist entities e_1, \dots, e_{n-1} in the KG such that $P_1(e, e_1), P_2(e_1, e_2), \dots, P_n(e_{n-1}, e')$ are facts in the KG. A pair of entities (e'', e) satisfies the body of r , denoted $P_i(e'', e)$, if $P_i(e'', e)$ is a fact in the KG. The inverse open path support, inverse open path standard confidence, and inverse open path head coverage of r are given respectively by*

$$\text{IOPsupp}(r) = |\{e : \exists e', e'' \text{ s.t. } \text{head}_r(e, e') \text{ and } P_i(e'', e)\}|$$

$$\text{IOPSC}(r) = \frac{\text{IOPsupp}(r)}{|\{e : \exists e'' \text{ s.t. } P_i(e'', e)\}|}$$

$$\text{IOPHC}(r) = \frac{\text{IOPsupp}(r)}{|\{e : \exists e' \text{ s.t. } \text{head}_r(e, e')\}|}$$

Notably, because any instantiated open path in a KG induces both an OP and IOP rule: the support for an IOP rule is the same as the support for its corresponding OP form; IOPSC is the same as the corresponding OPSC; and IOPHC is the same as the corresponding OPSC. This close relationship between OP and IOP rules helps us to mine both OP and IOP rules in the one process.

We show the relationship between an OP rule and its converse IOP version in the following example. Consider the OP rule, $P_1(x, z_0) \leftarrow P_2(z_0, z_1) \wedge P_3(z_1, z_2)$. Assume we have three entities $(\{e_3, e_4, e_5\})$ which can instantiate z_0 and satisfy both $P_1(x, z_0)$ and $P_2(z_0, z_1) \wedge P_3(z_1, z_2)$. Assume the number of entities that can instantiate z_0 to satisfy the head part is 5 $(\{e_1, e_2, e_3, e_4, e_5\})$ and the number of entities that can instantiate z_0 to satisfy the body part is 7 $(\{e_3, e_4, e_5, e_6, e_7, e_8, e_9\})$. Hence, we have the following for this rule, $\text{OPsupp} = 3$, $\text{OPSC} = 3/7$ and $\text{OPHC} = 3/5$. For the IOP version of the rule over

the same KG, $P_1(x, z_0) \rightarrow P_2(z_0, z_1) \wedge P_3(z_1, z_2)$, we have the same entities to instantiate z_0 while the predicate terms corresponding to the bodies and heads are swapped. Hence, we have the following for the IOP version of the rule, $IOP_{supp} = 3$, $IOPSC = 3/5$ and $IOPHC = 3/7$.

In many cases we need the rules to express not only the necessity of a chain of facts (the facts in the head of the IOP rule) but the number of different chains which should exist. For example, we may need a rule to express that each human has at least two parents. Hence, we introduce IOP rules annotated with a cardinality, Car . This gives us the following annotated form for each IOP rule.

$$IOPSC, IOPHC, Car : P'_t(x, z_0) \rightarrow P'_1(z_0, z_1) \wedge P'_2(z_1, z_2) \wedge \dots \wedge P'_n(z_{n-1}, y). \quad (4)$$

where $IOPSC$ and $IOPHC$ belong to $[0, 1]$ and denote those qualities of the rule. Car is an integer ≥ 1 .

Definition 5 (Cardinality of an IOP rule, Car). *Let r be an annotated IOP rule of the form (4) and let $Car(r)$ be the cardinality annotation for r . Then r satisfies $Car(r)$ iff for each entity $e \in \{e \mid \exists e'' \text{ s.t. } P_t(e'', e)\}$, $Car(r) \leq |\{e' \mid head_r(e, e')\}|$.*

The cardinality expresses a lower bound on the number of head paths which are satisfied in the KG for every instantiation of the variable that joins the body to the head. For example, if we have $0.8, 0.1, 2 : P_t(x, y) \rightarrow P_1(y, z) \wedge P_2(z, t)$ and have $P_t(e_1, e_2)$ satisfied, we should have at least two paths starting from e_2 that instantiate two distinct entities for variable t , of the form $P_1(e_2, z) \wedge P_2(z, e_3)$ and $P_1(e_2, z) \wedge P_2(z, e_4)$. There is no limitation on the instantiations for variable z which is scoped inside the head, so these two pairs of instantiations both satisfy cardinality of 2: (1) $P_1(e_2, e_5) \wedge P_2(e_5, e_3)$ and $P_1(e_2, e_5) \wedge P_2(e_5, e_4)$ and (2) $P_1(e_2, e_5) \wedge P_2(e_5, e_3)$ and $P_1(e_2, e_6) \wedge P_2(e_6, e_4)$.

Rules with the same head and the same body may have different cardinalities. In the given example we might have the following rule as well, $0.9, 0.1, 1 : P_t(x, y) \rightarrow P_1(y, z) \wedge P_2(z, t)$. While we have a rule with a cardinality, c_1 , we also have rules with lower cardinalities $1, \dots, (c_1 - 1)$ but their $IOPSC$ s should be as good or better than the rule with c_1 cardinality. This statement is quite intuitive since cardinality expresses a lower bound on the number of instances of the head in the KB.

Lemma 1 ($IOPSC$ is non-increasing with length). *Let r be an IOP rule of the form (3) with $n \geq 2$ and let r' be an IOP rule of the form $P'_t(x, z_0) \rightarrow P'_1(z_0, z_1) \wedge P'_2(z_1, z_2) \wedge \dots \wedge P'_n(z_{n-2}, y)$, being r shortened by the removal of the last head predicate. Then $IOPSC(r) \leq IOPSC(r')$.*

Proof. Observe that by definition 4, the denominator of $IOPSC()$ is not affected by the head of the rule, and so has the same value for both. Now, looking at the numerators, we have that

$$\forall e (\exists e_1, \dots, e_n (P'_1(e, e_1) \wedge P'_2(e_1, e_2) \wedge \dots \wedge P'_n(e_{n-2}, e_{n-1}) \wedge P'_n(e_{n-1}, e_n))) \implies \exists e_1, \dots, e_{n-1} (P'_1(e, e_1) \wedge P'_2(e_1, e_2) \wedge \dots \wedge P'_n(e_{n-2}, e_{n-1}))).$$

Therefore $\{e : \exists e' \text{ s.t. } head_r(e, e')\} \subseteq \{e : \exists e' \text{ s.t. } head_{r'}(e, e')\}$.

Therefore $\{e : \exists e', e'' \text{ s.t. } head_r(e, e') \text{ and } P_t(e'', e)\} \subseteq \{e : \exists e', e'' \text{ s.t. } head_{r'}(e, e') \text{ and } P_t(e'', e)\}$.

So $IOP_{supp}(r) \leq IOP_{supp}(r')$ as required. \square

This Lemma is useful in the process of rule learning since it shows that if we discard a rule with a length regarding its low $IOPSC$, v , we do not need to check the extension of this rule with more head atoms since those rules' $IOPSC$ s are equal or less than v . Hence the process of pruning such rules does not harm the completeness of the rule learning.

3.2. IOP Learning through Representation Learning

To mine IOP rules, we start with the open path rule learner OPRL [14], and adapt its embedding-based OP rule learning to learn annotated IOP rules. We call this new IOP rule learner, SHACLEARNER, shown in Algorithm 1.

SHACLEARNER (Algorithm 1) uses a sampling method `Sampling()`, Algorithm 2, to prune the entities and predicates that are less relevant to the target predicate to obtain a sampled KG. The sample is fed to an embedding learner, RESCAL [21, 22], in `Embeddings()`. Then in `PathFinding()`, SHACLEARNER uses the computed embedding representations of predicates and entities in heuristic functions that inform the generation of IOP rules bounded by a maximum length. Then, potential IOP rules are evaluated, annotated, and filtered in `Ev()` to produce annotated IOP rules. Eventually, a tree is discovered for each argument of each target predicate by aggregating mined IOP rules in `GreedySearch()`.

While the overall algorithm structure of SHACLEARNER is similar to OPRL [14], as is the embedding-

Algorithm 1 SHACLEARNER

Input: a KG K , a target predicate P_t
Parameters: max rule length l , max rule cardinality $MCar$, $MinIOPSC$, $MinIOPHC$, and $MinTreeSC$
Output: a set of IOP rules R and $Tree$

- 1: $K' := \text{Sampling}(K, P_t)$
- 2: $(\mathcal{P}, \mathcal{A}) := \text{Embeddings}(K')$
- 3: $R' := \emptyset$
- 4: **for** $2 \leq k \leq l$ **do**
- 5: Add PathFinding($K', P_t, \mathcal{P}, \mathcal{A}, k$) to R'
- 6: **end for**
- 7: $R := \text{Ev}(R', K, MCar, MinIOPSC, MinIOPHC)$
- 8: $Tree := \text{GreedySearch}(R, MinTreeSC)$
- 9: **return** $Tree$ and R

based scoring function, the following elements are novel in SHACLEARNER:

- OPRL cannot handle unary predicates while SHACLEARNER admits unary predicates both in the head and the body of IOP rules.
- SHACLEARNER can discover and evaluate IOP rules with minimum cardinality constraints in the head of the IOP rule, while OPRL is effectively limited to learning the special case of minimum cardinality 1 for all rules. For this reason, the evaluation method of SHACLEARNER, $\text{Ev}()$, differs from the OPRL evaluation module.
- The aggregation module that produces trees out of learnt IOP rules, ready for translation to SHACL, is novel in SHACLEARNER.

3.3. Sampling

In more detail, the $\text{Sampling}()$ method in line 1 of Algorithm 1 computes a fragment of the KG (K') consisting of a bounded number of entities that are related to the goal predicate (i.e., P_t). This sampling is essential since embedding learners (e.g., HOLE [21] and RESCAL) cannot handle massive KGs with millions of entities (e.g., YAGO2).

The sampling method, first introduced in [16], is shown in Algorithm 2. Since we search for IOP rules with up to l atoms (including the specific body target predicate, P_t), the entity set E_{sample} and corresponding fact set K' contains the information needed for learning such rules. Predicates less relevant to the target predicate are pruned in the sampling process and no facts about those predicates remain in K' . As discussed

Algorithm 2 Sampling

Input: a KG K , a target predicate P_t
Parameters: max rule length l
Output: K' a subgraph of the input graph

- 1: $E_1 = \{e \mid \exists e' : P_t(e, e') \vee P_t(e', e)\}$
- 2: **for** $2 \leq k \leq l$ **do**
- 3: $E_k = \{e \mid \exists e' : (e' \in E_{k-1}) \wedge (P_t(e, e') \vee P_t(e', e))\}$
- 4: **end for**
- 5: $E_{sample} = \bigcup_1^l E_i$
- 6: $K' := \{P_t(e, e') \mid (e \in E_{sample}) \wedge e' \in (E_{sample})\}$
- 7: **return** K'

in [23], consider E is the set of entities, and F is the set of facts, computing the set of sampled entities regarding a target predicate demands $2l \cdot |F| \cdot |E|$ in the worst case where l is the parameter of sampling which indicates the max length of the rules. Hence, the complexity of the sampling algorithm is square in terms of $|K|$, $K = (E, F)$. However, in the worst case, the sampled KG size is the same as the original KG, while in the real-world KGs, the number of entities that belong to E_0 is restricted and the neighbors of E_0 is restricted and so on. Hence, with the given fixed number of l the number of entities in sampled KG is far smaller than the original KG.

3.4. Embeddings

After sampling, in line 2 $\text{Embeddings}()$, we compute predicate embeddings as well as subject and object argument embeddings for all predicates in the sampled K' , as is done in RLvLR [16]. The embedding is obtained from RESCAL [21, 22] as it can provide an extensive representation of predicates and entities as is shown in previous heuristic rule learners like [23].

Briefly, we use RESCAL [21] to embed each entity e_i to a vector $\mathbf{E}_i \in \mathbb{R}^d$ and each predicate P_k to a matrix $\mathbf{P}_k \in \mathbb{R}^{d \times d}$ where \mathbb{R} is the set of real numbers and d is an integer (a parameter of RESCAL). For each given fact $P_0(e_1, e_2)$, the following scoring function is computed:

$$f(e_1, P_0, e_2) = \mathbf{E}_1^T \cdot \mathbf{P}_0 \cdot \mathbf{E}_2 \quad (5)$$

The scoring function indicates the plausibility of the fact that e_1 has relation P_0 with e_2 .

In Embeddings() we additionally compute argument embeddings according to the method proposed in [16]. To compute the subject (respectively object) argument embeddings of a predicate P_k , we aggregate the embeddings of entities that occur as the subject (respectively object) of P_k in the KG. Hence for each predicate P_k we have two vectors, \mathbf{P}_k^1 and \mathbf{P}_k^2 that present the subject argument and object argument of P_k respectively.

3.5. Generating and Pruning Rules

After that, in line 3 to line 7 of Algorithm 1, PathFinding() produces candidate IOP rules based on the embedding representation of the predicates which are involved in each rule. The candidate rules are pruned by the scoring function heuristic proposed in [14] for OP rules. Due to the close relationship between OP and IOP rules, a high-scoring candidate OP rules suggests both a good OP rule and a good IOP rule.

An IOP rule $P_t(x, y) \rightarrow P_1(y, z) \wedge P_2(z, t)$. acts to connect entities satisfying the subject argument of the body predicate, P_t , to entities forming the object argument of the last predicate, P_2 , along a path of entities that satisfy a chain of predicates in the rule. There is a relationship between the logical statement of the rule and the following properties in the embedding space [14]:

1. The predicate arguments that have the same variable in the rule should have similar argument emdeddings. For example we should have the following similarities, $\mathbf{P}_t^2 \sim \mathbf{P}_1^1$ and $\mathbf{P}_1^2 \sim \mathbf{P}_2^1$.
2. The whole path forms a composite predicate, like $P^*(x, t) = P_t(x, y) \wedge P_1(y, z) \wedge P_2(z, t)$. We compute the embedding representation of the composite predicate based on its components, $\mathbf{P}^*(x, t) = \mathbf{P}_t(x, y) \cdot \mathbf{P}_1(y, z) \cdot \mathbf{P}_2(z, t)$. Now we could check the plausibility of $P^*(x, t)$ for any pair of entities by equation 5. However, since we are interested in the existence of an entity-free rule, the following similarity will hold: $1 \approx \mathbf{P}_t^1 \cdot \mathbf{P}_t \cdot \mathbf{P}_1 \cdot \mathbf{P}_2 \cdot \mathbf{P}_2^2$.

Based on the above two properties, [14] defines two scoring functions that help us to heuristically mine the space of all possible IOP rules to produce a reduced set of candidate IOP rules.

The ultimate evaluation of an IOP rule will be done in the next step as described below.

3.6. Efficient Computation of Quality Measures

Now we focus our attention on the efficient matrix-computation of the quality measures that are novel for SHACLEARNER. Ev() in Algorithm 1 first evaluates candidate rules on the small sampled KG, and selects only the rules with $IOP_{supp}(r) \geq 1$. They may still include a large number of redundant and low quality rules and so are further downselected based on their IOPSC and IOPHC calculated over the full KG. We show in the following how to efficiently compute the measures over massive KGs using an *adjacency matrix* representation of the KG.

Let $K = (E, F)$ with $E = \{e_1, \dots, e_n\}$ be the set of all entities and $\mathbb{P} = \{P_1, \dots, P_m\}$ be the set of all predicates in F . Like RESCAL [21, 22], we represent K as a set of square $n \times n$ adjacency matrices by defining the function \mathcal{A} . Specifically, the $[i, j]$ th element $\mathcal{A}(P_k)[i, j]$ is 1 if the fact $P_k(e_i, e_j)$ is in F ; and 0 otherwise. Thus, $\mathcal{A}(P_k)$ is a matrix of binary values and the set $\{\mathcal{A}(P_k) \mid k \in \{1, \dots, m\}\}$ represents K .

We illustrate the computation of IOPSC and IO-PHC through an example. Consider the IOP rule $r : P_t(x, z_0) \rightarrow P_1(z_0, z_1) \wedge P_2(z_1, y)$. Let $E = \{e_1, e_2, e_3\}$ and

$$F = \{P_1(e_1, e_2), P_1(e_2, e_1), P_1(e_2, e_3), P_1(e_3, e_1), \\ P_2(e_1, e_2), P_2(e_3, e_2), P_2(e_3, e_3), P_t(e_1, e_3), \\ P_t(e_3, e_2), P_t(e_3, e_3)\}$$

The adjacency matrices for the predicates P_1 , P_2 and P_t are:

$$\mathcal{A}(P_1) : \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}, \mathcal{A}(P_2) : \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 1 & 1 \end{bmatrix},$$

$$\mathcal{A}(P_t) : \begin{bmatrix} 0 & 0 & 1 \\ 0 & 0 & 0 \\ 0 & 1 & 1 \end{bmatrix}$$

For IOPSC and IOPHC (definition 4) we need to calculate (1) the number of entities that satisfy the body of the rule, i.e.,

$\#e : \exists e'' \text{ s.t. } P_t(e'', e)$, (2) the number of entities that satisfy the head of a rule i.e., $\#e : \exists e' \text{ s.t. } head_r(e, e')$ and, (3) the number of entities that join the head of a rule to its body i.e.,

$\#e : \exists e', e'' \text{ s.t. } head_r(e, e') \text{ and } P_t(e'', e)$.

For (1) we can read the pairs (e'', e) directly from the matrix $\mathcal{A}(P_t)$. To find distinct es we sum each *column* (corresponding to each value for the second argument) and transpose to obtain the vector $\mathcal{V}^2(P_t)$. Each non-zero element of this vector indicates a satisfying e and the number of distinct es is given by simply counting the number of non-zero elements. In the example, the only non-zero element in $\mathcal{A}(P_t)$ is $\mathcal{A}(P_t)[1, 3]$ and after summing the columns and transposing we have

$$\mathcal{V}^2(P_t) = \begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix}$$

so we have only $\mathcal{V}^2(P_t)[3]$ non-zero and $\{e_2, e_3\}$ satisfies the head with count 2.

For (2) the pairs (e, e') satisfying the head are connected by the path P_1, P_2, \dots, P_m , and can be obtained, as for (1), directly from the matrix product $\mathcal{A}(P_1)\mathcal{A}(P_2)\dots\mathcal{A}(P_m)$, being the elements with a value $\geq Car$ (for rules with cardinality Car). To find distinct es we sum each *row* (corresponding to each value for the first argument) to obtain the vector $\mathcal{V}^1(\mathcal{A}(P_1)\mathcal{A}(P_2)\dots\mathcal{A}(P_m))$. Each element with $\geq Car$ value of this vector indicates a satisfying e and the number of distinct es is given by counting the number of elements in $\mathcal{V}^1(\mathcal{A}(P_1)\mathcal{A}(P_2)\dots\mathcal{A}(P_m))$.

In the example we have

$$\mathcal{A}(P_1)\mathcal{A}(P_2) = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 2 & 1 \\ 0 & 1 & 0 \end{bmatrix}, \mathcal{V}^1(P_1, P_2) = \begin{bmatrix} 0 \\ 3 \\ 1 \end{bmatrix}$$

with satisfying entities e_2 and e_3 and count of 2 for $Car = 1$. For $Car = 2$, e_2 satisfies the rule so the count is 1.

Computing (3) is now straightforward. We have that the row index of non-zero elements of $\mathcal{V}^2(P_t)$ indicates entities that satisfy the second argument of the body and that row index of elements with a value $\geq Car$ of $\mathcal{V}^1(\mathcal{A}(P_1)\mathcal{A}(P_2)\dots\mathcal{A}(P_m))$ indicates entities that satisfy the first argument of the head. Therefore we can find the entities that satisfy both of these conditions by pairwise multiplication. That is, the entities we need are $\{e_i \mid (\mathcal{V}^2(P_t)[i] > 0 \wedge \mathcal{V}^1(\mathcal{A}(P_1)\mathcal{A}(P_2)\dots\mathcal{A}(P_m))[i]) \geq Car \text{ and } 1 \leq i \leq n\}$, and the count of entities is the size of this set. For the example, for $Car = 1$, we have $\{e_2, e_3\}$ in the set with count 2 and for $Car = 2$, we have $\{e_2\}$ in the set with count 1.

Hence, we could have three versions of r , i.e., r^1, r^2 and r^3 with three different Car of 1, 2, and 3 respectively. For $Car = 1$, $IOPsupp(r^1) = |\{e_2, e_3\}| = 2$. From Definition 4 we can obtain $IOPHC(r^1) = 2/2$ and $IOPSC(r^1) = 2/2$. For $Car = 2$, $IOPsupp(r^2) = |\{e_2\}| = 1$. We can obtain $IOPHC(r^2) = 1/1$ and $IOPSC(r^2) = 1/2$. In this case, we have the same qualities for $Car = 3$ as we have for $Car = 2$.

3.7. From IOP Rules to Tree Shapes

Now in line 8 of Algorithm 1 we turn to deriving SHACL trees, as illustrated in Figure 3, from annotated IOP rules. This procedure is novel in SHACLLEARNER. We use a greedy search to aggregate the IOP rules for the subject argument and the object argument of each target predicate.

For example, the shape of Fig. 2 has the following tree:

$$\begin{aligned} human(x, x) \rightarrow & citizenOf(x, z_1) \wedge country(z_1, z_1) \\ & \wedge father(x, z_2) \wedge human(z_2, z_2) \\ & \wedge nativeLanguage(x, z_3) \\ & \wedge language(z_3, z_3). \end{aligned}$$

The general form of a tree is given by

$$\begin{aligned} P'_t(x, z_0) \rightarrow & \exists(z_*^*, y^*) \\ & P'_1(z_0, z_1^1) \wedge P'_2(z_1^1, z_2^1) \wedge \dots \wedge P'_n(z_{n-1}^1, y^1) \\ & \wedge P'_1(z_0, z_1^2) \wedge P'_2(z_1^2, z_2^2) \wedge \dots \wedge P'_m(z_{m-1}^2, y^2) \\ & \dots \\ & \wedge P'_1(z_0, z_1^q) \wedge P'_2(z_1^q, z_2^q) \wedge \dots \wedge P'_t(z_{t-1}^q, y^q) \end{aligned} \quad (6)$$

where each P'_i^j is either a predicate in the KG or its reverse with the subject and object bindings swapped. Free variables are universally quantified at the outside. In a tree we say the body of the shape is P_t and its head is the sequence of paths or branches, $Path^1 \wedge Path^2 \wedge \dots \wedge Path^q$. Hence we instantiate the atomic body to predict an instance of the head. All head branches and the body join in one shared variable, z_0 .

To assess the quality of a tree we follow the quality measures for IOP rules.

Definition 6 (Tree: Treesupp, TreeSC). Let r be a tree of the above form. Then a set of pairs of entities $(e, e^1), \dots, (e, e^q)$ satisfies the head of r , denoted $head_r(e)$, if there exist the following sequences of entities $e_1^1, \dots, e_{n-1}^1, e_1^2, \dots, e_{m-1}^2$ and e_1^q, \dots, e_{i-1}^q in the KG such that $P_1^1(e, e_1^1), P_2^1(e_1^1, e_2^1), \dots, P_n^1(e_{n-1}^1, e^1), P_1^2(e, e_1^2), P_2^2(e_1^2, e_2^2), \dots, P_m^2(e_{m-1}^2, e^2)$ and $P_1^q(e, e_1^q), P_2^q(e_1^q, e_2^q), \dots, P_i^q(e_{i-1}^q, e^q)$ are facts in the KG. A pair of entities (e'', e) satisfies the body of r , denoted $P_t(e'', e)$, if $P_t(e'', e)$ is a fact in the KG. The tree support and tree standard confidence of r are given respectively by

$$Treesupp(r) = |\{e : \exists e'' \text{ s.t. } head_r(e) \text{ and } P_t(e'', e)\}|$$

$$TreeSC(r) = \frac{Treesupp(r)}{|\{e : \exists e'' \text{ s.t. } P_t(e'', e)\}|}$$

To learn each tree we employ a greedy search, GreedySearch, in line 8 of Algorithm 1. To do so, we sort all rules that bind the subject argument (for the left hand tree in Fig. 3) in a non-increasing order with respect to IOPSC. Then we iteratively try to add the first rule in the list to the tree and compute the TreeSC. If the TreeSC drops below the defined threshold (i.e., $TreeSC_{MIN}$) we dismiss the rule otherwise we add it to the tree. For the right hand tree we do the same with the rules that bind the object argument of the target predicate. Since a conjunction of IOP rules form a tree, TreeSC is bounded above by the minimum IOPSC of its constituent IOP rules.

The uncertain shapes can be presented as standard SHACL shapes by applying the desired minimum confidence by the user. Using this threshold, we can generate hard shapes. Hard shapes have complete qualities and are not augmented with any properties like IOPSC or TreeSC. Aside from the cardinality, the tree may be straightforwardly interpreted as a set of SHACL shapes by reading off every path from the target predicate terminating at a node in the tree. The body predicate is declared a `sh:nodeShape` and the path of head predicates as nested `sh:path` declarations within a `sh:property` declaration. Cardinality of a path is read from the annotation of the branch at the terminating node, and declared by `sh:minCount` within the property declaration. See section 3.1 for an example.

SHACLEARNER supports the *SHACL Core*² [6] features (node and property shapes). The exceptions are: (1) it treats all properties (object and datatype) as "plain predicates" (there is no distinction of their nature), (2) it covers all forms of non-empty paths (paths with a length greater or equal to one), (3) only `min` cardinality, and (4) does not perform any kind of data type validation (related to the first exception). The limitation of SHACLEARNER is it does not mine SPARQL-like constraints (*SHACL-SPARQL*³).

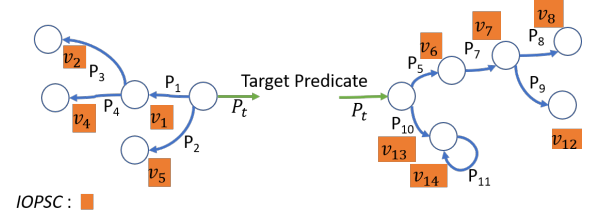


Fig. 3. Trees for a target predicate P_t . Each v_i indicates a TreeSC.

3.7.1. Tree Shapes are Useful for Human Interaction

Shapes offer KG documentation as readable patterns and also provide a way to validate a KG. Our novel tree shapes can additionally be used for KG-completion. While there are several methods proposed to complete KGs automatically (e.g., [16, 21, 22]) by predicting missing facts (also called *links*), all these methods traverse the KG in a breadth-first manner. In other words, they sequentially answer a set of independent questions such as *birthPlace(Trump, ?)* followed by *playsIn(Naomi_Watts, ?)*. Our proposed tree shapes instead provide an opportunity to work sequentially along a path of dependent questions such as *birthPlace(Trump, ?₁)* followed by *capitalOf(?₁, ?)*. The latter question cannot even be asked until we have an answer for the former question, and the existence of an answer to the former gives us the confidence to proceed to the next question along the path. This completion strategy is *depth-first* as it works through a shape tree. Importantly, when we want to ask such completion questions of a human, this depth-first questioning strategy will reduce the cognitive load due to the contextual connections between successive questions. This strategy for human-KG-completion is applied in the smart KG editor *Schímatos* [24] using trees that can be generated by SHACLEARNER.

²<https://www.w3.org/TR/shacl/#dfn-shacl-core>

³<https://www.w3.org/TR/shacl/#dfn-shacl-sparql>

1 Tree shapes can also help a human expert to ex-
 2 tract a more intuitive concise sub-tree out of a deeper,
 3 more complex tree when desired for human inter-
 4 pretation. If a tree with confidence $TreeSC_{orig}$ is pruned
 5 either by removing branches (width-wise) or by reduc-
 6 ing the length of branches (depth-wise), it remains a
 7 valid tree with confidence $TreeSC_{new}$ with the prop-
 8 erty that $TreeSC_{new} \geq TreeSC_{orig}$. Hence, by pruning
 9 a tree we obtain a simpler tree with higher confidence
 10 in the KG.

11 4. Related Work

12
 13
 14
 15 There are some exploratory attempts to address
 16 learning SHACL shapes from KGs [9, 25–28]. They
 17 are procedural methods without logical foundations
 18 and are not shown to be scalable to handle real-world
 19 KGs. They work with a small amount of data and the
 20 representation formalism they use for their output is
 21 difficult to compare with the IOP rules which we use
 22 in this paper. [28] carries out the task in a semi-automatic
 23 manner: it provides a sample of data to an off-the-shelf
 24 graph structure learner and provides the output in an
 25 interactive interface for a human user to create SHACL
 26 shapes.

27 In [8] the authors provide an interactive framework
 28 to define SHACL shapes with different complexities,
 29 including nested patterns that are similar to the trees
 30 that we use. A formalisation of the approach is given,
 31 but there are no shape quality measures that are es-
 32 sential for large scale shape mining. Because the pa-
 33 per does not provide a system that discovers patterns
 34 from massive KGs, we can not deploy their method for
 35 comparison purposes.

36 There are some works [29, 30] that use existing on-
 37 tologies for KGs to generate SHACL shapes. [29] uses
 38 two different kinds of knowledge to automatically gen-
 39 erate SHACL shapes: ontology constraint patterns as
 40 well as input ontologies. In our work we use the KG it-
 41 self to discover the shapes, without relying on external
 42 modelling artefacts.

43 From an application point of view, there are papers
 44 which investigate the application of SHACL shapes to
 45 the validation of RDF databases including [31, 32], but
 46 these do not contribute to the discovery of shapes.

47 [33] proposes an extended validation framework for
 48 the interaction between inference rules and SHACL
 49 shapes in KGs. When a set of rules and shapes are pro-
 50 vided, a method is proposed to detect which shapes
 51 could be violated by applying a rule.

1 There are some attempts to provide logical founda-
 2 tions for the semantics of the SHACL language in-
 3 cluding [34] that presents the semantics of recursive
 4 SHACL shapes. By contrast, in our work we approach
 5 SHACL semantics in the reverse direction. We start
 6 with logical formalisms with both well-defined seman-
 7 tics and motivating use cases to derive shapes that can
 8 be trivially expressed in a fragment of SHACL.

9 5. Experiments

10
 11
 12
 13 We have implemented our SHACLEARNER⁴ based
 14 on Algorithm 1, and conducted experiments to assess
 15 it. Our experiments are designed to prove the effective-
 16 ness of our SHACLEARNER at capturing shapes with
 17 varying confidence, length and cardinality from vari-
 18 ous real-world massive KGs. Since our proposed sys-
 19 tem is the first method to learn shapes from massive
 20 KGs automatically, we have no benchmark with which
 21 to compare. However, the performance of our system
 22 shows that it can handle the task satisfactorily and can
 23 be applied in practice. We demonstrate that:

- 24 1. SHACLEARNER is scalable so it can handle
 25 real-world massive KGs including DBpedia with
 26 over 11M facts.
- 27 2. SHACLEARNER can learn several shapes each
 28 for various target predicates.
- 29 3. SHACLEARNER can discover diverse shapes
 30 with respect to the quality measurements of
 31 IOPSC and IOPHC.
- 32 4. SHACLEARNER discovers shapes of varying
 33 complexity, and diverse with respect to length
 34 and cardinality.
- 35 5. SHACLEARNER discovers *every* high-quality
 36 rule (with $IOPSC \geq 0.9$) for a small complete
 37 KG, by comparison with an ideal learner.
- 38 6. SHACLEARNER discovers more complex shapes
 39 (trees) by aggregating learnt IOP rules efficiently.
 40

41 Our four benchmark KGs are described in Table 1.
 42 Three benchmarks, namely YAGO2s, Wikidata, and
 43 DBpedia, are common KGs and have been used in rule
 44 learning experiments previously [11, 16]. The fourth
 45 is a small synthetic KG, Poker, for analysing the com-
 46 pleteness of our algorithm. The Poker KG was intro-
 47 duced in [14], adapted from the classic version [35, 36]

48
 49
 50 ⁴Detailed experimental results can be found at
 51 <https://www.dropbox.com/sh/dn1kujeg0b6o9bw/AAAbKG9KfZ7e0eaNONz2zE93a?dl=0>

to be a rich and correct KG for evaluation experiments. Each poker hand comprises 5 playing cards drawn from a deck with 13 ranks and 4 suits. Each card is described using two attributes: suit and rank. Each hand is assigned to any or all of 9 different ranks, including High Card, One Pair, Two Pair, etc. We randomly generate 500 poker hands and all facts related to them to build a small but complete and correct KG as summarised in Table 1. 28 out of 35 predicates are unary predicates, such as `fullHouse(x)` where x is a specific poker hand.

Table 1
Benchmark specifications

KG	# Facts	# Entities	# Predicates
YAGO2s	4,122,438	2,260,672	37
Wikidata	8,397,936	3,085,248	430
Wikidata +UP	8,780,716	3,085,248	587
DBpedia 3.8	11,024,066	3,102,999	650
DBpedia 3.8 +UP	11,498,575	3,102,999	1,005
Poker	6,790	575	35

All experiments were conducted on an Intel Xeon CPU E5-2650 v4 @2.20GHz, 66GB RAM and running CentOS 8.

5.1. Transforming KGs with type predicates for experiments

Since real-world KG models treat predicates and entities in a variety of ways, we require a common representation for this work that clearly distinguishes entities and predicates. We employ an abstract model that is used in Description Logic ontologies, where classes and types are named unary predicates, and roles (also called properties) are named binary predicates.

Presenting the class and type information as unary predicates also allows us to learn fully abstracted (entity-free) shapes (e.g. $type_E(x) \rightarrow P1(x, y), P2(y, z)$) instead of partially instantiated shapes (e.g. $type(x, \mathbf{E}) \rightarrow P1(x, y), P2(y, z)$). This feature is important since learning partially instantiated shapes can cause an explosion in the space of possible shapes. Using self-loop links for unary predicates is a syntactic sugar to keep the presentation in the triple format, the same as the original input KGs.

In real-world KGs, concept or class membership may be modeled as entity instances of a binary fact. For example, DBpedia contains "`<dbo:type>(x, y)`"

and "`<dbo:class>(x, y)`" predicates where the second arguments of these predicates are types (e.g. "`<db:Album>`" or "`<db:City>`") or classes (e.g. "`<db:Reptile>`" or "`<db:Bird>`"). Instead, we choose to model types and classes with *unary* predicates. To do so we make new predicates like `<dbo:type>_<db:Album>(x)` from facts in the form `<dbo:type>(x, <db:Album>)`, where x is the name of an album. Then we produce new unary facts based on the new predicate and related facts. For example, for `<dbo:type>(Revolver, <db:Album>)`, we produce the new fact `<dbo:type>_<db:Album>(Revolver)`.

We manually choose two predicates from DBpedia 3.8 "`<dbo:type>`" and "`<dbo:class>`" and one from Wikidata "`<occupation_P106>`" to generate our unary predicates and facts. These predicates each have a class as their second argument. To prune the classes with few instances for which learning may be pointless, we consider only our unary predicates which have at least 100 facts. We retain the original predicates and facts in the KG as well as extending it with our new ones. In Table 1 we report the specifications of two benchmarks where we have added the unary predicates and facts (denoted as +UP).

5.2. Learning IOP Rules

We follow the established approach for evaluating KG rule-learning methods, that is, measuring the quantity and quality of distinct rules learnt. Rule quality is measured by Inverse open path standard confidence (IOPSC) and Inverse open path head coverage (IOPHC). We randomly select 50 target predicates from Wikidata and DBpedia unary predicates (157 and 355 respectively). We use all binary predicates of YAGO2s (i.e., 37) as target predicates. Each binary target predicate serves as two target predicates, once in the straight form (where the object argument of the predicate is the common variable to connect the head) and secondly in its reverse form (where the subject argument serves to connect). In this manner, we ensure that the results of SHACLEARNER on YAGO2s with its binary predicates as targets is comparable with the results for Wikidata and DBpedia that have unary predicates as targets. Hence for YAGO2s we have 74 target predicates.

A 10 hour limit is set for learning each target predicate.

Table 2 shows #Rules, the average numbers of quality IOP rules found; %SucTarget, the proportion of target predicates for which at least one IOP rule was found; and the running times in hours, averaged over

the targets. Only high quality rules meeting minimum quality thresholds are included in these figures, that is, with $IOPSC \geq 0.1$ and $IOPHC \geq 0.01$, thresholds established in comparative work [11]. These thresholds are pretty low, so rules with low qualities are also covered by choosing such low thresholds. Generally, selecting a low threshold at the time of learning is a safe choice since the user can decide later on pruning the rules by applying any qualities thresholds.

Table 2
Performance of SHACLEARNER on benchmark KGs

Benchmark	%SucTarget	#Rules	Time (hours)
YAGO2s	80	42	0.6
Wikidata+UP	82	67	1.8
DBpedia+UP	98	157	2.4

SHACLEARNER shows satisfactory performance in terms of both the runtime and the numbers of quality rules mined. Note that rules found have a variety of lengths and cardinalities. To better present the quality performance of rules we illustrate the distribution of rules with respect to the features, IOPSC, IOPHC, cardinality and length, and also IOPSC vs length. In the following, the *proportion* of mined rules having the various feature values is presented, to more evenly demonstrate the quality of performance over the three very different KGs.

The distribution of mined rules with respect to their IOPSC and IOPHC is shown in Figure 4. In the left-hand chart we observe a consistent decrease in the proportion of quality rules as the IOPSC increases. In the right hand chart we see a similar pattern for increasing IOPHC, but the decrease is not as consistent, showing there must be many relevant rules with many covered head instances. Over all benchmarks, the majority of learned rules have lower IOPSC and IOPHC. This is expected because the statistical likelihood of poor quality rules is much higher. Indeed, we show experimentally in section 5.3 that our pruning techniques, that are necessary for scalability, prune away predominantly lower quality rules.

With respect to IOPSC, proportionally more quality rules are learnt from DBpedia than from the other KGs, with Wikidata second, ahead of YAGO2s. This phenomenon might be caused by the way that we select the target predicates: for Wikidata and DBpedia we handpick unary predicates that represent classes, while for YAGO2 we consider all binary predicates.

The distribution of mined rules with respect to their cardinalities is shown in Figure 5. We observe that the largest proportion of rules has a cardinality of 1, as expected, as they have the least stringent requirements to be met in the KG. We observe an expected decrease with greater cardinalities as they demand tighter restrictions to be satisfied. YAGO2 demonstrates a tendency towards higher cardinalities than the other KGs, possibly a result of its more curated development.

The distribution of mined rules with respect to their lengths is shown in Figure 6. One might expect that as the length increases, the number of rules would increase since the space of possible rules grows, and this is what we see.

For a concrete example of SHACL learning, we show the following three IOP rules mined from DBpedia in the experiments. The IOPSC, IOPHC, and Cardinality annotations respectively prefix each rule.

0.64, 0.01, 1 : < dbo : type > _ < db : Song > (x) →
< dbo : album > (x, z₁) ∧ < dbo : recordLabel > (z₁, y).

0.63, 0.01, 1 : < dbo : type > _ < db : Song > (x) →
< dbo : producer > (x, y).

0.41, 0.01, 2 : < dbo : type > _ < db : Song > (x) →
< dbo : producer > (x, y).

< dbo : type > _ < db : Song > (x) expresses x is a song. The first rule indicates x should belong to an album (z_1) that has y as record label. The second rule requires a song (x) to have at least one producer while the third rule requires a song to have at least two producers, and these two rules are distinguished by the cardinality annotation. As we discussed in 3.1, the third rule is more constraining than the second, so the confidence of the third rule is lower than the confidence of the second, based on the KG data.

The rules can be trivially rewritten (through a script that we developed) into SHACL syntax as follows.

```
:s1 a sh:NodeShape ;
    sh:targetClass :<dbo:type>_<db:Song> ;
    sh:property [
        sh:minCount 1 ;
        sh:path :<dbo:album> ;
        [sh:path :<dbo:recordLabel> ; ]
    ] .
:s2 a sh:NodeShape ;
    sh:targetClass :<dbo:type>_<db:Song> ;
```

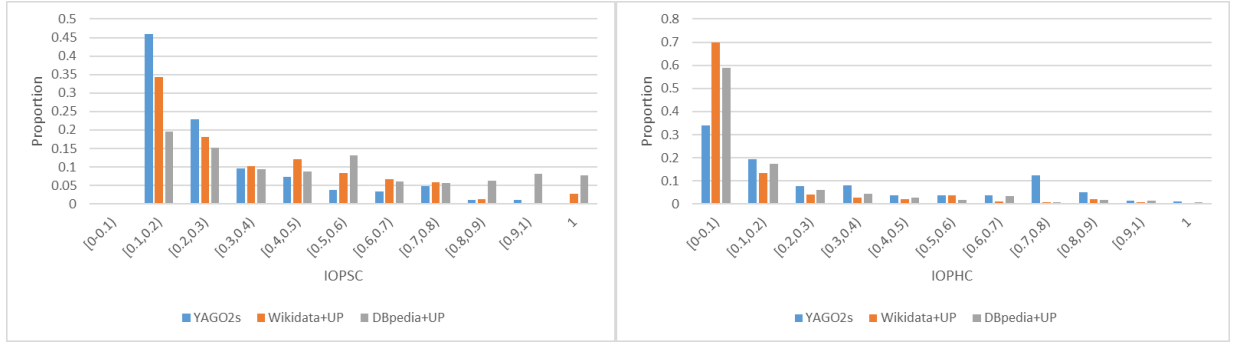


Fig. 4. Distribution of mined rules with respect to IOPSC and IOPHC

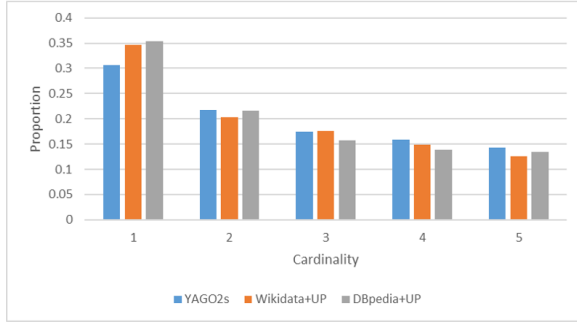


Fig. 5. Distribution of mined rules with respect to their cardinalities

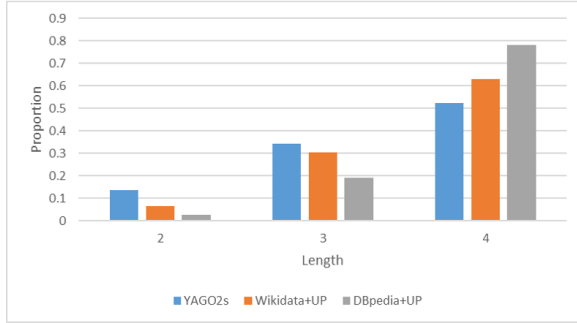


Fig. 6. Distribution of mined rules with respect to their lengths

```

43 sh:property [
44   sh:minCount 1 ;
45   sh:path :<dbo:producer> ;
46 ] .
47 :s3 a sh:NodeShape ;
48 sh:targetClass :<dbo:type>_<db:Song> ;
49 sh:property [
50   sh:minCount 2 ;
51   sh:path :<dbo:producer> ;
52 ] .

```

5.2.1. IOPSC vs IOPHC

Using rules found in the experiments, we further illustrate the practical meaning of the IOPSC and IOHC qualities. While IOPSC determines the confidence of a rule based on counting the proportion of target predicate instances for which the rule holds true in the KG, IOPHC indicates the proportion of rule consequent instances that are justified by target predicate instances in the KG, thereby indicating the relevance of the rule to the target. In Wikidata+UP, all unary predicates are occupations such as singer or entrepreneur, so all the entities which have these types turn out to be persons even though there is no explicit person type in our KG. Thus, the occupations all have very similar IOP rules about each of them with high IOPSC and low IOPHC, like the following one, for example.

$$0.47, 0.01, 1 : \langle \text{occupation} \rangle _ \langle \text{singer} \rangle (x) \rightarrow \langle \text{country_of_citizenship} \rangle (x, y).$$

On the other hand, for these unary occupation predicates there are also some IOP rules with high IOPHC that apply only to one specific unary predicate, such as the following one.

$$0.15, 0.16, 5 : \langle \text{occupation} \rangle _ \langle \text{singer} \rangle (x) \rightarrow \langle \text{performer_musical_artist} \rangle (x, y).$$

5.3. Completeness Analysis of IOP Rule Learning

SHACLEARNER uses two tricks to significantly reduce the search space for IOP rules in PathFinding() of Algorithm 1, namely the prior Sampling() and the heuristic pruning used inside PathFinding() that uses

the embedding-based scoring function. We conduct experiments to explore how these two pruning methods affect SHACLEARNER with regard to the quality and quantity of learnt rules. To do so we create three variants of SHACLEARNER as follows.

(-S+H): SHACLEARNER that does not sample and so uses the complete input KG in all components, including embedding learning, heuristic pruning and ultimate evaluation.

(+S-H): SHACLEARNER that samples but does not use heuristic pruning and so generates rules based on the sampled KG and evaluates rules on the complete KG.

(-S-H): SHACLEARNER that does not use sampling nor heuristic pruning. This system is an ideal IOP rule learner that generates and evaluates all possible rules up to the maximum length parameter. Since this method is a brute-force search, it cannot handle any real-world KG such as those we used in the performance evaluation (section 5.2).

(+S+H): SHACLEARNER with its full functionality.

Since we use only the small Poker KG for this experiment, we can handle the task without sampling or heuristic mechanisms. We use all 28 unary predicates as the target predicates. The first row of table 3 shows

Table 3

IOPSC performance of the ideal no-pruning SHACLEARNER variant on Poker, showing the percentage reduction in rules found by IOPSC intervals for variants with sampling and/or heuristic pruning

Learner/IOPSC	[0.1,0.3)	[0.3,0.5)	[0.5,0.7)	[0.7,0.9)	[0.9,1]
#(-S-H)	163	44	74	36	46
(-S+H)%	-10%	0%	-32%	-25%	0%
(+S-H)%	-93%	-93%	-12%	-11%	0%
(+S+H)%	-98%	-93%	-35%	-22%	0%

the number of IOP rules that are learnt by ideal, modified SHACLEARNER (-S-H) with no pruning, for all target predicates, separated into various IOPSC intervals. The latter rows show, for each variant, the percentage reduction in the number of rules found, relative to the first row. The last row corresponds to unmodified SHACLEARNER of Algorithm 1. For example, consider the first column (the number of learnt rules with IOPSC in range of [0.1,0.3)), in the first row (#(-S-H)) we have the number of rules learnt by ideal rule learner (SHACLEARNER (-S-H)) that is inefficient yet complete which is 163. In the second row,

we have the performance of the SHACLEARNER (-S+H) (the system without sampling and with heuristic rule learning module) in comparison with the ideal rule learner (SHACLEARNER (-S-H)) which is -10% that means, SHACLEARNER (-S+H) learns 146 rules (with IOPSC in range of [0.1,0.3)) which is 17 or 10% less than the performance of ideal learner, so its performance is -10%.

We observe that SHACLEARNER does not miss any rules of the highest quality, i.e., with $IOPSC \geq 0.9$. SHACLEARNER's pruning methods cause it to fail to discover more rules of lower quality, with the number of missing rules increasing as quality decreases. This is a reassuring property, since the goal of pruning is to improve the computational performance without missing high-quality rules. In real applications we will typically retain and action only the highest quality rules.

We observe that, unlike the other pruning variants, using heuristic pruning alone in (-S+H) does not uniformly increase in effectiveness with decreasing rule quality. This may be because using the complete KG for learning rules about all target predicates could harm the quality of the learnt embeddings for using them in the scoring function of SHACLEARNER. The better quality of embeddings extracted from the sampled KG comes from our sampling method that creates a customized KG WRT the target predicate. All entities in sampled KG are connected via target predicates (directly related to target predicate) or closed neighbours of directly related entities as we described in 3.3.

5.4. Learning Trees from IOP Rules

Now we turn to presenting results for the trees that are built based on the IOP rules discovered in the experiments. We report the characteristics of discovered trees in Table 4. We use a value of 0.1 for the $TreeSCMIN$ parameter and show average TreeSC for each KG, along with the average number of branches in the trees and the average runtime building each tree. The number of trees for each KG is defined by the number of target predicates for which we have at least one IOP rule (see %SucTarget in Table 2).

Table 4

Tree-learning performance of SHACLEARNER on benchmark KGs

Benchmark	TreeSC	#branches	Time (hours)
YAGO2s	0.13	21	0.06
Wikidata+UP	0.19	21	0.1
DBpedia+UP	0.2	88	0.14

The results show the running time for aggregating IOP rules into trees is lower than the initial IOP mining time by a factor greater than 10. If, on the other hand, we wanted to discover such complex shapes from scratch it would be exhaustively time consuming due to the sensitivity of rule learners to the maximum length of rules. The number of potential rules in the search space grows exponentially with regards to the maximum number of predicates of the rules. The average number of branches in the mined trees are 50%, 31%, and 56% of the corresponding number of mined rules respectively from Table 2. Hence, by imposing the additional tree-shaped constraint over the basic IOP-shaped constraint, at least 44% of IOP rules are pruned.

For an example of tree shape learning, in the following, we show a fragment of a 39-branched tree mined from DBpedia by aggregating IOP rules in the experiments.

$$\begin{aligned}
 0.13 : & \langle \text{dbo} : \text{type} \rangle _ \langle \text{db} : \text{Song} \rangle (x) \rightarrow \\
 & 1 : \langle \text{dbo} : \text{album} \rangle (x, z_1) \wedge \\
 & \langle \text{dbo} : \text{recordLabel} \rangle (z_1, y_1) \wedge \\
 & 2 : \langle \text{dbo} : \text{album} \rangle (x, z_2) \wedge \\
 & \langle \text{dbo} : \text{producer} \rangle (z_2, y_2) \wedge \\
 & 1 : \langle \text{dbo} : \text{album} \rangle (x, z_3) \wedge \\
 & \langle \text{dbo} : \text{genre} \rangle (z_3, y_3) \wedge \\
 & 3 : \langle \text{dbo} : \text{artist} \rangle (x, z_4) \wedge \\
 & \langle \text{dbo} : \text{musicalArtist} \rangle (z_4, y_4).
 \end{aligned}$$

Here, the first annotation value (0.13) presents the SC of the tree and the subsequent values at the beginning of each branch indicate the branch cardinality. This tree can be read as saying that a song has an album with a record label, an album with two producers, an album with a genre, and an artist who is a musical artist.

As can be seen here, there remains an opportunity for improving tree shapes for simplicity and easier interpretation by unifying some variables that occur in predicates that occur in multiple branches. We plan to investigate this potential post-processing step in future work.

6. Conclusion

In this paper we propose a method to learn SHACL shapes from KGs as a way to describe KG patterns, to validate KGs, and also to support new data entry. For entities that satisfy target predicates, our shapes describe conjunctive paths of constraints over properties, enhanced with minimum cardinality constraints.

We reduce the SHACL learning problem to learning a novel kind of rules, Inverse Open Path rules (IOP). We introduce rule quality measures IOP Standard Confidence, IOP Head Coverage and Cardinality which augment the rules. IOPSC effectively extends SHACL with shapes, representing the quantified uncertainty of a candidate shape to be selected for interestingness or for KG verification. We also propose a method to aggregate learnt IOP rules in order to discover more complex shapes, trees.

The shapes support efficient and interpretable human validation in a depth-first manner and are employed, for example, in an editor *Schímatos* [24] for manual knowledge graph completion. The shapes can also be used to complete information triggered by entities with only a type or class declaration by automatically generating dynamic data entry forms. In this manual mode, they can also be used more traditionally to complete missing facts for a target predicate, as well as other predicates related to the target, while enabling the acquisition of facts about entities that are entirely missing from the KG.

To learn such rules we adapt an embedding-based Open Path Rule Learner (OPRL) by introducing the following novel components: (1) we propose an IOP rule language which allow us to mine rules with open variables with one predicate forming the body and a chain of predicates as the head, while keeping the complexity of the learning phase manageable; (2) we introduce cardinality constraints and tree shapes for more expressive patterns; and (3) we propose an efficient method to evaluate IOP rules and trees by exactly computing the quality measures of each rule using fast matrix and vector operations.

Our experiments show that SHACLEARNER can mine IOP rules of various lengths, cardinalities, and qualities from three massive real-world benchmark KGs including Yago, Wikidata and DBpedia.

Learning shape constraints from schema-free knowledge-bases, such as most modern KGs, is a challenging task. The challenge starts from the formalism of constraints that determine the scope of knowledge that could be acquired. The next challenge, which is depen-

1 dent on the first choice, is to design an efficient learning
 2 method. Dealing with uncertainty in the constraints
 3 and the learning process adds an extra dimension of
 4 challenge, but also adds to utility. A good learning algorithm
 5 should scale gracefully so that discovered constraints
 6 are relatively more certain than those that are
 7 missed. SHACLEARNER establishes a benchmark for
 8 this problem.

9 In future work we will validate the shapes we learn
 10 with SHACLEARNER via formal human-expert evaluation
 11 and further extend the expressivity of the shapes
 12 we can discover. We also propose to redesign the
 13 SHACLEARNER algorithm for a MapReduce implementation
 14 to handle extremely massive KGs with tens
 15 of billions of facts, such as the most recent version of
 16 Wikidata.

17 Acknowledgments

18 The authors acknowledge the support of the Australian
 19 Government Department of Finance and the Australian
 20 National University for this work.

21 References

- 22 [1] S. Auer, C. Bizer, G. Kobilarov, J. Lehmann, R. Cyganiak
 23 and Z.G. Ives, DBpedia: A Nucleus for a Web of Open Data,
 24 in: *ISWC*, Lecture Notes in Computer Science, Vol. 4825,
 25 Springer, 2007, pp. 722–735.
- 26 [2] F.M. Suchanek, G. Kasneci and G. Weikum, Yago: a core of
 27 semantic knowledge, in: *WWW*, ACM, 2007, pp. 697–706.
- 28 [3] N. Noy, Y. Gao, A. Jain, A. Narayanan, A. Patterson and J. Taylor,
 29 Industry-Scale Knowledge Graphs: Lessons and Challenges,
 30 *Commun. ACM* **62**(8) (2019), 36–43–.
- 31 [4] E.S. Callahan and S.C. Herring, Cultural bias in Wikipedia
 32 content on famous persons, *Journal of the American Society
 33 for Information Science and Technology* **62**(10) (2011), 1899–
 34 1915.
- 35 [5] D. Vrandečić and M. Krötzsch, Wikidata: a free collaborative
 36 knowledgebase, *Commun. ACM* **57**(10) (2014), 78–85.
- 37 [6] H. Knublauch and D. Kontokostas, Shapes Constraint Language
 38 (SHACL), 2017. <https://www.w3.org/TR/shacl/>.
- 39 [7] Z. Huo, K. Taylor, X. Zhang, S. Wang and C. Pang, Generating
 40 multidimensional schemata from relational aggregation
 41 queries, *World Wide Web* **23** (2020).
- 42 [8] I. Boneva, J. Dusart, D. Fernández-Álvarez and
 43 J.E.L. Gayo, Semi Automatic Construction of ShEx
 44 and SHACL Schemas, *CoRR* **abs/1907.10603** (2019).
 45 <http://arxiv.org/abs/1907.10603>.
- 46 [9] P. Ghiasnezhad Omran, K. Taylor, S. Rodríguez Méndez and
 47 A. Haller, Towards SHACL Learning from Knowledge Graphs,
 48 in: *{ISWC2020} Posters & Demonstrations*, CEUR Proceedings,
 49 2020.
- 50 [10] H. Knublauch, D. Allemang and S. Steyskal, SHACL Advanced
 51 Features, 2017. <https://www.w3.org/TR/shacl-af/>.
- [11] L. Galárraga, C. Teflioudi, K. Hose and F.M. Suchanek, Fast
 rule mining in ontological knowledge bases with AMIE+, *The
 International Journal on Very Large Data Bases* (2015), 707–
 730. ISBN 1066-8888.
- [12] L. Bellomarini, E. Sallinger and G. Gottlob, The Vadalog system:
 Datalogbased reasoning for knowledge graphs, in: *VLDB*,
 Vol. 11, 2018, pp. 975–987. ISSN 21508097.
- [13] W. Fan, C. Hu, X. Liu and P. Lu, Discovering graph functional
 dependencies, in: *SIGMOD*, ACM, 2018, pp. 427–439. ISSN
 07308078. ISBN 9781450317436.
- [14] P. Ghiasnezhad Omran, K. Taylor, S. Rodríguez Méndez and
 A. Haller, Active Knowledge Graph Completion, Technical
 Report, ANU Research Publication, 2020.
- [15] P. Ghiasnezhad Omran, K. Taylor, S. Rodríguez Mendez and
 A. Haller, Active Knowledge Graph Completion, in: *ISWC
 Satellite Tracks (Posters and Demonstrations, Industry, and
 Outrageous Ideas)*, 2020, pp. 89–94.
- [16] P.G. Omran, K. Wang and Z. Wang, Scalable Rule Learning
 via Learning Representation, in: *IJCAI*, 2018, pp. 2149–2155.
 ISBN 9780999241127.
- [17] Y. Chen, D.Z. Wang and S. Goldberg, ScaLeKB: scalable
 learning and inference over large knowledge bases, *The International
 Journal on Very Large Data Bases* (2016), 893–918.
- [18] R. Agrawal and R. Srikant, Fast algorithms for mining association
 rules, in: *VLDB*, Vol. 1215, 1994, pp. 487–499.
- [19] K. Taylor, Generalization by absorption of definite clauses, *The
 Journal of Logic Programming* **40**(2–3) (1999), 127–157.
- [20] S. Staworko, I. Boneva, J.E. Labra Gayo, S. Hym,
 E.G. Prud’hommeaux and H. Solbrig, Complexity and Expressiveness
 of ShEx for RDF, in: *ICDT*, 2015, pp. 195–211.
<http://linkeddata.org/>.
- [21] M. Nickel, L. Rosasco and T. Poggio, Holographic Embeddings
 of Knowledge Graphs, in: *AAAI*, 2016, pp. 1955–1961.
- [22] M. Nickel, V. Tresp and H.-P. Kriegel, A three-way model for
 collective learning on multi-relational data, in: *ICML*, 2011,
 pp. 809–816.
- [23] P.G. Omran, K. Wang and Z. Wang, An Embedding-based
 Approach to Rule Learning in Knowledge Graphs, *TKDE* **33**
 (2021), 1348–1359. doi:10.1109/TKDE.2019.2941685. <https://ieeexplore.ieee.org/document/8839576/>.
- [24] J. Wright, S. Rodríguez Méndez, A. Haller, K. Taylor and
 P. Omran, Schimatos: a SHACL-based Web-Form Generator for
 Knowledge Graph Editing, in: *ISWC*, 2020.
- [25] N. Mihindukulasooriya, M. Rifat, A. Rashid, G. Rizzo,
 R. García-Castro, O. Corcho and M. Torchiano, RDF Shape
 Induction using Knowledge Base Profiling, in: *Annual {ACM}
 Symposium on Applied Computing, {SAC}*, Vol. 8, 2018,
 p. pages. ISBN 9781450351911.
- [26] D. Fernández-Álvarez, H. García-González, J. Frey, S. Hellmann
 and J.E.L. Gayo, Inference of latent shape expressions associated
 to DBpedia ontology, in: *ISWC Posters*, Vol. 2180,
 2018. ISSN 16130073.
- [27] B. Spahiu, A. Maurino and M. Palmonari, Towards improving
 the quality of knowledge graphs with data-driven ontology
 patterns and SHACL, in: *Workshop on Ontology Design and
 Patterns*, Vol. 2195, 2018, pp. 52–66. ISSN 16130073.
- [28] I. Boneva, J. Dusart, D.F. Álvarez and J.E. Labra Gayo, Shape
 designer for ShEx and SHACL constraints, in: *ISWC Posters*,
 Vol. 2456, 2019, pp. 269–272. ISSN 16130073.

- [29] A. Cimmino, A. Fernández-Izquierdo and R. García-Castro, Astrea: Automatic Generation of SHACL Shapes from Ontologies, in: *ESWC*, Vol. 12123 LNCS, 2020, pp. 497–513. ISSN 16113349. ISBN 9783030494605.
- [30] H.J. Pandit, D. O’Sullivan and D. Lewis, Using ontology design patterns to define SHACL shapes, in: *CEUR Workshop Proceedings*, Vol. 2195, 2018, pp. 67–71. ISSN 16130073.
- [31] J.-E. Labra-Gayo, E. Prud’hommeaux, H. Solbrig and I. Boneva, Validating and describing linked data portals using shapes, *CoRR* (2017).
- [32] I. Boneva, J.E. Labra Gayo and E.G. Prud’hommeaux, Semantics and validation of shapes schemas for RDF, in: *ISWC*, Vol. 10587 LNCS, 2017, pp. 104–120. ISSN 16113349. ISBN 9783319682877. <http://www.w3.org/Submission/spin-modeling/>.
- [33] P. Pareti, G. Konstantinidis, T.J. Norman and S. Murat, SHACL Constraints with Inference Rules, in: *ISWC*, 2019.
- [34] J. Corman, J.L. Reutter and O. Savković, Semantics and validation of recursive SHACL, in: *ISWC*, Vol. 11136 LNCS, 2018, pp. 318–336. ISSN 16113349. ISBN 9783030006709.
- [35] R. Cattral, F. Oppacher and D. Deugo, Evolutionary Data Mining with Automatic Rule Generalization., *Recent Advances in Computers, Computing and Communications* (2002), 296–300.
- [36] D. Dua and C. Graff, UCI Machine Learning Repository, 2017, <archive.ics.uci.edu/ml> Retrieved Nov 2019.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
511
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51

Author responses to the comments

Article title: Learning SHACL Shapes from Knowledge Graphs

Tracking Number: 2906-4120

Authors: Pouya Ghasnezhad Omran, Kerry Taylor, Sergio Rodriguez Mendez, and Armin Haller

We appreciate the reviewers' comments which have assisted us to improve the manuscript. Please find below our detailed response to the comments made by the reviewers. Adjusted parts are highlighted.

Editor comments

	comments	responses
0.1	<p>Though the authors clarified most of the issues in this version of the paper, some clarifications would further improve the paper:</p> <ol style="list-style-type: none">1.Lack of discussion about the SHACL features covered theoretically2.Lack of justification for the binary predicate used with types3.Lack of a qualitative analysis on discovered trees which can give readers more insights about the challenges and the potential future works4.Issue with table 3 to be read5.Further discussion on the choice of the sample and the efficiency of the experiments.	<p>We provided the explanations and revisions for these points. We have addressed these issues in our paper as follows:</p> <ol style="list-style-type: none">1. Refer to comment 3.12. Refer to comment 3.23. Refer to comment 3.44. Refer to comment 2.15. Refer to comments 2.2 and 2.3

Reviewer 1 comments

	comments	responses
1.1	<p>The authors clarified my doubts on the adoption of RESCAL. For future work, I suggest further investigating the impact of KG embeddings in learning SHACL shapes.</p>	<p>Agree. We plan to perform such an investigation in our future work. We included a paragraph in the end of the conclusion section where we state this aim.</p>

Reviewer 2 comments

	comments	responses
2.1	<p>First of all, regarding the experiments in section 5.3, table 3 is difficult to read.</p> <p>What is the probability that a rule satisfied in the KG is indeed discovered?</p> <p>Did you perform any analysis of what kind of rules does each of the configurations fail? Why?</p>	<p>We added the following explanation to make the content clear. (Page:15 column:1 line: 46)</p> <p>For example, consider the first column (the number of learnt rules with IOPSC in range of [0.1,0.3]), in the first row ($\#(-S-H)$) we have the number of rules learnt by ideal rule learner ($\backslash\text{SHACLearner}(-S-H)$) that is inefficient yet complete which is 163. In the second row, we have the performance of the $\backslash\text{SHACLearner}(-S+H)$ (the system without sampling and with heuristic rule learning module) in comparison with the ideal rule learner ($\backslash\text{SHACLearner}(-S-H)$) which is -10% that means, $\backslash\text{SHACLearner}(-S+H)$ learns 146 rules (with IOPSC in range of [0.1,0.3]) which is 17 or 10% less than the performance of ideal learner, so its performance is -10%.</p> <p>Based on our experiments, SHACLearner discovers all rules with high IOPSC ($\text{IOPSC} \geq 0.9$). As explained in the two last paragraphs of subsection 5.3. (Page:15 column:2 line:9)</p> <p>We observe that SHACLearner does not miss any rules of the highest quality, i.e., with $\text{IOPSC} \geq 0.9$. SHACLearner's pruning methods cause it to fail to discover more rules of lower quality, with the number of missing rules increasing as quality decreases. This is a reassuring property, since the goal of pruning is to improve the computational performance without missing high-quality rules. In real applications we will typically retain and action only the highest quality rules.</p> <p>We observe that, unlike the other pruning variants, using heuristic pruning alone in $(-S+H)$ does not uniformly increase in effectiveness with decreasing rule quality. This may be because using the complete KG for learning rules about all target predicates could harm the quality of the learnt embeddings for using them in the scoring function of SHACLearner. The better quality of embeddings extracted from sampled KG comes from our sampling method that creates a customized KG WRT the target predicate. All entities in sampled KG are connected via target predicates (directly related to target predicate) or closed neighbours of directly related entities as we described in 3.3.</p>

2.2	<p>Second, the authors claim that using the complete KG for learning rules about all target predicates could harm the quality of the learned embeddings. So, in case we want to apply SHACLEARNER to DBpedia, how should the sampling work?</p>	<p>To explain why using a sampled KG might improve the quality of embedding for learning shapes we revised the following part of the last paragraph of 5.3: (Page:15 column:2 line:18) We observe that, unlike the other pruning variants, using heuristic pruning alone in (-S+H) does not uniformly increase in effectiveness with decreasing rule quality. This may be because using the complete KG for learning rules about all target predicates could harm the quality of the learnt embeddings for using them in the scoring function of SHACLearner. The better quality of embeddings extracted from sampled KG comes from our sampling method that creates a customized KG WRT the target predicate. All entities in sampled KG are connected via target predicates (directly related to target predicate) or closed neighbours of directly related entities as we described in 3.3.</p> <p>We used our system on massive KGs like DBpedia as we reported in section 5.2. In DBpedia we have 650 original binary predicates and we added 355 unary predicates that we generated from type and class binary predicates. While we used 50 randomly selected generated unary predicates in DBpedia as target predicates, the system can be fed with any of these 1005 predicates as target predicate and it mines correspondent IOP rules and tree. We can put the minimum IOPSC and IOPHC to make sure we mine the rules that have enough support and generality in the KG.</p>
2.3	<p>Third, do you have any evidence of how much does the sampling take? For example, if we consider Dbpedia, and run SHACLearner how much will it take to prune all predicates and the relative entities?</p>	<p>We do not have the experimental results of running WRT different modules of SHACLearner (e.g., sampling, embedding learning, path finding, ...) to extract the running time of the sampling procedure. However, we added the following part about the theoretical complexity of the sampling method at the end of 3.3. (Page:7 column:1 line:51) As discussed in [23], consider E is the set of entities, and F is the set of facts, computing the set of sampled entities regarding a target predicate demands $2l \cdot F \cdot E$ in the worst case where l is the parameter of sampling which indicates the max length of the rules. Hence, the complexity of the sampling algorithm is square in terms of K , $K =$</p>

		(E,F). However, in the worst case, the sampled KG size is the same as the original KG, while in the real-world KGs, the number of entities that belong to E0 is restricted and the neighbors of E0 is restricted and so on. Hence, with the given fixed number of l the number of entities in sampled KG is far smaller than the original KG.
2.4	There are some typos still present in the revised version, some further proofreading will be welcome.	We addressed typos that we could find. We plan to use a professional proofreading service for the camera-ready version.

Reviewer 3 comments

	comments	responses
3.1	Regarding the response to 3.4, while I agree that there are various formalisms for shapes defined in the literature to express diverse patterns, as the title of this paper, "Learning SHACL Shapes from Knowledge Graphs" suggests it focuses on SHACL Shapes. Thus, I believe it's important to provide some insight on which SHACL features are covered, which can be covered theoretically but not covered due to time constraints, and which are fundamentally can't be covered by the approach out of the constraint types such as value type, cardinality, value range, string-based, property-pair etc defined in https://www.w3.org/TR/shacl/ .	We added the explanation at the last paragraph of 3.7 as follows. (Page:10 column:1 line:35) The uncertain shapes can be presented as standard SHACL shapes by applying the desired minimum confidence by the user. With the use of such a threshold, we can generate hard shapes. Hard shapes have complete qualities and are not augmented with any properties like IOPSC or ThreeSC. Aside from the cardinality, the tree may be straightforwardly interpreted as a set of SHACL shapes by reading off every path from the target predicate terminating at a node in the tree. The body predicate is declared a sh:nodeShape and the path of head predicates as nestedsh:path declarations within a sh:property declaration. Cardinality of a path is read from the annotation of the branch at the terminating node, and declared by sh:minCount within the property declaration. See section 3.1 for an example. SHACLEARNER supports the SHACL Core[6] features (node and property shapes). The exceptions are: (1) it treats all properties (object and datatype) as "plain predicates" (there is no distinction of their nature), (2) it covers all forms of non-empty paths (paths with a length greater or equal to one), (3) only min cardinality, and (4) does not perform any kind of datatype validation (related to the first exception). The limitation of SHACLearner is it does not mine SPARQL-like constraints (SHACL-SPARQL).
3.2	Regarding the response 3.6, while I understand the synaptic difference of encapsulating the type in an entity or in	In many KG modelling methods, including in rule learning methods like AMIE+[11], their language bias does not allow a predicate in the rule that has

	<p>a predicate, I still fail to see the benefit of that. Why not follow the RDF/OWL convention of using binary predicates for defining types? What problem that is being solved by unary predicates that was not possible to solve using the binary predicates (similar to how the authors have done in Yago2)? Isn't the notation $P(e,e)$ for unary predicates conflict with reflexive properties in the KG.</p> <p>Including unary predicates includes a lot additional details to the method and makes the SHACL Shapes have unconventional target classes such as <code>`sh:targetclass class:_;</code>. Thus, I believe still a bit clear description why binary predicates were not possible to use with types has to be included in the paper.</p>	<p>the same entity as its subject and object. "AMIE omits also reflexive rules, i.e., rules with atoms of the form $r(x, x)$, as they are typically of less interest in real-world KBs." [11] page 711 last part of first column.</p> <p>We followed the same principle, i.e. our shape do not include such atoms for binary predicates and we reserve the reflexive atom for presenting the unary predicates in our learnt shapes.</p> <p>To make it clear why using unary predicate is important we added the following part to section 5.1: (Page:12 column:1 line:39) Presenting the class and type information as unary predicates also able us to learn fully abstracted (entity-free) shapes (e.g., $\text{type_E}(x) \rightarrow P1(x,y), P2(y,z)$) instead of partially instantiated shapes (e.g., $\text{type}(x,E) \rightarrow P1(x,y), P2(y,z)$). This feature is important since learning partially instantiated shapes can cause an explosion in the space of possible shapes. Using self-loop links for unary predicates is a syntactic sugar to keep the presentation in the triple format, the same as the original input KGs.</p>
3.3	<p>In general, there are several points in the answer letter that are explained there but not explained in the paper such as 3.8, 3.14, 3.21. It might be beneficial to incorporate some of those explanations to the paper itself.</p>	<p>Fixed. We added the following parts:</p> <p>(Page:6 column:2 line:21) This Lemma is useful in the process of rule learning since it shows if we discard a rule with a length regarding its low IOPSC, v, we do not need to check the extension of this rule with more head atoms since those rules' IOPSCs are equal or less than v. Hence the process of pruning such rules does not harm the completeness of the rule learning.</p> <p>(Page:13 column:1 line:4) These thresholds are pretty low, so rules with low qualities are also covered by choosing such low thresholds. Generally, selecting a low threshold at the time of learning is a safe choice since the user can decide later on pruning the rules by applying any qualities thresholds</p>
3.4	<p>Regarding the response 3.25, even though I understand a fully-fledged human evaluation is out of the scope of</p>	<p>We provided the learnt trees and IOP rules in the external folder of the paper.</p>

	<p>this paper, a qualitative analysis by the authors of a smaller sample of discovered trees would still be beneficial for the reader to understand the usefulness of the discovered trees and also the limitations. Such analysis will really validate the results and give readers ideas about the current challenges and potential future work.</p>	<p>We plan to include the learnt shapes in the next release of Schimatos and evaluate those shapes WRT a human user evaluation.</p>
--	--	---