

Few-Shot Knowledge Validation using Rules

Michael Loster
michael.loster@hpi.de
Hasso Plattner Institute
Germany

Davide Mottin
davide@cs.au.dk
Aarhus University
Denmark

Paolo Papotti
paolo.papotti@eurecom.fr
EURECOM
France

Jan Ehmüller
jan.ehmuller@student.hpi.de
Hasso Plattner Institute
Germany

Benjamin Feldmann
benjamin.feldmann@student.hpi.de
Hasso Plattner Institute
Germany

Felix Naumann
felix.naumann@hpi.de
Hasso Plattner Institute
Germany

ABSTRACT

Knowledge graphs (KGs) form the basis of modern intelligent search systems – their network structure helps with the semantic reasoning and interpretation of complex tasks. A KG is a highly dynamic structure in which facts are continuously updated, added, and removed. A typical approach to ensure data quality in the presence of continuous changes is to apply logic rules. These rules are automatically mined from the data using frequency-based approaches. As a result, these approaches depend on the data quality of the KG and are susceptible to errors and incompleteness.

To address these issues, we propose COLT, a few-shot rule-based knowledge validation framework that enables the interactive quality assessment of logic rules. It evaluates the quality of any rule by asking a user to validate only a few facts entailed by such rule on the KG. We formalize the problem as learning a validation function over the rule’s outcomes and study the theoretical connections to the generalized maximum coverage problem. Our model obtains (i) an accurate estimate of the quality of a rule with fewer than 20 user interactions and (ii) 75% quality (F_1) with 5% annotations in the task of validating facts entailed by any rule.

CCS CONCEPTS

• **Computing methodologies** → *Neural networks*; **Active learning settings**; *Gaussian processes*; **Rule learning**; **Reasoning about belief and knowledge**; **Semantic networks**; **Supervised learning by classification**.

KEYWORDS

knowledge validation, logic rules, knowledge graph, knowledge base, neural networks, deep kernel learning, few-shot learning

ACM Reference Format:

Michael Loster, Davide Mottin, Paolo Papotti, Jan Ehmüller, Benjamin Feldmann, and Felix Naumann. 2021. Few-Shot Knowledge Validation using Rules. In *Proceedings of the Web Conference 2021 (WWW ’21)*, April 19–23, 2021, Ljubljana, Slovenia. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3442381.3450040>

This paper is published under the Creative Commons Attribution 4.0 International (CC-BY 4.0) license. Authors reserve their rights to disseminate the work on their personal and corporate Web sites with the appropriate attribution.

WWW ’21, April 19–23, 2021, Ljubljana, Slovenia

© 2021 IW3C2 (International World Wide Web Conference Committee), published under Creative Commons CC-BY 4.0 License.

ACM ISBN 978-1-4503-8312-7/21/04.

<https://doi.org/10.1145/3442381.3450040>

1 INTRODUCTION

In recent years, systems such as DeepDive [40] and Knowledge Vault [11] have made it increasingly easy to automatically construct vast knowledge graphs (KGs) from structured and unstructured data. These systems often consist of many different components designed to extract and integrate facts from multiple sources. Many of these components are based on machine learning techniques, which are rarely error-free. Thus, despite dramatic improvements using deep learning techniques, it is not possible to ensure the correct extraction, linkage, and discovery of relationships between textual entities, which form the facts of a knowledge graph. Knowledge graphs suffer not only from incorrect but also from missing facts, which in their sum negatively affect all downstream applications. Correcting these errors by adding missing or removing incorrect facts represents a key task that can be approached from two different directions.

The first, often used to add missing facts, leverages *knowledge graph embeddings* (KGE) [44], such as TransE [6] or HolE [32], in training link prediction models. These models are widely studied and can be used to predict new facts; however, their quality depends on the correctness and completeness of the underlying KG. While correctly capturing large portions of the graph, they fail to model noisy or underrepresented entities and relationships [35]. Additionally, KGEs are difficult to interpret as they cannot *explain* why a fact should be included in the KG.

The second approach is applying *logic rules*, which can be derived by rule learning systems, such as AMIE [13, 14, 25] and RuDiK [34]. By executing such rules on a KG, missing facts can be added and incorrect facts removed. For example, the rule

$$\text{isMarriedTo}(a, c) \wedge \text{livesIn}(c, b) \Rightarrow \text{livesIn}(a, b)$$

states that if a person a is married to another person c , they most likely live in the same place (b). By applying this rule, i.e., when the left-hand-side of the rule is satisfied, a new `livesIn` fact between a person and a place can be added to the KG. While rules are easier to interpret than embeddings, their discovery also suffers from noisy and missing KG data. Moreover, rules are rarely completely correct or completely incorrect. Commonly, a rule applies only to a certain percentage of KG facts [13, 34], without providing exact and complete information in which cases it can be safely applied.

Ideally, a rule should only be applied if it contributes to the quality improvement of a KG. Although rule learning systems come with a statistical measure for the support of a discovered rule, this measure is derived from the KG facts and assumes their correctness.

In practice, this dependence often leads to incorrect estimates of a rule’s confidence value [2]. Most systems refer to the confidence of a rule as a value between 0 and 1, describing how likely the rule in question implies valid facts, which in turn defines the validity of the rule itself.

Addressing such data quality issues of a KG requires *external information* to make necessary corrections. The natural way to verify facts is to consult experts on the topics covered by the KG at hand. As KGs often consist of millions of facts, a comprehensive manual verification becomes infeasible. Thus, we aim to radically reduce the manual effort by efficiently utilizing existing domain knowledge.

To retain the understandability of rule-based approaches, we focus on a solution that (i) computes reliable rule confidence measures and (ii) complements the rules by a probabilistic model that validates their resulting facts.

To address the data quality problems of KGs, we train a classifier for each rule that balances the exploration of new facts and the exploitation of already learned knowledge. This allows us to effectively reduce the number of annotations to a few interactions also referred to as a *few shots*. Unlike logic rules, the classifiers contain additional knowledge in the form of (i) information captured by the KGEs and (ii) user feedback, enabling the classifiers to make better predictions than the original rule. The learned classifiers allow the conditional application of each rule, according to the predicted validity of its facts. Compared to a data-driven estimation approach, these predictions enable us to calculate a rule’s confidence with far greater accuracy. Although the learned classifier can be used as a replacement for the learned rule, it reduces the role of the rule only to a limited extent. As the classifier itself is quite complex and therefore difficult to interpret, the rule serves as an explanation of the classifier, making its behavior more comprehensible. Beyond its use in curating data quality problems, the creation of precise rule confidence estimates also plays a fundamental role in reasoning systems, such as those used in fact checking [1]. In these applications, the quality of the conclusion depends on the precise measurement of the rule quality. However, verifying the correctness of a large number of facts to precisely assess rule quality is cost-intensive. We therefore reduce the amount of work by utilizing a sampling strategy. This strategy limits the effort for manually verifying facts to an adjustable budget, while simultaneously maximizing the benefit from the verified facts.

Our core contribution constitutes the development of a model that theoretically and empirically aims to validate the confidence and precision of logical rules by posing as few questions as possible to a human oracle. In doing so, we answer the key question of whether there is a model capable of achieving maximum precision, provided that 100% of the annotated data is available. By reducing the problem to maximum coverage, we show that such a model exists, but has its drawbacks. As active learning methods guarantee no perfect precision, which is confirmed by our experiments, we turn to Gaussian processes (GP), that theoretically ensure maximum precision as the number of interactions approaches infinity. Our proposed approach learns the characteristics of a particular rule from humans who annotate a small number of its generated facts and addresses the limitations of GPs (e.g., instability of matrix inversion) by means of deep kernel learning.

We summarize our *contributions* as follows:

- (1) We propose the COLT framework capable of assessing the quality and confidence of a particular rule by using KG embeddings and expert-validated facts (Section 3).
- (2) Using a classifier, we enable the conditional application of declarative rules and compute their confidence (Section 4).
- (3) We establish a connection between our problem, the weighted coverage problem (Section 4.1), and quality-preserving Gaussian processes (Section 4.2).
- (4) We show how our interactive learning approach effectively exploits embeddings to improve the classifier with minimal human effort. Using only 20 user interactions, we halve the error in confidence obtained with rule learning systems (Section 6).
- (5) We publish our dataset consisting of 26 rules with more than 23 000 annotated facts.

2 RELATED WORK

This section covers the topics related to our method, namely, rule discovery and their interactive application, followed by graph embeddings and hybrid methods using both rules and embeddings.

Rule discovery. The derivation of logic rules from KGs has been investigated for many years [8]. Recent systems, such as AMIE [13, 14] and RuDiK [34], can derive rules from large KGs using structural information, such as frequently occurring graph patterns. The mined rules can be used to derive new facts, find errors, derive conclusions, and better understand the underlying data. However, it is difficult to automatically estimate the quality of a discovered rule [2]. Although AMIE and RuDiK report statistical measures for every rule, such as the standard and PCA confidence, these metrics are based on pattern frequencies that are affected by data quality issues in the KG. We improve the confidence assessment of a rule by including a user as an additional source of knowledge. As such, the quality assessment of a rule does not depend solely on the KG.

Interactive rule execution. Different approaches have studied how to execute KG rules for a specific task with user involvement. For data repair, potential KG corrections are derived from a set of rules and suggested to the user [3, 12]. The user validates these corrections, thus improving the underlying KG until it is transformed into a consistent state. This line of work assumes that the given rules are correct and reliable, such as those manually written by experts. Our approach handles rules that are automatically generated, thus possibly approximate or erroneous. Our solution can be seen as a preliminary step to create an understanding of the rules to be fed to the data repair step.

Interactive error detection has also been studied for relational data [17, 19, 30]. However, such methods cannot be directly applied to graph data; while a rule can be materialized into a relational model by flattening a portion of the KG, the resulting “view” disregards other information in the graph.

Knowledge graph embeddings. Knowledge graph embeddings (KGEs) provide a compact representation of a KG by projecting its entities and relations into an n -dimensional space while retaining its structural information. These embeddings are used to improve tasks [44], such as relation extraction [49] and link prediction [27].

Besides semantic matching models [33], translational methods, such as TransE [6] and its extensions [21, 45], are widely used for

the generation of embeddings. These methods maximize a score associated with each fact. Often interpreted informally as the plausibility of a fact, the score is a distance measure between the entities of the evaluated predicate. As these methods use the KG facts to produce a representation of the graph, data quality problems, such as incorrect facts, directly affect the quality of the generated embeddings [35] and thus also downstream tasks.

We address quality issues by learning classifiers that incorporate human knowledge during their training and do not solely rely on the KG data. By integrating user feedback, the classifier learns how to counteract poor embeddings and mitigate their negative effects. This limits the impact of noisy and missing KG information.

KGes and rules. Although all methods use the KG facts for computing KGEs, some approaches can consume additional information, such as textual descriptions [48], literals [15], or logical rules [16, 50]. We focus on the approaches that rely on rules during their learning process, as they are closest to our work.

One approach jointly models logical rules and fact triples in a unified framework, representing triples as atomic and rules as complex formulae, while learning takes place by minimizing a global loss function that spans both formulae [16]. Another system creates high-quality rules from noisy KGs using KGEs to control the rule generation process [20]. A third approach uses rules to generate new triples for the training of the KGEs, which are used to derive new rules [50]. These works address a problem different from ours: we train a classifier on top of existing KGEs and existing rules that counteracts poor embedding quality by incorporating user feedback into its learning process. Our goal is to interactively learn a model that can both assess the quality of a given rule and predict when the rule should be applied.

Active learning. Methods that actively sample data points to be labeled are called active learning methods [38]. They exploit the properties of model and data and perform well with human assessors and exploratory tasks [24, 28]. Our work builds on the general idea of active learning but establishes important connections with weighted coverage and GP-UCB algorithm based on Gaussian processes [41] to devise quality-preserving solutions. For consistency, we evaluate the most common active learning strategies described in Section 4.4.

3 BACKGROUND AND PROBLEM

A *knowledge graph* (KG) is a triple $G : \langle \mathcal{E}, \mathcal{R}, \mathcal{F} \rangle$, where \mathcal{E} is a set of entities, \mathcal{R} a set of relationships, and $\mathcal{F} \subseteq \mathcal{E} \times \mathcal{R} \times \mathcal{E}$ is a set of facts. For instance, a general-purpose KG contains the entities "Max Planck" and "Germany", connected via "isCitizenOf" relationship. A triple (e_1, r, e_2) with $e_1, e_2 \in \mathcal{E}$, and $r \in \mathcal{R}$ is called a *fact*; as r is a relationship in the pair (e_1, e_2) , each fact can be equivalently represented as an *atom* $r(e_1, e_2)$. The set of triples constitutes the knowledge graph, also known as knowledge base [9], information graph [29], or heterogeneous information network [39].

A knowledge graph might contain errors, such as ("Max Planck", "isCitizenOf", "China"). The set of all true facts (knowledge) is denoted by $\mathcal{K} \subseteq \mathcal{E} \times \mathcal{R} \times \mathcal{E}$, and we assume that at least a part of the knowledge graph is correct, i.e., $\mathcal{F} \cap \mathcal{K} \neq \emptyset$. Partial correctness constitutes a reasonable assumption regarding the soundness of

the KG construction process. An *atom* is the smallest logic statement composed by facts, including variables. For example, the atom $hasChild(x, "John")$ denotes all entities (by the variable x) that have a child named John.

DEFINITION 1. A rule $A_1 \wedge A_2 \wedge \dots \wedge A_n \Rightarrow C$ is a *logical implication* consisting of a head and a body. The head is a single atom C , and the body is a conjunction of atoms $A_1 \wedge \dots \wedge A_n$.

The rule $\rho_0 : hasChild(a, c) \wedge isCitizenOf(a, b) \Rightarrow isCitizenOf(c, b)$ describes that if a has a child c and is a citizen of country b , their child must also be a citizen of b . This is a *positive* rule, which asserts the existence of specific facts in the head, given the body. We also consider *negative* rules to identify inconsistencies in the data [3], which denote sets of facts that lead to contradictions. For example, the rule $\rho_1 : hasChild(x, y) \wedge isMarriedTo(x, y) \Rightarrow \perp$ identifies possible inconsistencies in the KGs. The rule can also be expressed as $\rho_1' : hasChild(x, y) \Rightarrow \neg isMarriedTo(x, y)$ to derive false facts from the KG.

The *instances* $I_G(\rho) \subseteq \mathcal{E} \times \mathcal{R} \times \mathcal{E}$ of a rule ρ in the knowledge graph G are the set of facts expressed by the right-hand-side of the rule, assuming the left-hand-side is true. For example, the instances of rule ρ_0 are facts of the form $isCitizenOf(c, b)$, given that both $hasChild(a, c)$ and $isCitizenOf(a, b)$ are true.

Ideally, if we knew the true facts \mathcal{K} , we could validate the *confidence* of the rule on the knowledge graph G by computing the instances in \mathcal{K} that are correctly captured by a rule ρ via the ratio $\frac{|I_{\mathcal{K}}(\rho) \cap I_G(\rho)|}{|I_{\mathcal{K}}(\rho)|}$. As \mathcal{K} is, in reality, unavailable, we model the validation as a *benefit function* $f : I_G(\rho) \rightarrow [0, 1]$ over the instances $I_G(\rho)$. f returns a value that converges to 1 if the fact $(e_1, r, e_2) \in I_G(\rho)$ is confidently correct, i.e., $(e_1, r, e_2) \in \mathcal{K}$, otherwise it converges to 0. Access to (parts of) \mathcal{K} can be obtained by an external oracle, e.g., by a domain expert who verifies whether a fact is true or not. Assessing the validity of facts by consulting an oracle can be costly, which is why we limit the number of evaluations to a budget B , thus bounding the evaluation effort.

The *rule-driven knowledge validation* problem aims at finding a benefit function f within a budget B that maximizes the knowledge over the instances $I_G(\rho)$ of rule ρ :

PROBLEM 1. Given a knowledge graph $G : \langle \mathcal{E}, \mathcal{R}, \mathcal{F} \rangle$, a rule ρ with instances $I_G(\rho) \subseteq \mathcal{E} \times \mathcal{R} \times \mathcal{E}$, and a budget B , find a set of facts L

$$\arg \max_{L \subseteq I_G(\rho)} \sum_{x \in L} f(x) \quad \text{subject to } C(L) \leq B$$

where C represents the sum of costs of all the facts in L .

In many practical scenarios, the cost of a fact is 1, as the budget quantifies the number of questions to the oracle. In our setting, we make this assumption but note that the problem, even if the benefit f is known, is NP-hard [43] with a non-constant cost.

4 THE COLT FRAMEWORK

In practice, the benefit function f in Problem 1 is unknown. As such, excerpting the knowledge f from an oracle makes the direct optimization in Problem 1 impossible. To circumvent this issue, COLT makes two assumptions. **First**, the benefit function can be equivalently represented by a *classifier* \tilde{f} on the facts trained on the evaluations seen so far. **Second**, since facts are related to one another,

we assume the existence of a *similarity* $\kappa: (\mathcal{E} \times \mathcal{R} \times \mathcal{E}) \times (\mathcal{E} \times \mathcal{R} \times \mathcal{E}) \mapsto [0, 1]$ among pairs of facts. Such similarities naturally manifest in knowledge graphs by means of relationships among entities. As facts themselves are too sparse and fail to capture implicit KG relationships, it takes a different approach to capture their similarity. Section 5 provides more insights on how to compute expressive similarities.

Our COLT framework, presented in Algorithm 1, proceeds iteratively by using the information in the classifier and asking the user to evaluate facts until a specified budget B is depleted. In the beginning, the classifier is initialized (Line 2) using the instances and the similarity function. At every iteration t , the algorithm selects a new fact x_t using the classifier's knowledge (Line 5). The user then validates the fact (Line 6) providing a value $y_t = 1$ if the fact is true and $y_t = 0$ otherwise. Finally, the algorithm updates its beliefs by incorporating the true value y_t for the fact x_t . The classifier serves as a proxy for the unknown function f , which encodes the true facts \mathcal{K} for the instances $\mathcal{I}_G(\rho)$ of rule ρ . This classifier is then used to (i) estimate a confidence value for the current rule and (ii) determine for which facts a rule is likely to be correct.

The classifier \tilde{f} is of crucial importance, as it should guarantee a high classification accuracy with few evaluated facts. Section 4.1 describes a theoretical connection to maximum set coverage and a greedy baseline that fully relies on the correctness of the similarity measure. In Section 4.2, we relax the constraint on the similarity with a model that combines neural networks and probabilistic Gaussian processes to overcome the rigidity of maximum coverage.

Algorithm 1 The COLT framework

Input: knowledge graph G ; rule instances $\mathcal{I}_G(\rho)$
Input: similarity κ ; budget B
Output: A classifier \tilde{f}

- 1: $U \leftarrow \mathcal{I}_G(\rho)$ ▷ Unlabeled facts
- 2: $\tilde{f} \leftarrow \text{INITIALIZE}(U, \kappa)$ ▷ Initialize classifier
- 3: $L \leftarrow \emptyset$ ▷ Initialize set of labeled facts
- 4: **for** $t = 1 \dots B$ **do**
- 5: $x_t \leftarrow \text{SELECT}(\tilde{f}, U, \kappa)$ ▷ Select next fact
- 6: $y_t \leftarrow f(x_t)$ ▷ User validates x_t
- 7: $L \leftarrow L \cup \{(x_t, y_t)\}$ ▷ Add validated fact to L
- 8: $\text{UPDATE}(\tilde{f}, \kappa, L)$ ▷ Update classifier using new fact
- 9: $U \leftarrow U \setminus \{x_t\}$ ▷ Remove x_t from U
- 10: **return** \tilde{f}

4.1 COLT-MC: A maximum coverage baseline

We observe a connection between Problem 1 and the *maximum coverage* (MAXCOVER) problem [31] and devise a greedy baseline for our problem. Given some sets over a universe of elements and a fixed number B , the MAXCOVER problem finds B sets that overall contain the maximum number of elements from the universe.

To appreciate the commonalities between Problem 1 and MAXCOVER, first note that every fact x identifies a set $\kappa(x, \cdot)$ in which each of its elements has a different similarity weight. Second, the benefit $f(x)$ of each fact x corresponds to the sum of similarities $f(x) = \sum_{x'} \kappa(x, x')$. As such, Problem 1 translates into finding a set of facts that *maximizes the weighted coverage* identified by each fact's similarity set $\kappa(x, \cdot)$. The variant of MAXCOVER in which every element's weight varies with respect to the set

containing the element is called *generalized maximum coverage* (GENMAXCOVER) [7]. Generalized maximum coverage admits a greedy $(1-1/e)$ -approximate solution. Also, the generalized version allows each element's cost to be different, which is not our case as the cost is determined by the number of facts evaluated by a user. We relax the GENMAXCOVER assumptions and devise a greedy solution for the problem with fixed cost. Such a greedy scheme selects at each iteration the fact that maximizes the *weighted marginal gain*,

$$\Delta(x|L) = \sum_{x' \in \mathcal{I}_G(\rho)} \kappa(x, x') - \max_{x'' \in L} \kappa(x', x'') \quad (1)$$

which quantifies the increment in similarity if the fact x is added to the set L . As soon as the set L contains B facts, the algorithm stops. We refer to the greedy baseline in Algorithm 2 as COLT-MC.

The classifier \tilde{f} is a 1-nearest-neighbor (1-NN) classifier that provides a binary label for each fact. The evaluation of \tilde{f} on unlabeled facts returns the label of the most similar fact among those evaluated so far. In the end, the approximation of the benefit function f provided by the classifier \tilde{f} is

$$\tilde{f}(x) = \begin{cases} f(x) & \text{if } x \in L \\ f(\arg \max_{x' \in L} \kappa(x, x')) & \text{otherwise} \end{cases}$$

One convenient property of COLT-MC is that it ensures maximum quality if all the facts in the instances of the rule are evaluated. We observe this property also empirically in Section 6.3.

Algorithm 2 COLT-MC

- 1: **function** INITIALIZE(U, κ)
- 2: $\tilde{f} \leftarrow 0$ ▷ Initialize uniform classification
- 3: **function** SELECT(f, U, κ)
- 4: **return** $\arg \max_{x \in U} \Delta(x|L)$ ▷ See Eq. 1
- 5: **function** UPDATE(\tilde{f}, κ, L)
- 6: $(x_t, y_t) \leftarrow L_t$ ▷ Last labeled fact
- 7: $\tilde{f}(x_t) \leftarrow y_t$

4.2 COLT-GP: A learning-based solution

The COLT-MC baseline implicitly assumes that both the similarity function κ on the facts as well as the user's validation are correct and reliable. To relax these assumptions, we propose a method based on deep kernel learning (DKL) that combines the advantages of neural networks and Gaussian processes (GPs) (Section 4.3).

To estimate the uncertainty value of \tilde{f} , our proxy of the benefit function, we assume that the benefit function at step t is a Bernoulli-distributed random variable sampled from a logit-transformed Gaussian distribution. The Gaussian random variable is represented by a Gaussian process [5, Ch. 6], which models points in space as individual Gaussians related through a kernel function representing the similarity κ between facts.

A GP is determined by the mean μ and covariance matrix, C in which each element $C_{ij} = \kappa(x_i, x_j)$ represents the similarity between a triple pair x_i and x_j . GPs define a distribution over functions where the *prior* of such distribution is a Gaussian $\mathcal{N}(0, C)$ with mean 0 and covariance $C(\mathcal{I}_G(\rho), \mathcal{I}_G(\rho))$ among all instances in $\mathcal{I}_G(\rho)$. In practice, the prior is never computed. Instead, we are interested in the *posterior* probability $P(x \in \mathcal{K}|L)$ of a fact x belonging to the true facts \mathcal{K} knowing the labeled facts L . The posterior

of a GP is Gaussian; however, in order to classify the facts, the function sampled from the posterior needs to be transformed to return a value between 0 and 1, with high probability. This transformation is made by a *logistic sigmoid* function $s(x) = \frac{1}{1+e^{-x}}$. The final posterior is the expectation of the predictions over the sample of logistic-transformed Gaussian functions:

$$P(x \in \mathcal{K}|L) = \int s(f_*)P(f_*|L) df_* \quad (2)$$

The integral in Eq. 2 is analytically intractable but can be approximated with sampling methods or analytical approximations, such as the Stochastic Variational Inference (SVI) used in Section 4.3.

The choice of GPs for classification has two significant implications. First, a GP is a Bayesian model that allows sampling using the predictions from the posterior. As the choice of the sampling strategy is critical for such models, we analyze multiple choices for different sampling strategies in Section 4.4. Second, GPs are non-parametric models; thus, they can easily update their beliefs by using Bayes' theorem on the posterior computation.

We are now ready to describe COLT-GP (Algorithm 3), represented in Figure 1. For the sake of completeness, Algorithm 3 reports on the prior initialization of the random variable (Line 2), even though it is never explicitly computed. The model parameters (μ, C) are updated using the posterior inference in Eq. 2. Each iteration step samples an unlabeled fact using the *SELECT* function, which retrieves an element according to one of the strategies in Section 4.4. The sampled fact is then presented to a user who evaluates its correctness. For example, given the supposed fact \langle "Jon Landau", "produced", "Avatar" \rangle , a user verifies whether or not Jon Landau produced the movie Avatar. The labeled triple is then added to the training set L .

The posterior calculates the probability of a fact being true, which is estimated by Eq. 2. To convert this probability into a label of 0 or 1 for each fact x , $\tilde{f}(x)^1$ returns 1 if the probability ≥ 0.5 , otherwise 0.

$$\tilde{f}(x) = \begin{cases} 1 & \text{if } P(x \in \mathcal{K}|L) \geq 0.5 \\ 0 & \text{otherwise} \end{cases}$$

We experimented with different probability thresholds and observed good results by choosing a threshold of 0.5. This is also confirmed by our experiments in Section 6.1.

Algorithm 3 COLT-GP

- 1: **function** INITIALIZE(\tilde{f}, U, κ)
 - 2: $\tilde{f} \sim s(\mathcal{N}(0, C))$ \triangleright Gaussian prior, not explicitly computed
 - 3: **function** SELECT(f, U, κ)
 - 4: **return** x_t sampled with one strategy in Sec. 4.4
 - 5: **function** UPDATE(\tilde{f}, κ, L)
 - 6: $\tilde{f} \leftarrow$ Update the posterior using SVI [47] for Eq. 2
-

4.3 Learning with deep kernels

The inference step for GPs requires the inversion of their covariance matrix, which grows with the number of evaluated facts. The numerical instability of this inversion coupled with modest performance on high dimensional data [10] is detrimental to the final

¹With a little abuse of notation \tilde{f} indicates both the posterior and the classifier.

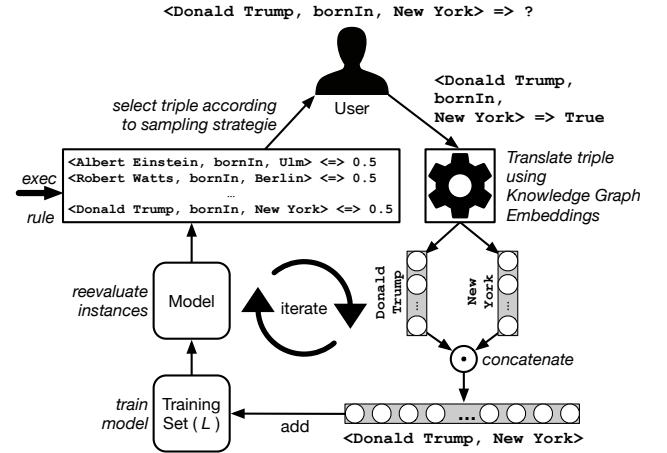


Figure 1: Overview of our interactive model-building.

classification quality. To circumvent the GPs' deficiencies, we leverage deep kernel learning (DKL) [46], which scales linearly with the number of evaluated facts and combines the strengths of both neural networks and GPs in one unified model. We use DKL to learn a flexible similarity function, while at the same time incorporating the user's evaluations of the facts. To learn the model parameters and perform posterior inference, we use stochastic variational inference (SVI) [47], which enables the use of non-Gaussian likelihoods.

Equation 3 defines the general model architecture as a composition of functions. It consists of two main components: a set of fully connected network layers (FC_1 - FC_4) and one GP layer applied to the network's output vector. This combination results in a deep probabilistic neural network that meets our requirements.

$$\phi(\mathbf{x}) = GP(FC_3(FC_2(FC_1(\mathbf{x})))) \quad (3)$$

The neural network part of the model (ϕ) consists of three fully connected layers $FC_n(\mathbf{x}) = \sigma(\mathbf{W}\mathbf{x} + b)$, where \mathbf{W} and b correspond to the weight and bias terms, and σ denotes the activation function, which in our case is set to ReLU. The network part of the model effectively performs a dimensionality reduction, by transforming n -dimensional input vectors $\mathbf{x} = (x_1 \dots x_n)$ into a low-dimensional vector representation, which is then passed to the GP layer.

The GP part of the model is connected to the output layer of the network and consists of j Gaussian Processes $g_1 \dots g_j$, each having their own corresponding kernel $k_1 \dots k_j$ operating on a subset of the vector generated by the network. When selecting the number of GPs in this layer, we follow Wilson et al. [47] and use one GP for each dimension in the output vector ($j = 4$). To obtain the final result, the individual GPs are first additively combined and then transformed by an observation model, which in our case is a Bernoulli likelihood, to receive the final result Y .

4.4 Sampling strategies

As stated in Section 4.2, the sampling strategy is at the heart of our COLT-GP algorithm, as it is responsible for the incremental selection of new facts. A good sampling strategy aims to select training examples of high information density that are likely to improve

classifier performance when used for training. In the following, we present the sampling strategies considered in our experiments.

Random. Random sampling is the simplest sampling technique: it randomly selects training examples from the data and adds them to the training dataset, ignoring any prior knowledge.

Maximum uncertainty. Maximum uncertainty sampling (MAXU) uses the model itself to assess how valuable the labeling of each data point is in terms of confidence gain [26]. The sampling uses the model’s uncertainty on unlabeled data points: the data points for which the classifier is most uncertain are those closest to its decision boundary. The incremental labeling of the data points with maximum uncertainty leads to the refinement of the decision boundary. In binary classification, the binary entropy measure determines the classifier’s uncertainty with respect to its classification where p is the uncertainty provided by the model on the class label $X = 1$ (e.g., logistic regression provides p as output):

$$H(X) = -p \log_2(p) - (1-p) \log_2(1-p)$$

$$\text{with } X = \begin{cases} 1 & \text{with probability } p \\ 0 & \text{with probability } 1-p \end{cases}$$

As mentioned in [36], a disadvantage of this strategy can be its sensitivity to outliers, which is tackled by GP-UCB.

GP-UCB. The Gaussian process upper confidence bound algorithm (GP-UCB) by Srinivas et al. [41] defines a sampling criterion that aims at finding a trade-off between exploring the function space and exploiting maximal function areas. To this end, the authors propose the following sampling strategy:

$$x_t = \arg \max_{x \in U} \mu_{t-1}(x) + \sqrt{\beta_t \sigma_{t-1}(x)}$$

where U denotes the unlabeled facts as defined in Algorithm 1. The exclusive selection of elements $x \in U$ where the model has either a high variance (σ_{t-1}) or a high expected reward (μ_{t-1}) cannot be regarded as the optimal strategy, as useful information is neglected and the model’s ability to generalize is compromised. To overcome this issue, the proposed GP-UCB sampling strategy strikes a balance between *exploration* by selecting elements where the model has a certain degree of uncertainty (large σ_{t-1}), and *exploitation* by choosing elements with high belief (large μ_{t-1}). By simultaneously optimizing both criteria, the strategy seeks to achieve an equilibrium between exploration and exploitation.

The parameter β_t trades-off exploration and exploitation. A large β_t corresponds to more exploration, while a small β_t leads to more exploitation. The GP-UCB strategy ensures a logarithmic growth of the regret, the deviation in quality with respect to the optimal sampling, by setting $\beta_t = 2 \log(|I_G(\rho)| t^2 \pi^2 / 6\delta)$ with $\delta \in (0, 1)$. In other words, the regret on the evaluations tends to 0 when the number of validated facts tends to infinity.

According to the above formula, β_t starts at a relatively small value, which encourages it to be more exploitative. Along with the increasing number of sampling steps t , β_t also increases, making it more exploratory. We use this policy to update β_t during our experiments.

5 COMPUTING SIMILARITIES

The choice of an adequate similarity measure κ depends on the specifics of the data and the application domain. We avoid the need to design an ad-hoc similarity measure by proposing a flexible one based on *knowledge graph embeddings* [44] to represent a graph in a low-dimensional space. These methods implicitly capture structural and semantic information and are agnostic to the application domain. Although errors and missing data in the knowledge graph affect the quality of the generated embeddings [35], they constitute a solid starting point. The design of such embedding methods is outside of the scope of this work, thus we refer to [44] for a comprehensive survey on these methods.

For our purposes, we employ HypER [4], an expressive and fast embedding method. For each fact $\langle s, r, o \rangle$ with subject s , object o , and relationship r , HypER computes 200-dimensional vectors \mathbf{s} , \mathbf{o} , and \mathbf{r} . HypER maximizes the likelihood of the existence of a fact in the knowledge graph. Additionally, the method aggregates information from the subject and the relationship employing convolution operators. The object \mathbf{o} in the vector space is then computed as a non-linear transformation of the subject and relationship vector.

These embeddings are calculated for all subjects, relationships, and objects in a knowledge graph before any rule is processed. The final embedding \mathbf{x} of a fact $x = \langle s, r, o \rangle$ is the 400-dimensional vector concatenation $\mathbf{x} = \mathbf{s} \circ \mathbf{o}$ of the subject and object vectors. Since we consider one rule at a time, the relationship r in the rule body always remains the same and can be omitted.

COLT builds on HypER embeddings, but any choice of superior embedding technique would only enhance COLT’s performance.

COLT-MC: The similarity score between fact $x = \langle s, r, o \rangle$ and fact $x' = \langle s', r', o' \rangle$ is the cosine similarity among the fact vectors:

$$k_{\text{cos}}(x, x') = \frac{\mathbf{x} \cdot \mathbf{x}'}{\|\mathbf{x}\| \|\mathbf{x}'\|} = \frac{(\mathbf{s} \circ \mathbf{o}) \cdot (\mathbf{s}' \circ \mathbf{o}')}{\|\mathbf{s} \circ \mathbf{o}\| \|\mathbf{s}' \circ \mathbf{o}'\|}$$

COLT-GP: The Gaussian processes included in COLT-GP are kernel methods, which naturally incorporate similarity scores. The default kernel for Gaussian processes is the RBF kernel:

$$k_{\text{RBF}}(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|}{2\sigma^2}\right)$$

The RBF kernel provides smoothness as the similarity exponentially decreases with an increase of the Euclidean distance. The parameter σ is a learned parameter in the optimization. GPs with RBF kernels implicitly project the points into an infinite-dimensional space so that the GP layer represents a hidden layer with an infinite number of neurons, thereby significantly increasing the model’s flexibility. The transformation induced by deep kernel learning (DKL) in Section 4.3 further modifies the embedding space through a parametric non-linear function g . The final kernel is:

$$k_{\text{DKL}}(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|g(\mathbf{x}) - g(\mathbf{x}')\|}{2\sigma^2}\right)$$

6 EXPERIMENTAL EVALUATION

In the experimental evaluation, we answer the following questions.

Q1: Is COLT-GP able to learn the truth value for each rule (Sec. 6.1)? **Q2:** How do different sampling strategies affect the training process (Sec. 6.2)? **Q3:** How does COLT-GP fare in comparison to our strong

baseline COLT-MC and active learning (Sec. 6.3)? **Q4:** What is the relationship between user interactions, real confidence, and rule support (Sec. 6.4)?

Datasets. We use YAGO [42], an open-source knowledge base, to derive logical rules and graph embeddings. At the time of writing, YAGO2 comprises 948 358 triples, 36 relationship types, and 470 485 entities². We ran RuDiK [34] and AMIE [13] with default parameters on YAGO to obtain positive and negative rules. Out of 1 517 mined rules, 928 produced more than 5 500 triples. Each rule can be thought of as a query asking for triples in the KG that satisfy the rule’s left-hand side.

To keep the annotation process manageable, we selected from all generated rules those that did not exceed 5 500 generated instances, reducing the number of candidate rules to 589. Out of these 589 rules, we selected a set of rules that exhibit a variety in their properties: the number of atoms, the set of triplets induced, and the rule confidence. We decided to include mainly rules with at least three atoms for two reasons. First, they are more difficult for users to analyze and validate. Second, rules with fewer atoms typically express simple facts, such as inverse relations ($\text{isFatherOf}(x, y) \Rightarrow \text{hasFather}(y, x)$), that are easily captured by standard algorithms. Employing this selection strategy resulted in the 26 rules (22 positive; 4 negative) listed in Table 1. For negative rules whose execution resulted in a very high number of triples satisfying the rule, we randomly selected 1,000 triples for annotation. To cover a range of different output sizes, we annotated rules with a small output, as few as 40 triples, up to a rule with 5 269 triples. Similarly, we annotated rules that are mostly right, rules that are nearly always wrong, and rules with mixed confidence. For the listed 26 rules, we manually annotated a total of 23 324 triples, of which 5 524 corrected errors and missing facts in the underlying KG. To encourage further research in this area, we provide the annotated dataset for public access³.

Algorithms. In addition to COLT-MC (Section 4.1), we evaluate the COLT-GP algorithm (Section 4.2) against two active learning baselines. Both baselines follow the basic active learning procedure as described in [38]. While the first baseline employs a linear regression model (AL-LogReg), the second uses a more expressive SVM with RBF kernel (AL-SVM) [37]. Following the results of Section 6.2, we choose maximum uncertainty as the selection strategy for both active learning approaches.

Training details. To train the model presented in Section 4.3, we use stochastic variational inference (SVI) and gradient descent. We employ the ADAM optimizer [23] with default parameters $\alpha = 0.001$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, and $\epsilon = 10^{-8}$. We jointly train the weights of the neural network and the parameters of the GPs with batch size 32 until convergence or a maximum of 1000 epochs. To initialize the weight matrices, we used the initialization scheme proposed by He et al. [18]. We train the models on a Linux machine with 2×2.20GHz CPUs, each with 10 cores, 251 GB RAM and one Nvidia Titan X Pascal GPU.

Runtime. Recall that this work focuses on the study of methods for quality assessment and confidence estimation of logic rules

using a few human-validated rule instances. While a thorough evaluation of the runtime deviates from the main focus of this work, we observe that our methods run in real-time in less than 10 seconds on 20–100 evaluated instances. This result, obtained on a commodity machine, endorses the use of the COLT framework and its most expressive COLT-GP method on production systems. While we did not explicitly measure the annotation time of the users, we expect it to vary depending on the user’s experience and several external factors that are independent of the proposed algorithm. In our experience, most rules could be annotated within a time frame of 30 seconds to 5 minutes, depending on their complexity.

6.1 Q1: Learning rule characteristics

Here we aim to show that COLT-GP is able to learn the characteristics of a rule from the instances it generates. We divide the annotated instances into training and test sets, using a ratio of 70 to 30. We train COLT-GP on the training and test the ability to predict whether each rule’s instance holds for the unseen instances in the test set. This scenario primarily aims to test the learned model’s generalization capabilities and, at the same time, represents the most challenging classification scenario.

Figure 2 summarizes the results by showing the ROC curves for the annotated rules as well as the mean and variance of each of the rule’s ROC curves. The ROC curves of the individual rules and, hence, also the average ROC curve are all significantly above the dashed line representing the behavior of a random classifier. COLT-GP attains 89.78% average AUC, which testifies the correctness of the classifier in predicting unknown facts. Even the lower end of the variance line attains 79.38% AUC, which is significantly above the performance of a random classifier. We conclude that COLT-GP can learn a meaningful model by using individual facts.

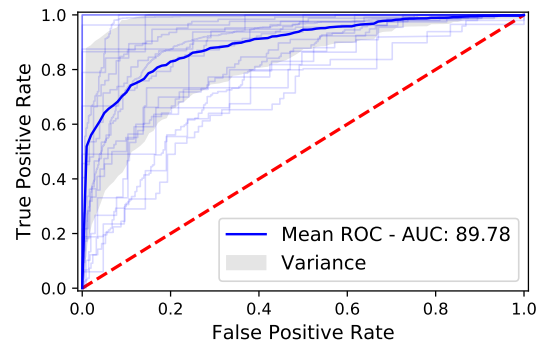


Figure 2: ROC curve for classifying relations.

6.2 Q2: Selecting the sampling strategy

The sampling strategies are responsible for collecting the data used for training. The purpose of this experiment is to investigate how the performance of our approach depends on the sampling strategy.

Experimental setup. With $\mathcal{I}_G(\rho)$ being the set of all instances created by the rule ρ , we follow Algorithm 1, as specified in Section 4. To test the trained models, we make a prediction on all instances in $\mathcal{I}_G(\rho)$, which includes the instances added to the training set L . This creates a realistic evaluation scenario where instances used for training also remain in the KG. Without limiting the budget

²http://resources.mpi-inf.mpg.de/yago-naga/amie/data/yago2/yago2core_facts.clean.notypes.tsv.7z

³<https://hpi.de/naumann/projects/repeatability/datasets/colt-dataset.html>

| ID | Rule | #KG | #correct | #total | type |
|----|--|-------|----------|--------|------|
| 1 | $\text{dealsWith}(a, v_0) \wedge \text{isLeaderOf}(v_1, b) \wedge \text{wasBornIn}(v_1, v_0) \Rightarrow \text{dealsWith}(a, b)$ | 9 | 24 | 99 | pos |
| 2 | $\text{dealsWith}(v_0, b) \wedge \text{isCitizenOf}(v_1, a) \wedge \text{isLeaderOf}(v_1, v_0) \Rightarrow \text{dealsWith}(a, b)$ | 42 | 42 | 42 | pos |
| 3 | $\text{dealsWith}(v_0, b) \wedge \text{isCitizenOf}(v_1, a) \wedge \text{isLocatedIn}(v_1, v_0) \Rightarrow \text{dealsWith}(a, b)$ | 108 | 134 | 134 | pos |
| 4 | $\text{isCitizenOf}(v_0, b) \wedge \text{livesIn}(v_0, a) \Rightarrow \text{dealsWith}(a, b)$ | 18 | 664 | 734 | pos |
| 5 | $\text{diedIn}(a, v_0) \wedge \text{isKnownFor}(v_1, v_0) \wedge \text{livesIn}(v_1, b) \Rightarrow \text{diedIn}(a, b)$ | 808 | 815 | 838 | pos |
| 6 | $\text{diedIn}(a, v_0) \wedge \text{isLeaderOf}(v_1, v_0) \wedge \text{livesIn}(v_1, b) \Rightarrow \text{diedIn}(a, b)$ | 2,143 | 3,135 | 4,292 | pos |
| 7 | $\text{actedIn}(a, b) \wedge \text{created}(a, b) \Rightarrow \text{directed}(a, b)$ | 384 | 388 | 1,003 | pos |
| 8 | $\text{actedIn}(v_0, b) \wedge \text{created}(a, b) \wedge \text{directed}(v_0, b) \Rightarrow \text{directed}(a, b)$ | 421 | 414 | 804 | pos |
| 9 | $\text{hasWonPrize}(a, v_0) \wedge \text{hasWonPrize}(b, v_0) \wedge \text{hasAcademicAdvisor}(b, a) \Rightarrow \neg \text{hasAcademicAdvisor}(a, b)$ | 85 | 85 | 85 | neg |
| 10 | $\text{influences}(a, b) \Rightarrow \neg \text{hasAcademicAdvisor}(a, b)$ | 1,000 | 1,000 | 1,000 | neg |
| 11 | $\text{isCitizenOf}(v_0, a) \wedge \text{livesIn}(v_0, b) \Rightarrow \text{hasCapital}(a, b)$ | 43 | 50 | 734 | pos |
| 12 | $\text{dealsWith}(a, v_0) \wedge \text{hasCurrency}(v_0, b) \Rightarrow \text{hasCurrency}(a, b)$ | 7 | 10 | 293 | pos |
| 13 | $\text{hasCapital}(v_0, v_1) \wedge \text{hasCurrency}(v_0, b) \wedge \text{isLocatedIn}(v_1, a) \Rightarrow \text{hasCurrency}(a, b)$ | 18 | 46 | 49 | pos |
| 14 | $\text{dealsWith}(v_0, a) \wedge \text{hasOfficialLanguage}(v_0, b) \Rightarrow \text{hasOfficialLanguage}(a, b)$ | 52 | 72 | 668 | pos |
| 15 | $\text{hasOfficialLanguage}(v_0, b) \wedge \text{isLocatedIn}(v_0, a) \Rightarrow \text{hasOfficialLanguage}(a, b)$ | 18 | 62 | 126 | pos |
| 16 | $\text{hasOfficialLanguage}(v_0, b) \wedge \text{isLocatedIn}(v_1, a) \wedge \text{livesIn}(v_1, v_0) \Rightarrow \text{hasOfficialLanguage}(a, b)$ | 20 | 52 | 65 | pos |
| 17 | $\text{hasOfficialLanguage}(v_0, b) \wedge \text{isLocatedIn}(v_1, v_0) \wedge \text{livesIn}(v_1, a) \Rightarrow \text{hasOfficialLanguage}(a, b)$ | 20 | 16 | 40 | pos |
| 18 | $\text{influences}(a, v_0) \wedge \text{isCitizenOf}(v_0, b) \Rightarrow \text{isCitizenOf}(a, b)$ | 146 | 547 | 1,382 | pos |
| 19 | $\text{hasChild}(a, b) \Rightarrow \neg \text{isMarriedTo}(a, b)$ | 991 | 990 | 1,000 | neg |
| 20 | $\text{isLocatedIn}(v_0, b) \wedge \text{livesIn}(a, v_0) \Rightarrow \text{isPoliticianOf}(a, b)$ | 130 | 2,257 | 5,269 | pos |
| 21 | $\text{isCitizenOf}(v_0, b) \wedge \text{isLeaderOf}(v_0, v_1) \wedge \text{livesIn}(a, v_1) \Rightarrow \text{livesIn}(a, b)$ | 416 | 863 | 882 | pos |
| 22 | $\text{isMarriedTo}(a, v_0) \wedge \text{livesIn}(v_0, b) \Rightarrow \text{livesIn}(a, b)$ | 181 | 442 | 537 | pos |
| 23 | $\text{actedIn}(a, b) \wedge \text{created}(a, b) \Rightarrow \text{produced}(a, b)$ | 207 | 261 | 1,003 | pos |
| 23 | $\text{actedIn}(v_0, b) \wedge \text{directed}(a, b) \wedge \text{produced}(v_0, b) \Rightarrow \text{produced}(a, b)$ | 246 | 281 | 702 | pos |
| 25 | $\text{hasAcademicAdvisor}(v_0, a) \wedge \text{worksAt}(v_0, b) \Rightarrow \text{worksAt}(a, b)$ | 47 | 163 | 543 | pos |
| 26 | $\text{actedIn}(a, b) \Rightarrow \neg \text{wroteMusicFor}(a, b)$ | 1000 | 998 | 1,000 | neg |

Table 1: Generated rules and their statistics: number of KG facts that satisfy the rule, number of correct facts, total number of facts, and type of the rule (colors are for ease of reading).

B, L eventually contains the same elements as $I_G(\rho)$. We use this validation approach as it is widely used in the evaluation of active learning methods, where each iteration adds more instances to the training set. This setup enables us to show that our model improves incrementally, ideally leading to perfect quality (F1 score) if all instances are available and the information capacity of the proposed model is sufficiently large to fully capture the problem. If a model cannot achieve an F1-score of 100% when training and predicting on identical data, it is likely that its capacity is not high enough to perfectly separate the classification space. We examine this scenario in more detail in Section 6.3.

Evaluation. Figure 3 shows how the different sampling strategies affect precision, recall, and F-measure. As expected, the random selection of training instances proves to be the worst strategy. Although its precision initially increases to 78%, it then decreases and remains, on average, 3.93% and 6.22% below the precision value of GP-UCB and max-uncertainty (MAXU), respectively. In terms of F-measure, the random sampling strategy lags 4.42% and 6.72% behind the GP-UCB and MAXU strategies, regardless of the percentage of training instances. GP-UCB takes up to 6% of training instances to achieve a competitive precision value of 77.44% and remains below the precision value of MAXU throughout the experiment. Overall, the MAXU strategy provides the best results in F-measure, outperforming the other two strategies at all times.

GP-UCB’s inability to achieve a better performance than MAXU can be attributed to the β_t parameter discussed in Section 4.4. While

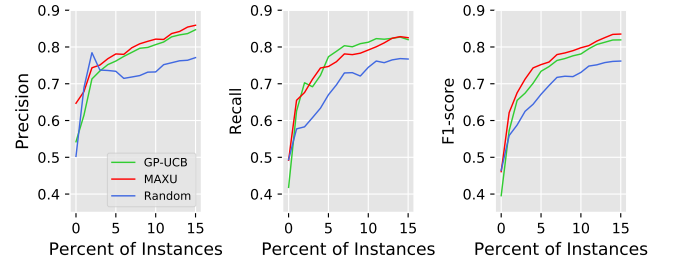


Figure 3: Impact of sampling strategies on predictions.

the formula in [41] for selecting β_t provides a theoretical bound on prediction quality, it appears to be too conservative for our problem. As a result, GP-UCB makes too conservative choices and is unable to outperform MAXU due to an unfavorable trade-off between exploration and exploitation. As the MAXU strategy achieves 75% F1 with only 5% training instances, we employ this strategy in all subsequent experiments.

6.3 Q3: Model performance

This experiment analyzes the behavior of COLT-GP w.r.t. the amount of training data and compares the performance to our COLT-MC and two active learning baselines. We aim to answer question (Q3) on the necessity of our learning strategy as opposed to active learning

methods and direct use of similarities between KG embeddings in COLT-MC.

Figure 4 shows precision, recall, and F1-score curves of all algorithms for an increasing amount of training data. AL-LogReg and AL-SVM rapidly fall behind the COLT algorithms. Although AL-SVM, leveraging its nonlinear RBF kernel, provides better results than AL-LogReg, neither model possesses sufficient capacity to meet the complexity of the problem. This is reflected in the performance metrics of Figure 4, where after using 15% (AL-SVM) and 20% (AL-LogReg) of the training data, improvements begin to stagnate for both models, so that their respective F-measures fluctuate around 84% and 80%. In contrast, COLT-GP and COLT-MC outperform the active learning approaches and continue to improve as training data increases.

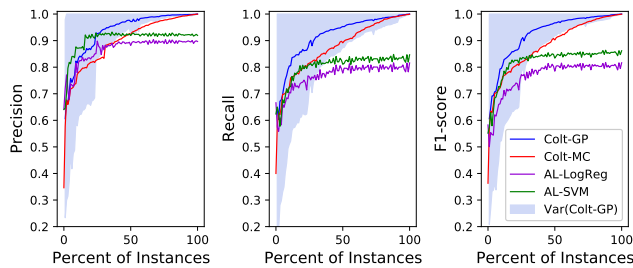


Figure 4: COLT compared to active learning.

While COLT-GP exceeds COLT-MC in all plots, it exhibits a faster increase than COLT-MC. Regarding the F1-score, COLT-GP grows on average 1.54% faster than COLT-MC between 5-8% of used training instances and slows down to 0.85% on 22-25% range. In contrast, COLT-MC increases almost linearly between 12% and 100% F1-score.

Variance in prediction for COLT-GP decreases with the amount of training data, as the model becomes increasingly confident. In terms of F-measure, COLT-GP achieves 100% when all training data is used, proving that the model is expressive enough to gain the necessary knowledge. Using only 5% and 10% of the training data, COLT-GP achieves 75% and 80% F-measure. We measured the biggest improvement over COLT-MC when using 31% of the training data, resulting in an F-measure of 92.8% compared to 80.2% produced by COLT-MC. This represents an improvement of 12.6% in F-measure. Looking at the evaluation of all four approaches, we conclude that it is neither sufficient to solely rely on the similarities of knowledge graph embedding nor to apply simple active learning techniques to solve the problem adequately.

6.4 Q4: Rule confidence estimation

We compare the rule confidence estimation of COLT-GP with the data-driven standard confidence measure (*SCF*) as used by AMIE [13]. It is defined as $\frac{\#KG}{\#total}$ in Table 1. We first determine the correct confidence value for each rule based on its annotated facts and use it to calculate the average estimation error over all rules. The continuous red line (*SCF*) in Figure 5 shows that the average *SCF* estimate is 19.6% away from the correct confidence value. To calculate the estimation error of COLT-GP, we use an increasing number of instances to train a model for each rule and report the average error of these models. For this calculation, we only use rules from Table 1, that produce at least 100 instances. The blue line (COLT-GP)

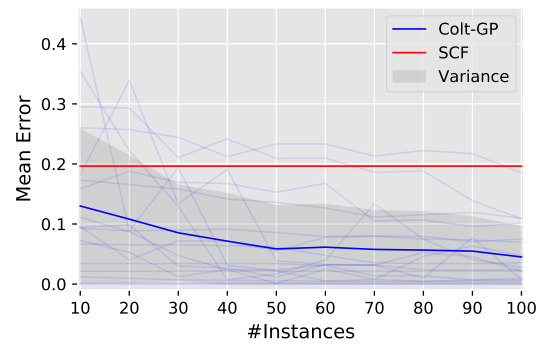


Figure 5: Error reduction in rule confidence estimation using an increasing number of training instances. Light blue lines indicate individual rules error.

in Figure 5 shows that, from the beginning, COLT-GP estimates the average rule confidence more accurately than AMIE’s *SCF* estimate. As a result, 10 to 20 training instances are sufficient to achieve an average estimated confidence error of 12.9 to 10.8 percent, which corresponds to an improvement of 6.7 and 8.8 percentage points over the *SCF* estimate. As the number of training data increases, the estimation error decreases and reaches its lowest value of 4.5 percentage points when using 100 training instances. Figure 5 also shows that, although the estimated error between 50-100 training instances improves by only 1.3 percentage points, the variance decreases by 3.5%, indicating that COLT-GP gains more confidence in its predictions.

7 CONCLUSION & FUTURE WORK

We introduced COLT, a few-shot rule-based knowledge validation framework for interactive quality assessment of rules. Our algorithm, COLT-GP, uses Gaussian processes and deep kernel learning to interactively learn classifiers on rules. COLT-GP benefits from knowledge graph embeddings and user feedback to remedy data quality issues in KGs. We juxtapose COLT-GP with a strong baseline based on maximum coverage. COLT-GP attains 10% error in confidence estimation with only 20 user-validated facts and 75% prediction quality in rule validation with only 5% labeled instances. We release the code and dataset of 26 rules and 23 000 labeled facts.

While this work focuses on learning the characteristics of single rules, we plan to extend our approach to efficiently learn the confidence of multiple rules simultaneously. A promising step towards this goal might be to harness the meta-learning capabilities of neural processes [22]. Furthermore, the COLT framework could generate refined rules on the facts validated by the model on each rule to accomplish a fine-grained maintenance of KGs. We note that COLT’s focus on rule-driven validation does not prevent its adoption to triples returned by other query mechanisms.

REFERENCES

- [1] N. Ahmadi, J. Lee, P. Papotti, and M. Saeed. Explainable fact checking with probabilistic answer set programming. In *Conference for Truth and Trust Online (TTO)*, 2019.
- [2] N. Ahmadi, T.-T.-D. Truong, L.-H.-M. Dao, S. Ortona, and P. Papotti. Rulehub: A public corpus of rules for knowledge graphs. *Journal of Data and Information Quality (JDIQ)*, 12(4):1–22, 2020.

- [3] A. Arioua and A. Bonifati. User-guided repairing of inconsistent knowledge bases. In *Proceedings of the International Conference on Extending Database Technology (EDBT)*, 2018.
- [4] I. Balazevic, C. Allen, and T. M. Hospedales. Hypernetwork knowledge graph embeddings. In *Proceedings of the International Conference on Artificial Neural Networks (ICANN)*, pages 553–565, 2019.
- [5] C. M. Bishop. *Pattern recognition and machine learning*. Springer, 2006.
- [6] A. Bordes, N. Usunier, A. García-Durán, J. Weston, and O. Yakhnenko. Translating embeddings for modeling multi-relational data. In *Proceedings of the International Conference on Neural Information Processing Systems (NIPS)*, pages 2787–2795, 2013.
- [7] R. Cohen and L. Katzir. The generalized maximum coverage problem. *Information Processing Letters*, 108(1):15–22, 2008.
- [8] L. Dehaspe and H. Toivonen. Discovery of frequent datalog patterns. *Data Mining and Knowledge Discovery*, 3(1):7–36, Mar 1999.
- [9] O. Deshpande, D. S. Lamba, M. Tourn, S. Das, S. Subramaniam, A. Rajaraman, V. Harinarayan, and A. Doan. Building, maintaining, and using knowledge bases: a report from the trenches. In *Proceedings of the International Conference on Management of Data (SIGMOD)*, pages 1209–1220, 2013.
- [10] J. Djolonga, A. Krause, and V. Cevher. High-dimensional Gaussian process bandits. In *Proceedings of the International Conference on Neural Information Processing Systems (NIPS)*, pages 1025–1033, 2013.
- [11] X. Dong, E. Gabrilovich, G. Heitz, W. Horn, N. Lao, K. Murphy, T. Strohmann, S. Sun, and W. Zhang. Knowledge vault: a web-scale approach to probabilistic knowledge fusion. In *Proceedings of the International Conference on Knowledge discovery and data mining (SIGKDD)*, pages 601–610, 2014.
- [12] W. Fan, P. Lu, C. Tian, and J. Zhou. Deducing certain fixes to graphs. *PVLDB*, 12(7):752–765, 2019.
- [13] L. Galárraga, C. Teflioudi, K. Hose, and F. M. Suchanek. Fast rule mining in ontological knowledge bases with AMIE+. *Vldb Journal*, 24(6):707–730, 2015.
- [14] L. A. Galárraga, C. Teflioudi, K. Hose, and F. M. Suchanek. AMIE: association rule mining under incomplete evidence in ontological knowledge bases. In *Proceedings of the International World Wide Web Conference (WWW)*, pages 413–422, 2013.
- [15] G. A. Gesese, R. Biswas, and H. Sack. A comprehensive survey of knowledge graph embeddings with literals: Techniques and applications. In *Proceedings of the Workshop on Deep Learning for Knowledge Graphs (DL4KG)*, pages 31–40, 2019.
- [16] S. Guo, Q. Wang, L. Wang, B. Wang, and L. Guo. Jointly embedding knowledge graphs and logical rules. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 192–202, 2016.
- [17] J. He, E. Veltri, D. Santoro, G. Li, G. Mecca, P. Papotti, and N. Tang. Interactive and deterministic data cleaning. In *Proceedings of the International Conference on Management of Data (SIGMOD)*, pages 893–907, 2016.
- [18] K. He, X. Zhang, S. Ren, and J. Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 1026–1034, 2015.
- [19] A. Heidari, J. McGrath, I. F. Ilyas, and T. Rekatsinas. HoloDetect: Few-shot learning for error detection. In *Proceedings of the International Conference on Management of Data (SIGMOD)*, pages 829–846, 2019.
- [20] V. T. Ho, D. Stepanova, M. H. Gad-Elrab, E. Kharlamov, and G. Weikum. Rule learning from knowledge graphs guided by embedding models. In *Proceedings of the International Semantic Web Conference (ISWC)*, pages 72–90, 2018.
- [21] G. Ji, S. He, L. Xu, K. Liu, and J. Zhao. Knowledge graph embedding via dynamic mapping matrix. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 687–696, 2015.
- [22] H. Kim, A. Mnih, J. Schwarz, M. Garnelo, S. M. A. Eslami, D. Rosenbaum, O. Vinyals, and Y. W. Teh. Attentive neural processes. In *International Conference on Learning Representations (ICLR)*, 2019.
- [23] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015.
- [24] N. Klyuchnikov, D. Mottin, G. Koutrika, E. Müller, and P. Karras. Figuring out the user in a few steps: Bayesian multifidelity active search with cokriging. In *SIGKDD*, pages 686–695, 2019.
- [25] J. Lajus, L. Galárraga, and F. M. Suchanek. Fast and exact rule mining with AMIE 3. In *Proceedings of the European Semantic Web Conference (ESWC)*, volume 12123, pages 36–52, 2020.
- [26] D. D. Lewis and W. A. Gale. A sequential algorithm for training text classifiers. In *Proceedings of the International Conference on Information retrieval (SIGIR)*, pages 3–12, 1994.
- [27] D. Liben-Nowell and J. Kleinberg. The link-prediction problem for social networks. *Journal of the American society for information science and technology*, 58(7):1019–1031, 2007.
- [28] M. Lissandrini, D. Mottin, T. Palpanas, and Y. Velegrakis. Data exploration using example-based methods. *Synthesis Lectures on Data Management*, 10(4):1–164, 2018.
- [29] M. Lissandrini, D. Mottin, D. Papadimitriou, T. Palpanas, and Y. Velegrakis. Unleashing the power of information graphs. *SIGMOD Record*, 43(4), 2014.
- [30] M. Mahdavi, Z. Abedjan, R. C. Fernández, S. Madden, M. Ouzzani, M. Stonebraker, and N. Tang. Raha: A configuration-free error detection system. In *Proceedings of the International Conference on Management of Data (SIGMOD)*, pages 865–882, 2019.
- [31] G. L. Nemhauser, L. A. Wolsey, and M. L. Fisher. An analysis of approximations for maximizing submodular set functions-I. *Mathematical programming*, 14(1):265–294, 1978.
- [32] M. Nickel, L. Rosasco, and T. A. Poggio. Holographic embeddings of knowledge graphs. In *Proceedings of the National Conference on Artificial Intelligence (AAAI)*, pages 1955–1961, 2016.
- [33] M. Nickel, V. Tresp, and H. Kriegl. A three-way model for collective learning on multi-relational data. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 809–816, 2011.
- [34] S. Ortona, V. V. Meduri, and P. Papotti. Robust discovery of positive and negative rules in knowledge bases. In *Proceedings of the International Conference on Data Engineering (ICDE)*, pages 1168–1179, 2018.
- [35] J. Pujara, E. Augustine, and L. Getoor. Sparsity and noise: Where knowledge graph embeddings fall short. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1751–1756, 2017.
- [36] N. Roy and A. McCallum. Toward optimal active learning through sampling estimation of error reduction. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 441–448, 2001.
- [37] G. Schohn and D. Cohn. Less is more: Active learning with support vector machines. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 839–846, 2000.
- [38] B. Settles. *Active Learning*. Synthesis Lectures on Artificial Intelligence and Machine Learning. Morgan & Claypool Publishers, 2012.
- [39] C. Shi, Y. Li, J. Zhang, Y. Sun, and S. Y. Philip. A survey of heterogeneous information network analysis. *IEEE Transactions on Knowledge and Data Engineering (TKDE)*, 29(1):17–37, 2016.
- [40] J. Shin, S. Wu, F. Wang, C. D. Sa, C. Zhang, and C. Ré. Incremental knowledge base construction using DeepDive. *PVLDB*, 8(11):1310–1321, 2015.
- [41] N. Srinivas, A. Krause, S. Kakade, and M. W. Seeger. Gaussian process optimization in the bandit setting: No regret and experimental design. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 1015–1022, 2010.
- [42] F. M. Suchanek, G. Kasneci, and G. Weikum. Yago: a core of semantic knowledge. In *Proceedings of the International World Wide Web Conference (WWW)*, pages 697–706, 2007.
- [43] H. P. Vanchinathan, A. Marfurt, C.-A. Robelin, D. Kossmann, and A. Krause. Discovering valuable items from massive data. In *Proceedings of the International Conference on Knowledge discovery and data mining (SIGKDD)*, pages 1195–1204, 2015.
- [44] Q. Wang, Z. Mao, B. Wang, and L. Guo. Knowledge graph embedding: A survey of approaches and applications. *IEEE Transactions on Knowledge and Data Engineering (TKDE)*, pages 2724–2743, 2017.
- [45] Z. Wang, J. Zhang, J. Feng, and Z. Chen. Knowledge graph embedding by translating on hyperplanes. In *Proceedings of the National Conference on Artificial Intelligence (AAAI)*, pages 1112–1119, 2014.
- [46] A. G. Wilson, Z. Hu, R. Salakhutdinov, and E. P. Xing. Deep kernel learning. In *Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 370–378, 2016.
- [47] A. G. Wilson, Z. Hu, R. Salakhutdinov, and E. P. Xing. Stochastic variational deep kernel learning. In *Proceedings of the International Conference on Neural Information Processing Systems (NIPS)*, pages 2586–2594, 2016.
- [48] R. Xie, Z. Liu, J. Jia, H. Luan, and M. Sun. Representation learning of knowledge graphs with entity descriptions. In *Proceedings of the National Conference on Artificial Intelligence (AAAI)*, pages 2659–2665, 2016.
- [49] P. Xu and D. Barbosa. Connecting language and knowledge with heterogeneous representations for neural relation extraction. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, pages 3201–3206, 2019.
- [50] W. Zhang, B. Paudel, L. Wang, J. Chen, H. Zhu, W. Zhang, A. Bernstein, and H. Chen. Iteratively learning embeddings and rules for knowledge graph reasoning. In *Proceedings of the International World Wide Web Conference (WWW)*, pages 2366–2377, 2019.