

Constraint Rules in Knowledge Graphs

Pablo Mollá Chárlez

February 9, 2025

Contents

1	Rule Evaluation Metrics	2
1.1	Support	2
1.2	Support of rule in sets	2
1.3	Confidence	2
1.4	Head Coverage	2
1.5	Weight of a Rule	3
2	SHACL: Shape-Based Constraint Language	3
2.1	SHACL Constraints	3
2.2	SHACK's Validation Steps	5

1 Rule Evaluation Metrics

1.1 Support

Support of a rule R indicates how many facts in the knowledge base K the rule correctly accounts for. Formally:

$$\text{support}(R) = |\{p \mid (K \wedge R \models p) \wedge p \in K\}|.$$

As rules become more specialized (adding more conditions in the premise), support often decreases because fewer facts match. In practice, a rule with very low support may be discarded (its applicability is minimal). Support addresses the **relevance of the rule** (how many facts it covers).

1.2 Support of rule in sets

The coverage of a rule R in a given set S (Support of the rule in set S) is defined as:

$$\text{support}(R, S) = |p : (K \wedge R \models p) \wedge p \in S|$$

1.3 Confidence

Confidence of a rule R measures the proportion of the rule's predictions that are actually true in K :

$$\text{confidence}(R) = \frac{\text{support}(R)}{\text{support}(R) + |\text{cex}(R)|},$$

where $\text{cex}(R)$ is the **set of counterexamples**-predictions the rule makes that are not found in K . Confidence represents the **accuracy of the rule** (fraction of correct predictions). Because knowledge graphs typically lack explicit negative facts, we need **assumptions to infer counterexamples**:

- **Closed World Assumption (CWA)**: Anything not stated is considered false. If the rule predicts a fact not in K , that prediction counts as a counterexample.

$$\text{cex}_{\text{CWA}} = |\{(x, y) \in B : (x, y) \notin H\}| \text{ where } R : B \longrightarrow H$$

- **Open World Assumption (OWA)**: Anything not stated is unknown. A missing fact is not necessarily a counterexample; it could still be true but not yet recorded.

$$\text{cex}_{\text{OWA}} = |B \cap \neg H| \text{ where } R : B \longrightarrow H$$

- **Partial Completeness Assumption (PCA)**: If any relationship r is known for a subject s , then all such r -relationships for s are expected to be recorded. This assumption tries to balance the extremes of CWA and OWA, treating some parts of the knowledge base as “closed” while leaving others “open.”

$$\text{cex}_{\text{PCA}} = |\{(x, y) \in B : x \in H \wedge (x, y) \notin H\}| \text{ where } R : B \longrightarrow H$$

1.4 Head Coverage

Head coverage for a rule $B \Rightarrow r(x, y)$ captures what fraction of the known facts for the predicate r (in the knowledge base K) is “explained” by that rule. Formally:

$$\text{hc}(B \Rightarrow r(x, y)) = \frac{\text{support}(B \Rightarrow r(x, y))}{|\{(x, y) \mid r(x, y) \in K\}|}.$$

Where $\text{support}(B \Rightarrow r(x, y))$ is the number of true pairs (x, y) for which B holds and $r(x, y)$ is actually in K , and $|\{(x, y) \mid r(x, y) \in K\}|$ is the total count of true instances of r in K .

Suppose we have a knowledge base about people and their living locations, and we are examining the rule:

$$\text{marriedTo}(x, y), \text{livedIn}(x, z) \Rightarrow \text{livedIn}(y, z).$$

In K , there are 10 known facts of the form `livedIn`(y, z). By checking the premises (`marriedTo`(x, y) and `livedIn`(x, z)), we find that the rule correctly predicts 4 of those 10 facts. In other words, the rule’s support is 4 (it matches 4 existing “livedIn” facts). So the **head coverage** is:

$$\text{hc} = \frac{\text{support}}{\text{total known facts for } r} = \frac{4}{10} = 0.4.$$

Thus, 40% of the `livedIn`(y, z) relationships in the KB are “covered” (i.e., explained) by that rule.

1.5 Weight of a Rule

The **weight of a rule** $w(R)$ is a composite metric that tries to balance how well R covers the true facts in the knowledge graph (akin to recall) and how many non-facts or potential negatives it predicts (related to precision). One way to define it is:

$$w(R) = \alpha \left(1 - \frac{\text{support}(R, \mathcal{G})}{|\mathcal{G}|} \right) + \beta \left(\frac{\text{support}(R, \mathcal{V})}{|\mathcal{V}|} \right),$$

where

- \mathcal{G} is the **set of grounded facts** (true statements) in the knowledge graph (finite).
- \mathcal{V} is a set representing **potential negative or unobserved facts** (which can be extremely large). We can also use $|U(R, \mathcal{V})|$, $U(R, \mathcal{V}) \subseteq \mathcal{V}$ is chosen so that $\text{support}(R, \mathcal{V}) \subseteq U(R, \mathcal{V})$, ensuring proper normalization.
- $\text{support}(R, \mathcal{G})$ counts **how many true facts R covers**, and $\text{support}(R, \mathcal{V})$ counts **how many “negative” or unobserved facts R covers**.
- α, β let you tune how much each term contributes.

In words, the first term $\left(1 - \frac{\text{support}(R, \mathcal{G})}{|\mathcal{G}|} \right)$ decreases when R correctly covers more real facts (increasing recall), whereas the second term $\frac{\text{support}(R, \mathcal{V})}{|\mathcal{V}|}$ increases if R also covers many negative or unobserved facts (worsening precision).

For instance, suppose that \mathcal{G} has 100 known true facts relevant to R . The rule covers 30 of them: $\text{support}(R, \mathcal{G}) = 30$. Moreover, \mathcal{V} has 300 potential negative facts; the rule covers 45 of those: $\text{support}(R, \mathcal{V}) = 45$. We choose $\alpha = 0.5$ and $\beta = 0.5$. We would get:

$$w(R) = 0.5 \left(1 - \frac{30}{100} \right) + 0.5 \left(\frac{45}{300} \right) = 0.5(0.70) + 0.5(0.15) = 0.35 + 0.075 = 0.425.$$

A lower weight in this scheme would typically mean the rule either fails to cover enough true facts or covers too many negatives-hence is less desirable.

2 SHACL: Shape-Based Constraint Language

SHACL (**Shapes Constraint Language**) is a **W3C Recommendation (2017)** that provides a standard schema language for RDF. It **allows one to define validation constraints on classes, nodes, properties, and literals** by expressing them directly in RDF. The validation process takes two inputs—a **data graph** (the RDF dataset to be validated) and a **shape graph** (the RDF-based definitions of the constraints)-and **produces a validation report**. A **shape graph can contain multiple shapes**, each described by RDF triples using the SHACL vocabulary (`sh:`). **Shapes come in two forms: Node Shapes**, which apply constraints directly to a node, and **Property Shapes**, which specify constraints on the nodes reached via a particular property path. **The goal of SHACL is thus to enforce and verify that RDF data conforms to the intended structure and semantics.**

2.1 SHACL Constraints

Table 1: SHACL Constraint Forms

Constraint	Shape Type	Description
Node Kind and Datatype Constraints		
<code>sh:class c</code>	Node / Property	Target nodes must be instances of class <code>c</code> . Multiple <code>sh:class</code> values are conjunctive.
<code>sh:datatype t</code>	Node / Property	Target nodes must be literals of datatype <code>t</code> ; at most one <code>sh:datatype</code> .
<code>sh:nodeKind k</code>	Node / Property	At most one <code>sh:nodeKind</code> declaration for the shape.
Numeric Constraints		
<code>sh:minInclusive / sh:minExclusive /</code> <code>sh:maxInclusive / sh:maxExclusive</code>	Node / Property	All target nodes must satisfy the numeric comparison (often combined with <code>sh:datatype xsd:integer</code> or other numeric types).
String Length Constraints		
<code>sh:minLength n; sh:maxLength m</code>	Node / Property	Target nodes must be literals or IRIs. One <code>sh:minLength</code> and one <code>sh:maxLength</code> allowed.
Regular Expression Constraints		
<code>sh:pattern regex; sh:flags flag</code>	Node / Property	Target nodes must match the given <code>regex</code> ; <code>flag</code> often includes <code>"i"</code> (case-insensitive).
Language Tag Constraints		
<code>sh:languageIn (lang_list)</code>	Node / Property	All target nodes must have a language tag in the given list.
<code>sh:uniqueLang true</code>	Property shape	Ensures at most one value per language tag among value nodes.
Cardinality Constraints		
<code>sh:maxCount n; sh:minCount m</code>	Property	On <code>sh:path p</code> , bounds the number of value nodes from <code>p</code> . <code>n</code> and <code>m</code> must be <code>xsd:integer</code> .
Fixed Value Constraints		
<code>sh:hasValue v</code>	Property	All values from the path must be equal to <code>v</code> .
<code>sh:in (v1 v2 ... vn)</code>	Property	All value nodes must be one of <code>v1</code> , <code>v2</code> , ..., <code>vn</code> .
Property Pair Constraints		
<code>sh:equals r</code>	Property	Values from <code>p</code> must match values from <code>r</code> exactly.
<code>sh:disjoint r</code>	Property	Values from <code>p</code> must not overlap values from <code>r</code> .
<code>sh:lessThan r</code>	Property	$\max(V_p) < \min(V_r)$.
<code>sh:lessThanOrEquals r</code>	Property	$\max(V_p) \leq \min(V_r)$.
Logical Operators		
<code>sh:not [shape]</code>	Node / Property	Negates the constraints defined in the nested shape.
<code>sh:and (shapes...)</code>	Node / Property	Target node must satisfy <u>all</u> listed shapes.
<code>sh:or (shapes...)</code>	Node / Property	Target node must satisfy <u>at least one</u> listed shape.
<code>sh:xone (shapes...)</code>	Node / Property	Target node must satisfy <u>exactly one</u> listed shape.

2.2 SHACK's Validation Steps

1. Identify Target Nodes. Determine which nodes the constraint applies to (e.g., via **sh:targetClass**, **sh:targetSubject** or **sh:targetObjectsOf**).
2. Define the Constraint Condition in SPARQL. Typically a SELECT query that returns any nodes violating the rule.
3. Check for Violations. If the query returns any results, the graph fails to conform. Otherwise, it passes.

This SPARQL-based method is quite flexible—constraints can incorporate arbitrary graph patterns and even aggregations. Additional Constraints Expressible in SHACL include:

- Unary Inclusion Dependencies: Restricting which nodes can appear as property values.
- (Conditional) Inclusion Dependencies (IND/CIND): Ensuring that values in one property subset must appear in another, potentially under certain conditions.
- Relation Functionality: Enforcing a property can have at most one value (i.e., functional).
- Inverse Relation Functionality: Ensuring that the property's inverse is functional (each value node points back to at most one subject).