

# Discovering Top-k Interesting Rules

Paper ID: XXX

## ABSTRACT

This paper studies two questions about discovery of data quality rules. What rules are interesting? How can we discover top-ranked interesting rules? We consider *entity enhancing rules*, denoted by REEs, which can embed machine learning classifiers as predicates in logic rules, and subsume conditional functional dependencies, denial constraints and matching dependencies as special cases. We propose an interestingness model in terms of both objective measures such as support, confidence and diversity, and subject measures for user preference and interests; we show how to learn the subjective model and weight vectors via active learning. Based on the interestingness model, we develop a top- $k$  algorithm to discover interesting REEs, and an any-time algorithm for successive interesting-rule discovery via lazy evaluation. We parallelize these algorithms such that they guarantee to reduce runtime when more processors are used. Using real-life and synthetic datasets, we show that the algorithms are able to find interesting rules and speed up conventional rule-discovery methods by 8.4 times on average.

## 1 INTRODUCTION

Poor data quality has been a longstanding challenge for decades. Real-life data is often dirty, *e.g.*, attributes of the same entity may carry inconsistent or conflict values, different entities may be mismatched into one, and the same entity bears different representations and is treated as distinct ones. Dirty data is costly, *e.g.*, Gartner reports that poor data quality is responsible for an average of \$15 million per year in losses for organizations [27], and IBM estimates that dirty data costs the US \$3.1 trillion in 2016 alone (cf. [58]).

To improve data quality, a variety of logic rules have been studied for *entity resolution* (ER) and *conflict resolution* (CR), *e.g.*, matching dependencies (MDs) [4, 6, 16] for ER, and denial constraints (DCs) [3] and conditional functional dependencies (CFDs) [17, 19] for CR. Here ER is to identify tuples that refer to the same real-world entity, and CR is to resolve inconsistencies pertaining to the entities. There has also been recent effort to embed machine learning (ML) classifiers as predicates in logic rules [21], to unify ML-based and rule-based methods. When provided with high-quality rules, one can detect value conflicts and entity mismatches, fix the errors, and interpret the reasons behind the errors and their corrections.

To make practical use of the rule-based methods, it is a must to be able to discover interesting rules from real-life datasets. In the big data era, it is unrealistic to solely rely on domain experts to write or find rules manually. Several algorithms have been developed for discovering functional dependencies (FDs) [5, 34, 37, 52, 74, 78], DCs [12, 46, 54], MDs [65, 69] and CFDs [10, 18]. However, a major problem encountered in practice is that rule discovery often yields an excessive number of rules. For example, on a small dataset with only 27 attributes and 368 tuples, there are 128,726 FDs found [51]. The practitioners are often overwhelmed by the large number of rules and have to spend a significant amount of time to manually inspect and select rules that are interesting to them. This staggering cost hampers the applicability of the rule-based methods.

To tackle this issue, two questions have to be answered.

(1) *Interestingness*. What rules are interesting? Conventional measures include, *e.g.*, support and confidence, for how often the rules can be applied and how strong the associations between their preconditions and consequences are. In practice, however, these “universal” objective measures often do not suffice. In our experience, practitioners in industry often have accumulated “business rules” from their practice; they already know quite a few rules that have high support and confidence. What they want from rule discovery are rules that are “surprising” to them, *i.e.*, sensible rules that they do not yet know, to complement those they have already got. This suggests that we consider “subjective measures” in the interestingness model, to capture the preference and interests of different users.

(2) *Top-k algorithms*. Assume that the interestingness model is in place. The next question is how to discover top-ranked rules based on the model? A brute-force approach is to first find all the rules that hold on a dataset, and then sort the rules and return top- $k$  ones for a predefined constant  $k$ . Can we do better? When we use a search engine, we do not want it to first enumerate all matches of our keywords and then return the top-ranked ones. Analogous to search, users of a rule discovery algorithm typically want to efficiently find top- $k$  rules, and if these do not meet their need, continue to find the next top- $k$  ones. To our knowledge, however, no such discovery algorithm is yet available for data quality rules such as functional dependencies (FDs), DCs, CFDs and MDs.

**Contributions & organizations.** This paper tackles these questions. As a testbed of our proposed methods, we consider entity enhancing rules (REEs) [21], since REEs (a) embed ML classifiers in logic rules, (b) unify ER and CR with the same rules, (c) are collectively defined across multiple relations, and (d) subsume DCs, CFDs and MDs as special cases. We will review REEs in Section 2.

(1) *Interestingness model* (Section 3). We propose a model to characterize the interestingness of REEs, in terms of both conventional objective measures and new subjective measures for users preference and interests, bearing various weights. We also propose an active learning approach to determine the subjective model and weight vectors. In addition, we show that the conventional notion of support does not fit collective rules such as REEs, *i.e.*, it is no longer anti-monotonic; in light of this, we introduce a revision of the notion for REEs and show its anti-monotonicity.

(2) *Top-k algorithm* (Section 4). Based on the model, we formulate top- $k$  discovery problem for REEs, and develop a top- $k$  algorithm as a solution. The algorithm maintains an interestingness upper bound and adopts two effective processing orders to speed up the discovery. We also show how to accommodate the collective nature of REEs and its embedded ML predicates in the process.

(3) *Anytime algorithm* (Section 5). Extending the top- $k$  algorithm, we provide another algorithm that allows users to continue to find the next top- $k$  REEs when needed, via lazy evaluation.

(4) *Parallelization* (Section 6). To scale with large datasets, we parallelize the top- $k$  algorithm and the anytime algorithm using a cluster of machines. We show that the algorithms are parallelly

scalable [39], *i.e.*, they guarantee to reduce the parallel runtime when more machines are used. Hence in theory, the algorithms are able to efficiently discover REEs even when the datasets grow big, as long as we add more computing resources when needed.

(5) *Experimental study* (Section 7). Using real-life and synthetic datasets, we empirically find the following. On average, (a) top- $k$  REEs discovery speeds up conventional methods by 8.4 times. It takes 667s on DBLP with 18 attributes and 1,799,559 tuples for  $k = 10$  with 20 machines, versus 8726s by traditional methods. (b) The lazy evaluation strategy makes the anytime algorithm 4.85 times faster than the top- $k$  one when users want the 4th top-10 REEs. (c) Our interestingness model is 25.5% more accurate than the state-of-the-art language models, and its subjective measures improves the accuracy from 0.48 to 0.76. (d) The algorithms are parallel scalable; they are 3.6 times faster using 20 machines instead of 4.

Since REEs subsume DCs, FDs, CFDs and MDs, our algorithms can also be used to discover such top-ranked rules. For discovering such rules, this work is among the first efforts to train and employ subjective models, and provides the first top- $k$  and anytime algorithms, to the best of our knowledge. These discovery algorithms are also among the first with the parallel scalability for such rules.

**Related work.** The related work is categorized as follows.

*Interestingness models.* The need for considering both objective measures and subjective measures has long been recognized in the data mining community, to reflect users’ universal and individualistic preference, respectively (see [29] for a survey). However, while many rule-mining algorithms rank rules using some particular objective measures (*e.g.*, [12, 54]), we are not aware of any discovery algorithms for FDs, DCs, CFDs and MDs that take subjective measures into account. For instance, while MDedup [36] employs several objective features for MDs in feature engineering and adopts the ML regression model Gaussian Process to learn a score for each MD, it aims to discover MDs with high F-measures, not for users’ surprisingness or unexpectedness. While ML models [14, 67] were designed to learn the representations of rules, they do not consider how to learn rule interestingness. Recently, RNNLogic [56] proposes to learn logic rules for graphs and select high-quality ones, but it is also based on objective features such as support and confidence.

This work differs from the prior work as follows. (1) We propose an interestingness model designated for data quality rules that combines both objective and subjective measures, the first of the kind for REEs, DCs, CFDs and MDs. (2) We propose to use active learning and pairwise ranking strategies to learn the interestingness of rules, to catch the users’ background knowledge.

*Rule discovery.* A number of discovery algorithms have been developed for data quality rules, classified as follows. (1) *Levelwise search.* TANE [34], FUN [50], FD\_mine [76] discover FDs based on a lattice structure. Depmine [47] develops agree sets to accelerate the search, HyFD [52] combines sampling techniques and lattice traversal to prune non-FDs, and DynFD [66] mines FDs in dynamic datasets. The levelwise methods have also been extended to CFDs discovery, including CTANE [18] and tableau generation [31]. For ER rules, [69] mines thresholds of similarity functions for discovering MDs. SMFD [28] proposes a top-down based framework to discover and

validate FDs in a secured multi-party scenario. (2) *Depth-first search.* DFD [1] adopts the depth-first search (random walk) in the lattice structure. FastFDs [74] improves Depmine by using different sets to mine FDs. For CFDs and DCs, depth-first search approaches also apply, *e.g.*, FastCFDs [18] extends FastFDs to mine CFDs, while FastDC [12], Hydra [8], DCFinder [54] and ADCMiner [46] find valid DCs based on evidence sets and min-cover algorithms. Although DCs support multiple relation atoms [12] by definition, all its discovery algorithms [8, 46, 54] restrict their settings to bi-variable DCs discovery only. (3) *Hybrid approaches.* HyMD [65] makes use of both levelwise and depth-first search, and discovers MDs with pre-computed index structures. MDedup [36] approaches ER by selecting useful MDs discovered by HyMD. (4) *Learning-based approaches.* [23] utilizes inductive learning to find FDs. AutoFD [78] adopts structure learning (*i.e.*, graphical lasso) to discover FDs. [68] and [35] adopt the rule learning strategies to find MDs with similarity functions (see [11] for more ER solvers). (5) *Top- $k$  discovery.* Closer to this work are top- $k$  algorithms for discovering association rules in relations [72] and graphs [22].

This work differs from the prior work in the following. (1) We study top- $k$  rule discovery based on a new interestingness model, not to find all rules on a dataset subject to *e.g.*, support. It differs from [22, 72] in the use of different ranking criteria (interestingness) and thus different early-termination strategies. As remarked earlier, even the notion of support is different for REEs. (2) We extend our top- $k$  algorithm to an anytime algorithm, in order for users to find the next top- $k$  results when needed, without re-discovering starting from scratch. No existing work has considered anytime rule discovery. (3) It is more challenging to discover REEs than DCs, CFDs and MDs, since REEs can embed ML predicates, and are collectively defined across multiple relations. For instance, the traditional notion of support is no longer anti-monotonic for REEs.

*ML models.* Various ML models have been adopted for ER, CR, link prediction, and checking semantic similarity. (1) *Models for ER,* via SVM and logistic regression (see [30] for a survey on traditional ML models for ER), unsupervised learning, *e.g.*, ZeroER [73], and deep learning, *e.g.*, Ditto [44], DeepMatcher [49], DeepER [15], AutoEM [80], GraphER [41], MPM [25]. (2) *Models for CR,* including HoloClean [60], HoloDetect [32], Raha [48] and SLiMFast [62], for error detection, error correction and data fusion. (3) *Models for similarity checking,* by employing popular language models such as BERT-based models [13, 40, 45, 59], XLNet [75] and GPT [9, 57].

This work is not to develop another ML model. Instead, we show how to leverage existing well-trained ML models by embedding them as ML predicates in REEs. As will be seen in Example 1, one can even plug in a link prediction ML models into REEs to reveal the “hidden” associations between entities, to facilitate ER and CR.

*Parallel rule discovery.* Parallel algorithms have been developed for rule discovery. [26, 42] employ massive parallelism to discover FDs, but they do not consider communication cost. [43] discovers FDs in a distributed setting but only returns local FDs. [63] extends FastFD by minimizing communication cost. [64] proposes a distributed framework for discovering both FDs and DCs. SMFD [28] performs FD discovery based on the lattice structure in a distributed manner, but it only focuses on how to enforce privacy constraints.

tid	pid	address	zipcode	#beds	#baths	selling_price(\$)
$t_1$	$p_1$	67 old head RD, Sydney	2029	2	2	0.85 M
$t_2$	$p_2$	NO.67 old head road, SYD	2029	2	2	0.85 M
$t_3$	$p_3$	15 Pitt street, Pitt point	2022	5	3	1.50 M
$t_4$	$p_4$	15 Pitt street	2022	5	3	1.20 M

Table 1: Example Property relation  $D_1$

This work is among the first discovery algorithms for relational data quality rules that guarantee the parallel scalability, *i.e.*, the execution time is guaranteed to be reduced if more machines are used, when both computational and communication costs are considered.

## 2 ENTITY ENHANCING RULES

We next review entity enhancing rules (REEs) defined in [21].

We define REEs over a database schema  $\mathcal{R} = (R_1, \dots, R_m)$ . Each  $R_j$  is a schema of the form  $R(A_1 : \tau_1, \dots, A_n : \tau_n)$ , where  $A_i$  is an attribute of type  $\tau_i$ . A relation  $D$  of  $R$  is a set of tuples having attributes  $A_i$  of  $R$  ( $i \in [1, n]$ ) with values from the domains of  $\tau_i$ . We assume *w.l.o.g.* that each tuple  $t$  in  $D$  has an id attribute, which uniquely identifies the entity that  $t$  represents. An instance  $\mathcal{D}$  of  $\mathcal{R}$  is a collection  $(D_1, \dots, D_m)$ , where  $D_i$  is a relation of  $R_i$  ( $i \in [1, m]$ ).

**Predicates.** Following tuple relational calculus (see, *e.g.*, [2]), we define *predicates* over schema  $\mathcal{R}$  as follows:

$$p ::= R(t) \mid t.A \oplus c \mid t.A \oplus s.B \mid \mathcal{M}(t[\vec{A}], s[\vec{B}]),$$

where  $\oplus$  is an operator in  $\{=, \neq\}$ , and  $c$  is a constant. Here (a)  $R(t)$  says that  $t$  is a tuple of  $R$ , for  $R \in \mathcal{R}$ ; we refer to  $R(t)$  as a *relation atom* over  $\mathcal{R}$ , and to  $t$  as a tuple variable *bounded* by  $R(t)$ ; (b)  $t.A$  denotes an attribute of  $t$  when  $t$  is bounded by  $R(t)$  and  $A$  is an attribute in  $R$ ; (c)  $t.A \oplus c$  is a *constant predicate* when  $c$  is a value in the domain of  $A$ ; (d)  $t.A \oplus s.B$  compares *compatible* attributes  $t.A$  and  $s.B$ , *i.e.*, tuple  $t$  (resp.  $s$ ) is bounded by  $R(t)$  (resp.  $R'(s)$ ), and  $A \in R$  and  $B \in R'$  have the same type; (e)  $\mathcal{M}(t[\vec{A}], s[\vec{B}])$  is an *ML predicate* where  $t[\vec{A}]$  and  $s[\vec{B}]$  are vectors of pairwise compatible attributes. Here  $\mathcal{M}$  can be any ML classifiers (see below), which returns true if it determines that  $t[\vec{A}]$  and  $s[\vec{B}]$  are correlated, and false otherwise.

**REEs.** An *entity enhancing rule* (REE)  $\varphi$  over  $\mathcal{R}$  is defined as

$$\varphi : X \rightarrow p_0,$$

where  $X$  is a conjunction of *predicates* over  $\mathcal{R}$ , and  $p_0$  is a predicate over  $\mathcal{R}$  such that all tuple variables in  $\varphi$  are bounded in  $X$ . We refer to  $X$  as the *precondition* of  $\varphi$  and  $p_0$  as the *consequence* of  $\varphi$ .

**Example 1:** Consider the following example from a (simplified) realty database with self-explained schemas Property (pid, address, zipcode, #beds, #baths, selling\_price) and RentalContract (cid, pid, renter, tenant, monthly\_rent, rent\_start\_year, rent\_end\_year).

Below are example REEs over the database schema.

(1)  $\varphi_1$  :  $\text{Property}(t_a) \wedge \text{Property}(t_b) \wedge \mathcal{M}_{\text{addr}}(t_a.\text{address}, t_b.\text{address}) \wedge \text{RentalContract}(s_a) \wedge \text{RentalContract}(s_b) \wedge t_a.\text{pid} = s_a.\text{pid} \wedge t_b.\text{pid} = s_b.\text{pid} \wedge s_a.\text{rent\_start\_year} = s_b.\text{rent\_start\_year} \wedge s_a.\text{rent\_end\_year} = s_b.\text{rent\_end\_year} \wedge \mathcal{M}_{\text{name}}(s_a[\vec{A}], s_b[\vec{A}]) \rightarrow t_a.\text{pid} = t_b.\text{pid}$ . Here  $\mathcal{M}_{\text{addr}}$  is an ML model for checking the closeness of addresses,  $\mathcal{M}_{\text{name}}$  checks the semantic similarity on names, and  $\vec{A}$  denotes (renter, tenant) in relation RentalContract.

Intuitively,  $\varphi_1$  says that two properties can be identified if (a) their addresses match, (b) renters and tenants are “similar”, and (c) rental periods are the same. It uses ML models  $\mathcal{M}_{\text{addr}}$  and  $\mathcal{M}_{\text{name}}$  to check semantic closeness of text attributes, *e.g.*, “Michael” vs. “Mike”.

(2)  $\varphi_2$ :  $\text{Property}(t_a) \wedge \text{RentalContract}(s_a) \wedge \text{RentalContract}(s_b) \wedge$

tid	cid	pid	renter	tenant	monthly_rent	rent_start_year	rent_end_year
$t_5$	$c_1$	$p_1$	Michael	Leo	560	2016	2020
$t_6$	$c_2$	$p_2$	Mich.	Leo Z.	560	2016	2020
$t_7$	$c_3$	$p_3$	Ada	Mike	1,300	2014	2016
$t_8$	$c_4$	$p_3$	Ada	Connie	1,500	2016	2018
$t_9$	$c_5$	$p_4$	Ada	Bob	1,500	2018	2020

Table 2: Example RentalContract relation  $D_2$

$\mathcal{M}_{\text{rent}}(t_a, s_a) \wedge s_a.\text{rent\_start\_year} \neq s_b.\text{rent\_start\_year} \wedge t_a.\text{pid} = s_b.\text{pid} \rightarrow s_a.\text{monthly\_rent} \neq s_b.\text{monthly\_rent}$ . Here we use a link prediction model  $\mathcal{M}_{\text{rent}}(t_a, s_a)$  to predict whether the tenant in  $s_a$  may rent property  $t_a$  (note that RentalContract might not record the fact “ $s_a$  rents  $t_a$ ” yet). Intuitively, it says that rent fluctuates at different times, since  $t_a$  is the same property as in  $s_b$ .

(3)  $\varphi_3$  :  $\text{Property}(t_a) \wedge \text{Property}(t_b) \wedge X \rightarrow \mathcal{M}_{\text{addr}}(t_a.\text{address}, t_b.\text{address})$ , where  $\mathcal{M}_{\text{addr}}$  is the address ML model used in  $\varphi_1$ ,  $X = \bigcap_{A_s \in \mathcal{T}} t_a.A_s = t_b.A_s$  and  $\mathcal{T}$  denotes a designated set of attributes in Property (not shown in the simplified schema), including built-up area, floor, and neighborhood information. Here conditions in  $X$  interpret the prediction of  $\mathcal{M}_{\text{addr}}(t_a.\text{address}, t_b.\text{address})$  in logic.  $\square$

As shown in Example [21], REEs subsume as special cases MDs for ER (by using ML models to check similarity), and FDs, CFDs and DCs for CR. They may carry atoms  $R(t)$  from multiple relations  $R$  for, *e.g.*, collective ER [7]. Moreover, REEs may embed state-of-the-art ML classifiers, such as (1) NLP models, *e.g.*, Bert [13], for text classification; (2) ER models and link prediction models, *e.g.*, ditto [44] and DeepMatcher [49], to reveal “hidden” associations between tuples across relations, which are not connected by, *e.g.*, keys and foreign keys; and (3) models for data fusion, error detection and correction, *e.g.*, HoloClean [60] and HoloDetect [32] when they are taken as classifiers to return Boolean values.

**Semantics.** Consider an instance  $\mathcal{D}$  of  $\mathcal{R}$ . A *valuation*  $h$  of tuple variables of  $\varphi$  in  $\mathcal{D}$ , or simply a valuation of  $\varphi$ , is a mapping that instantiates  $t$  in each  $R(t)$  with a tuple in a relation  $D$  of  $\mathcal{D}$ .

We say that  $h$  *satisfies* a predicate  $p$ , written as  $h \models p$ , if the following are satisfied: (1) If  $p$  is a relation atom  $R(t)$ ,  $t \oplus c$  or  $t.A \oplus s.B$ , then  $h \models p$  is interpreted as in tuple relational calculus following the standard semantics of first-order logic [2]. (2) If  $p$  is  $\mathcal{M}(t[\vec{A}], s[\vec{B}])$ , then  $h \models p$  if  $\mathcal{M}$  predicts true on  $(h(t)[\vec{A}], h(s)[\vec{B}])$ .

Given a conjunction  $X$  of predicates, we say  $h \models X$  if for *all* predicates  $p$  in  $X$ ,  $h \models p$ . Given an REE  $\varphi$ , we write  $h \models \varphi$  such that if  $h \models X$ , then  $h \models p_0$ . An instance  $\mathcal{D}$  of  $\mathcal{R}$  *satisfies*  $\varphi$ , denoted by  $\mathcal{D} \models \varphi$ , if for *all* valuations  $h$  of tuple variables of  $\varphi$  in  $\mathcal{D}$ ,  $h \models \varphi$ . That is, REEs have a universal semantics. We say that  $\mathcal{D}$  *satisfies* a set  $\Sigma$  of REEs, denoted by  $\mathcal{D} \models \Sigma$ , if for all  $\varphi \in \Sigma$ ,  $\mathcal{D} \models \varphi$ .

**Example 2:** Continuing with Example 1, assume that  $\mathcal{D}$  consists of two relations  $D_1$  and  $D_2$  of schemas Property and RentalContract, shown in Tables 1 and 2, respectively. Consider valuation  $h_2$ :  $t_3 \mapsto t_a$ ,  $t_8 \mapsto s_a$  and  $t_7 \mapsto s_b$ . It satisfies REE  $\varphi_2$ , since Mike rented property  $p_3$  earlier than Connie and thus, he pays a cheaper rent.

As another example, consider the valuation  $h_1$  of REE  $\varphi_1$  that has mappings:  $t_1 \mapsto t_a$ ,  $t_2 \mapsto t_b$ ,  $t_5 \mapsto s_a$  and  $t_6 \mapsto s_b$ . It helps us identify that  $t_1$  and  $t_2$  are indeed the same property.  $\square$

## 3 LEARNING INTERESTINGNESS

In this section we study the interestingness of REEs. We first present a set of criteria for making REEs interesting, including objective measures and subjective measures (Section 3.1). We then propose

our interesting model (Section 3.2). Finally, we propose an active-learning approach to learn the subjective model and the weights of the measures in the interestingness model (Section 3.3).

### 3.1 Interestingness Measures

To find truly useful REEs for users, we consider (a) objective measures, which are based only on the datasets and are “universal” to different users; and (b) subjective measures, which are based on both the data and the users, including the users’ background, preference and interests, and may vary for different groups of users.

**3.1.1 Objective Measures.** We start with objective measures, including *support*, *confidence*, *attribute diversity*, *minimality* and *succinctness*. Note that other objective measures can be plugged in our interestingness model when needed. To simplify the discussion we just name a few that are commonly used in practice.

**Support.** Support measures how frequent an REE can be applied. For collective rules that are defined across multiple relations such as REEs, the conventional notion of support has to be revised.

To see this, we first define an order on REEs. Given two REEs  $\varphi : X \rightarrow p_0$  and  $\varphi' : X' \rightarrow p_0$  that have the same consequence  $p_0$ , we say that  $\varphi$  has a *lower order* than  $\varphi'$ , denoted by  $\varphi \leq \varphi'$ , if  $X \subseteq X'$ . Intuitively,  $\varphi$  is less restrictive than  $\varphi'$ .

One might think that given a dataset  $\mathcal{D}$  of database schema  $\mathcal{R}$ , support for REEs  $\varphi : X \rightarrow p_0$  over  $\mathcal{R}$  is simply the number of distinct valuations  $h$  of  $\varphi$  in  $\mathcal{D}$  such that  $h \models X$ . This is the conventional definition of support commonly used in data quality rules on a *single* relation, e.g., FDs, CFDs, etc. It satisfies the *anti-monotonicity* on single relations, i.e., if  $\varphi \leq \varphi'$ , then the support of  $\varphi$  is at least that of  $\varphi'$ . Unfortunately, for *collective* rules involving *multiple relations*, this definition does not work. To see why it is the case, consider the following from the realty database in Example 1.

**Example 3:** Let  $X$  be  $\text{Property}(t_a) \wedge t_a.\text{zipcode} = 2022$  and  $p_0$  be  $t_a.\text{selling\_price} = 1.50M$ . Consider two REEs  $\varphi$  and  $\varphi'$ ,  $\varphi : X \rightarrow p_0$  and  $\varphi' : X' \rightarrow p_0$ , where  $X' = X \wedge \text{RentalContract}(s_a) \wedge t_a.\text{pid} = s_a.\text{pid}$ . Clearly,  $\varphi \leq \varphi'$  since  $X \subseteq X'$ . However, if the conventional support definition is applied, then the support of  $\varphi$  is 1 by  $t_3 \mapsto t_a$ , while the support of  $\varphi'$  is 2, since both  $(t_3, t_7) \mapsto (t_a, s_a)$  and  $(t_3, t_8) \mapsto (t_a, s_a)$ , violating the anti-monotonicity. Intuitively, this is because in collective rules, a tuple in a relation can join with multiple tuples in another relation, yielding a larger “support”.  $\square$

To fix this, we revise the notion of support and establish its anti-monotonicity, a property that is critical to reducing the search space when discovering rules. To the best of our knowledge, this is among the first anti-monotonicity for collective data quality rules.

For the ease of illustration, we assume predicates in this section involve two tuple variables, i.e.,  $t.A \oplus s.B$  or  $\mathcal{M}(t[\bar{A}], s[\bar{B}])$ , while all notations extend naturally to predicates with one tuple variable.

We use the following notions. Given a predicate  $p$ , we define an REE  $\varphi_p$  to verify whether two tuples satisfy  $p$ :  $R(t) \wedge R'(s) \rightarrow p$ , where  $t$  and  $s$  (of relation schema  $R$  and  $R'$ , respectively) are the tuple variables used in  $p$ . Let  $H_p$  be the set of valuations of  $\varphi_p$  in  $\mathcal{D}$ . We define the *support set* of  $p$  on  $\mathcal{D}$ , denoted by  $\text{spset}(p, \mathcal{D})$ , as

$$\text{spset}(p, \mathcal{D}) = \{\langle h(t), h(s) \rangle \mid h \in H_p \wedge h \models \varphi_p\},$$

i.e., the set of tuple pairs satisfying  $p$ . Similarly, given a conjunction

$X$  of predicates, we define the *support set* of  $X$  as follows:

$$\text{spset}(X, \mathcal{D}) = \{\langle h(t), h(s) \rangle \mid \forall p \in X (\langle h(t), h(s) \rangle \in \text{spset}(p, \mathcal{D}))\},$$

i.e., the set of all tuple pairs satisfying *all* predicates in  $X$ .

Given  $\varphi : X \rightarrow p_0$ , assume that  $H$  is the set of all valuations of  $\varphi$  in  $\mathcal{D}$ , and  $t_0$  and  $s_0$  are the tuple variables used in  $p_0$ . Then the *support set* of  $\varphi$ , denoted by  $\text{spset}(\varphi, \mathcal{D})$ , is defined as

$$\text{spset}(\varphi, \mathcal{D}) = \{\langle h(t_0), h(s_0) \rangle \mid h \in H \wedge h \models X \wedge h \models \varphi\}.$$

To quantify the frequency of  $\varphi$ , we define the *support* of  $\varphi$  as

$$\text{supp}(\varphi, \mathcal{D}) = |\text{spset}(\varphi, \mathcal{D})|.$$

Similarly we define the notions of  $\text{supp}(p, \mathcal{D})$  and  $\text{supp}(X, \mathcal{D})$ .

For an integer  $\sigma$ , an REE is  $\sigma$ -frequent on  $\mathcal{D}$  if  $\text{supp}(\varphi, \mathcal{D}) \geq \sigma$ .

**Theorem 1:** For any instance  $\mathcal{D}$  of  $\mathcal{R}$  and REEs  $\varphi$  and  $\varphi'$ , if  $\varphi \leq \varphi'$ , then  $\text{spset}(\varphi', \mathcal{D}) \subseteq \text{spset}(\varphi, \mathcal{D})$  and  $\text{supp}(\varphi', \mathcal{D}) \leq \text{supp}(\varphi, \mathcal{D})$ .  $\square$

**Proof.** There are two cases to consider: (1)  $\varphi$  and  $\varphi'$  use the same set of tuple variables. In this case,  $\text{spset}(\varphi', \mathcal{D})$  is clearly a subset of  $\text{spset}(\varphi, \mathcal{D})$ , since the predicates that those tuple variables have to satisfy in  $\varphi$  make a subset of those in  $\varphi'$ , and hence more valuations can contribute to the support of  $\varphi$ . (2) REE  $\varphi'$  uses more tuple variables than  $\varphi$ . By the definition of support, the additional tuple variables used in  $\varphi'$  will not increase the support. Indeed, for each  $\langle h(t_0), h(s_0) \rangle$  in  $\text{spset}(\varphi', \mathcal{D})$ ,  $\langle h(t_0), h(s_0) \rangle$  must also be in  $\text{spset}(\varphi, \mathcal{D})$ , since otherwise  $h$  cannot satisfy  $\varphi'$ . In both cases,  $\text{spset}(\varphi', \mathcal{D}) \subseteq \text{spset}(\varphi, \mathcal{D})$  and thus  $\text{supp}(\varphi', \mathcal{D}) \leq \text{supp}(\varphi, \mathcal{D})$ .  $\square$

**Example 4:** In Example 3,  $\varphi \leq \varphi'$ . By Theorem 1,  $\text{supp}(\varphi, \mathcal{D}) \geq \text{supp}(\varphi', \mathcal{D})$ , since  $\text{spset}(\varphi, \mathcal{D}) = \text{spset}(\varphi', \mathcal{D}) = \{t_3 \mapsto t_a\}$ .  $\square$

**Confidence.** Confidence indicates how often an REE  $\varphi : X \rightarrow p_0$  has been found true, given that  $X$  is satisfied.

For an REE  $\varphi : X \rightarrow p_0$ , the *confidence* of  $\varphi$  on  $\mathcal{D}$  denoted by  $\text{conf}(\varphi, \mathcal{D})$ , is defined to be  $\text{conf}(\varphi, \mathcal{D}) = \frac{\text{supp}(X \wedge p_0, \mathcal{D})}{\text{supp}(X, \mathcal{D})}$ .

For a threshold  $\delta$ , an REE is  $\delta$ -confident on  $\mathcal{D}$  if  $\text{conf}(\varphi, \mathcal{D}) \geq \delta$ .

**Minimality.** An REE  $\varphi : X \rightarrow p_0$  over  $\mathcal{R}$  is said to be *trivial* if  $p_0 \in X$ . In the rest of this paper, we only consider non-trivial REEs.

An REE  $\varphi : X \rightarrow p_0$  is *left-reduced* on  $\mathcal{D}$  if  $\varphi$  is  $\sigma$ -frequent,  $\delta$ -confident and moreover, there exists no REE  $\varphi'$  such that  $\varphi' \leq \varphi$  and  $\varphi'$  is  $\sigma$ -frequent and  $\delta$ -confident. Intuitively, it means no predicate in  $X$  can be removed, i.e., the minimality of predicates.

A *minimal* REE  $\varphi$  on  $\mathcal{D}$  is a non-trivial and left-reduced REE. Intuitively, there is no redundancy in a minimal REE.

**Attribute diversity.** To deduce a consequence  $p_0$ , we want different preconditions  $X$  to include diverse attributes so that if noises appear in some attributes, we can still use other attributes to deduce  $p_0$ . To do this, we maintain a counter  $\text{ct}_{p_0}(A)$  for each  $A$  in  $\mathcal{R}$  and increase it by 1 whenever an REE whose precondition involves  $A$  is discovered. Then, the *diversity* of an REE  $\varphi : X \rightarrow p_0$  is  $\text{div}(\varphi) = 1/\max\{\text{ct}_{p_0}(A) \mid A \text{ is an attribute that appears in } X\}$ .

**Succinctness.** We define the *succinctness* of REEs as  $\text{suc}(\varphi) = 1/|X|$  where  $|X|$  denotes the number of predicates used in  $X$ . It is widely agreed that REEs with larger succinctness are easier to understand and are statistically more likely to be true [65].

**3.1.2 Subjective Measures.** Different users have diverse preferences from rule discovery. A domain expert often wants rules that are unexpected, while a freshman in the community might be more interested in finding simple and efficient rules. For ER, one wants rules of the form  $X \rightarrow t.id = s.id$  or  $X \rightarrow t.id \neq s.id$ , while for CR, users are more interested in deducing  $p_0$  on other attributes. Thus the assumption that users have a universal preference is not always grounded. In light of this, below we train an ML model  $\mathcal{M}_{\text{sub}}$  for capturing subjective user preference, to compensate objective ones.

**Model architecture.** We learn ML model  $\mathcal{M}_{\text{sub}}$  by first transforming each rule into an embedding, and then feeding it into a lightweight model, which outputs a scalar score, indicating the subjective preference of users. The details are explained as follows.

Embedding. Given an REE  $\varphi : X \rightarrow p_0$ , we create a rule embedding.

(1) Firstly, we embed each predicate  $p$  in  $X$  by tokenizing its operator and operands into three tokens, say  $T_1, T_2$ , and  $T_3$ : (1) if  $p$  is  $t.A \oplus c$ , we tokenize it as  $t.A, \oplus$  and  $c$ ; similarly for  $t.A \oplus s.B$ ; and (2) if  $p$  is an ML predicate  $\mathcal{M}(t[\bar{A}], s[\bar{B}])$ , we treat it as  $\mathcal{M}()$ ,  $t[\bar{A}]$  and  $s[\bar{B}]$ . Similar to NLP tasks, we treat tokens as words and construct a vocabulary  $\mathcal{V}$ . For each  $T_i \in \mathcal{V}$ , we use ELMo [55], a state-of-the-art word representation model, to transform  $T_i$  to a  $d$ -dimensional representation  $\mathbf{T}_i \in \mathbb{R}^{d \times 1}$ . Then we adopt a linear layer  $\mathbf{w}_{\text{emb}} \in \mathbb{R}^{3 \times 1}$  to generate the predicate embedding  $E_p \in \mathbb{R}^{d \times 1}$  of  $p$ :

$$E_p = [T_1; T_2; T_3] \mathbf{w}_{\text{emb}},$$

where  $[\cdot]$  denotes concatenation;  $\mathbf{w}_{\text{emb}}$  is shared by all predicates.

(2) After obtaining the predicate embeddings for each  $p$  in  $X$ , we compute the *precondition embedding*, denoted by  $E_X \in \mathbb{R}^{d \times 1}$  for  $X$ . Recall that  $X$  is a conjunction  $p_1 \wedge p_2 \wedge \dots \wedge p_{|X|}$  of predicates. The embedding  $E_X$  should satisfy the *permutation invariant* property, i.e.,  $E_X = E_{X'}$ , where  $X' = p_{\pi_1} \wedge \dots \wedge p_{\pi_{|X|}}$ , and  $\pi_1, \dots, \pi_{|X|}$  is a new permutation of  $1, \dots, |X|$ . To achieve the permutation invariant property, we adopt *deep sets* [77] based on the predicate embeddings, and obtain  $E_X = \rho(\sum_{i=1}^{|X|} \Phi(E_{p_i}))$ , where  $\rho$  and  $\Phi$  are two linear layers without activation functions. Moreover, we use predicate embedding of  $p_0$  as the consequence embedding, denoted by  $E_{p_0}$ .

(3) Finally, the precondition embedding  $E_X$  and the consequence embedding  $E_{p_0}$  are concatenated together to form the rule embedding, denoted by  $E_\varphi \in \mathbb{R}^{(d_r+d) \times 1}$ , as shown in Equation 1:

$$E_\varphi = [E_X^T; E_{p_0}^T]^T \quad (1)$$

Lightweight model. Given the rule embedding  $E_\varphi$ , our lightweight model employs a fully-connected layer with the Sigmoid activation function to output the final subjective interestingness score:

$$\mathcal{M}_{\text{sub}}(\varphi) = \text{Sigmoid}(\mathbf{w}_{\text{light}}^T E_\varphi + b_{\text{light}}),$$

where  $\mathbf{w}_{\text{light}} \in \mathbb{R}^{(d_r+d) \times 1}$  and  $b_{\text{light}} \in \mathbb{R}$  are parameters of the lightweight model. Since we adopt the Sigmoid activation function, the subjective score is guaranteed to be in the range  $[0, 1]$ ,

Remark. Note that we do not adopt existing language models to acquire the rule-level representations for training  $\mathcal{M}_{\text{sub}}$  for the following reasons: (1) REEs do not follow natural language structure and thus, directly applying language model (e.g., Bert [13]) does not work well; (2) in REEs, token segmentation possesses different characteristic as text; and (3) given  $\varphi : X \rightarrow p_0$ ,  $X$  is a conjunction of predicates, but not a sequence of predicates as in text for which

there is no guarantee for the permutation invariant property.

## 3.2 Modeling User Preference

Putting these together, we are ready to define the *interestingness* of an REE  $\varphi$ . Let  $F$  and  $G$  be the set of objective and subjective measures. The interestingness of  $\varphi$  is defined to be

$$\text{interest}(\varphi) = \sum_{f \in F} w_f f(\varphi) + \sum_{g \in G} w_g g(\varphi),$$

where  $w_f$  and  $w_g$  are non-negative weights associated with each measures and they are sum up to 1 (we apply Softmax on learned parameters to guarantee this property). Moreover, we assume *w.l.o.g.* that each measure  $f$  in  $F$  outputs a scalar value from 0 to 1, where a larger value is more preferable; similarly for  $G$ .

Denote the complete weight vector by  $\mathbf{w}_{\text{interest}}$ , representing the user preference. Here a larger weight indicates that the corresponding measurement is more important to the user. Once the weight vector is determined, we can compute the interestingness of each rule, based on which we can deduce the top- $k$  interesting rules.

Note that the weight vector also depends on users' preference. Some users might rank support the most important, while the others might give diversity a higher priority. This highlights the need for learning individualistic preference by interacting with users.

## 3.3 An Active Learning Approach

We propose to jointly learn the subjective model  $\mathcal{M}_{\text{sub}}$  and the weight vector  $\mathbf{w}_{\text{interest}}$ . Denote the overall model by  $\mathcal{M}_{\text{interest}}$ , which is the combination of  $\mathcal{M}_{\text{sub}}$  and  $\mathbf{w}_{\text{interest}}$ . In the following, we show how we learn  $\mathcal{M}_{\text{interest}}$  based on *active learning*.

We adopt a pairwise ranking setting for users to label the partial orders of a few rules, since it is impractical to ask the users to label the actual interestingness score  $\text{interest}(\varphi)$  for each individual rule  $\varphi$ . More specifically, we maintain a rule pool  $S_{\text{REEs}}$ . Given a pair of rules  $\langle \varphi_i, \varphi_j \rangle$  in  $S_{\text{REEs}}$ , a user may label 1 on  $\langle \varphi_i, \varphi_j \rangle$  if s/he thinks that  $\varphi_i$  is more interesting than  $\varphi_j$ , denoted by  $\varphi_j \ll \varphi_i$ ; otherwise, the user labels the pair 0. We denote the label by  $y^{(i,j)} \in \{0, 1\}$ . For each training instance  $\langle \varphi_i, \varphi_j, y^{(i,j)} \rangle$ , we adopt the Siamese neural network with shared parameters to separately compute the interestingness scores of  $\varphi_i$  and  $\varphi_j$ . Then we use the Cross Entropy loss function to train  $\mathcal{M}_{\text{interest}}$  as follows.

$$\begin{aligned} \Pr(\varphi_i \ll \varphi_j) &= \text{Sigmoid}(\text{interest}(\varphi_i) - \text{interest}(\varphi_j)) \\ \mathcal{L}_{\text{CE}} &= \sum_{i,j} y^{(i,j)} \log(\Pr(\varphi_i \ll \varphi_j)) + (1 - y^{(i,j)}) \log(1 - \Pr(\varphi_i \ll \varphi_j)). \end{aligned}$$

Since there are  $|S_{\text{REEs}}| \times |S_{\text{REEs}}|$  pairs of rules, it is impractical for users to label all of them. Thus we adopt an active learning approach, which selects high-quality pairs of rules for users to label.

As shown in Figure 1, the active learner starts by randomly sampling a few pairs of rules from  $S_{\text{REEs}}$ ; it lets the user label them to generate the initial training set  $C_{\text{train}}$  of rules and the initial validation set  $C_{\text{valid}}$ , where  $C_{\text{train}}$  is used to train  $\mathcal{M}_{\text{interest}}$  and  $C_{\text{valid}}$  is used to evaluate its accuracy (lines 1-3).

We then iteratively learn  $\mathcal{M}_{\text{interest}}$  (lines 4-15). More specifically, we first train  $\mathcal{M}_{\text{interest}}$  using the current training data  $C_{\text{train}}$ . Then we actively select more rule pairs from  $S_{\text{REEs}}$  for users to label, aiming to improve the accuracy. To do this, we use the current  $\mathcal{M}_{\text{interest}}$  to compute the interestingness score for each rule in  $S_{\text{REEs}}$ . Then we select rule pairs in  $S_{\text{REEs}}$  that  $\mathcal{M}_{\text{interest}}$  cannot distinguish

---

**Algorithm ActiveLearner**

*Input:* A rule pool  $S_{\text{REEs}}$ , accuracy  $\text{acc}_{\text{min}}$ , maximum # of iterations  $\text{iter}$ , # of initial rule pairs  $M$ , and # of rule pairs  $N$  to label in each iteration.

*Output:* The interestingness model  $M_{\text{interest}}$ .

1. Construct  $C_{\text{REEs}}$ , consisting of  $M$  pairs of rules sampled from  $S_{\text{REEs}}$ ;
  2. Ask the user to label  $\mathcal{D}_{\text{REEs}}$ ;
  3. Split  $C_{\text{REEs}}$  into  $C_{\text{train}}$  and  $C_{\text{valid}}$ ;  $\text{step} := 0$ ; Initialize ML model  $M_{\text{interest}}$ ;
  4. **while**  $\text{step} \leq \text{iter}$  **do**
  5.   Incrementally train  $M_{\text{interest}}$  using the current  $C_{\text{train}}$ ;
  6.    $\text{acc}_{\text{valid}} := \text{Evaluate}(M_{\text{interest}}, C_{\text{valid}})$ ;
  7.   **if**  $\text{acc}_{\text{valid}} > \text{acc}_{\text{min}}$  **then**
  8.     **break** ;
  9.   **for each**  $\varphi \in S_{\text{REEs}}$  **do**
  10.     Compute the interestingness of  $\varphi$ , i.e.,  $M_{\text{interest}}(\varphi)$ ;
  11.     Select a set  $\Delta C_1$  of top- $\lceil \frac{N}{2} \rceil$  rule pairs  $\langle \varphi_i, \varphi_j \rangle$  in  $S_{\text{REEs}}$  with the smallest differences of interestingness  $|M_{\text{interest}}(\varphi_i) - M_{\text{interest}}(\varphi_j)|$ ;
  12.     Randomly select a set  $\Delta C_2$  of  $\lfloor \frac{N}{2} \rfloor$  rule pairs from  $S_{\text{REEs}}$ ;
  13.      $\Delta C := \Delta C_1 \cup \Delta C_2$  and ask the user to label  $\Delta C$ ;
  14.      $C_{\text{train}} := C_{\text{train}} \cup \Delta C$ ;  $\text{step} := \text{step} + 1$ ;
  15. **return**  $M_{\text{interest}}$ ;
- 

**Figure 1: Algorithm ActiveLearner**

well, i.e., pairs that have the smallest differences in interestingness scores, and send them to users for labeling. Intuitively, asking the user to label such pairs is more beneficial than labeling pairs with a clear margin, since it provides more information. To increase the diversity, we also randomly select a few rule pairs from  $S_{\text{REEs}}$  (line 12) for users to label. Finally, we include the newly labeled data in  $C_{\text{train}}$  (note that none of the newly labeled rules can be in  $C_{\text{valid}}$ ).

The process proceeds until it reaches the maximum number of iterations or the accuracy of  $M_{\text{interest}}$ , evaluated by using  $C_{\text{valid}}$ , reaches the minimum accuracy specified by users. We find it often suffices for users to label 50 rule pairs in 5 rounds of interaction.

## 4 INTERESTING RULE DISCOVERY

Based on the interestingness model, we discover top- $k$  interesting rules (Section 4.1) and develop such an algorithm (Section 4.2).

### 4.1 Problem Statement

Denote by  $\Sigma_{\text{all}}$  the set of minimal REEs on  $\mathcal{D}$  that are  $\sigma$ -frequent and  $\delta$ -confident for thresholds  $\sigma$  and  $\delta$ . As remarked earlier,  $\Sigma_{\text{all}}$  may contain an excessive number of rules that are not very relevant to users' applications and interest. To reduce such rules, we adopt the following strategies. (1) We pick an application-dependent set of candidate consequences  $p_0$ , denoted by RHS, which pertain to a particular application of users. (2) For each  $p_0$  in RHS, we learn a subset  $\mathcal{P}_0$  of predicates, including all predicates correlated to  $p_0$ ; we only focus on discovering REEs  $\varphi: X \rightarrow p_0$  such that  $X \subseteq \mathcal{P}_0$ .

**Top- $k$  REEs discovery.** The top- $k$  discovery problem is as follows.

- *Input:* A schema  $\mathcal{R}$ , an instance  $\mathcal{D}$  of  $\mathcal{R}$ , a consequence set RHS, the support/confidence thresholds  $\sigma/\delta$ , an integer  $k$ , and a weight vector  $\mathbf{w}$  on objective/subjective measures.
- *Output:* A set  $\Sigma$  consisting of the top- $k$  interesting REEs w.r.t.  $\mathbf{w}$  such that for each  $\varphi: X \rightarrow p_0$  in  $\Sigma$ , (a)  $p_0 \in \text{RHS}$ ; (b)  $X \subseteq \mathcal{P}_0$ , where  $\mathcal{P}_0$  is a set of predicates correlated to  $p_0$ ; and (c)  $\varphi$  is minimal and moreover, it is  $\sigma$ -frequent and  $\delta$ -confident.

### 4.2 A Top- $k$ Discovery Algorithm

We start with *pruning strategies* to remove early those REEs that are unlikely to become top- $k$  interesting rules, and reduce the large

---

**Algorithm Topk-Miner**

*Input:*  $\mathcal{D}$ , RHS,  $k$ ,  $\sigma$  and  $\delta$ .

*Output:* A heap  $\Sigma$  of top- $k$  interesting REEs such that for each  $\varphi: X \rightarrow p_0$  in  $\Sigma$ , (1)  $p_0 \in \text{RHS}$ ; (2)  $X \subseteq \mathcal{P}_0$ , where  $\mathcal{P}_0$  is a set of predicates correlated to  $p_0$ .

1.  $\Sigma :=$  an empty max-heap of maximum size  $k$ , ordered by rule interestingness;
2. Build auxiliary structures, e.g., position list indexes (PLI) [54];
3. **for each**  $p_0 \in \text{RHS}$  **do**
4.    $\mathcal{P}_{\text{sel}} := \emptyset$ ;  $\mathcal{P}_{\text{re}} := \mathcal{P}_0$ ;
5.    $\Sigma := \text{Expand}(\mathcal{D}, \mathcal{P}_{\text{sel}}, \mathcal{P}_{\text{re}}, p_0, k, \delta, \sigma, \Sigma)$ ;
6. **return**  $\Sigma$ ;

#### Procedure Expand

*Input:*  $\mathcal{D}$ ,  $\mathcal{P}_{\text{sel}}$ ,  $\mathcal{P}_{\text{re}}$ ,  $p_0$ ,  $k$ ,  $\delta$ ,  $\sigma$  and the current heap  $\Sigma$  of interesting REEs.

*Output:* An updated heap  $\Sigma$  of interesting REEs.

7.  $Q :=$  an empty queue;  $Q.\text{add}(\langle \mathcal{P}_{\text{sel}}, \mathcal{P}_{\text{re}} \rangle)$ ;
  8. **while**  $Q \neq \emptyset$  **do**
  9.    $\langle \mathcal{P}_{\text{sel}}, \mathcal{P}_{\text{re}} \rangle := Q.\text{pop}()$ ;
  10.    $\varphi := \mathcal{P}_{\text{sel}} \rightarrow p_0$ ;  $T_k :=$  the  $k$ -th highest interestingness in  $\Sigma$ ;
  11.   **if**  $\varphi$  is minimal (and thus, it is  $\sigma$ -frequent and  $\delta$ -confident) **then**
  12.     **if**  $\text{interest}(\varphi) > T_k$  **then**
  13.       Update  $\Sigma$  using  $\varphi$ ;
  14.     **continue**;
  15.   UB  $:=$  the interestingness upper bound of rules expanded from  $\varphi$  (Lemma 2);
  16.   **if**  $\exists \varphi' \in \Sigma$  s.t.  $\varphi' \leq \varphi$  [P1] or  $\text{supp}(\varphi) < \sigma$  [P2] or  $\text{UB} < T_k$  [P3] **then**
  17.     **continue**; // Early termination of the current expansion
  18.   **for each**  $p \in \mathcal{P}_{\text{re}}$  **do** // Add predicates from  $\mathcal{P}_{\text{re}}$  to  $\mathcal{P}_{\text{sel}}$
  19.      $Q.\text{add}(\langle \mathcal{P}_{\text{sel}} \cup \{p\}, \mathcal{P}_{\text{re}} \setminus \{p\} \rangle)$
  20. **return**  $\Sigma$ ;
- 

**Figure 2: Algorithm Topk-Miner**

number of candidate rules to be examined in top- $k$  rule discovery.

**Pruning strategy.** Our strategies are based on *anti-monotonicity* and *interestingness* and thus, REEs with high orders [P1], low supports [P2] or low interestingness [P3] are pruned early.

We maintain a heap  $\Sigma$  of top- $k$  minimal REEs that are discovered so far. Denote the  $k$ -th highest interestingness of rules in  $\Sigma$  by  $T_k$ . Assume that we are checking whether a candidate REE  $\varphi: X \rightarrow p_0$  is one of the top- $k$  interesting rules; if not, we revise and expand  $X$  with more predicates from  $\mathcal{P}_0$  to make such an REE if possible.

**[P1] Prune non-minimal REEs:** Before we perform exact checking for  $\varphi$ , we first check whether there exists an REE  $\varphi'$  discovered so far such that  $\varphi' \leq \varphi$ . If so, we can skip the processing and expansion for  $\varphi$ , since  $\varphi$  and all of its subsequent expansions are not minimal.

**[P2] Prune REEs with low support:** If  $\text{supp}(\varphi, \mathcal{D}) < \sigma$ , we do not consider any  $\varphi'$  such that  $\varphi \leq \varphi'$  since by Theorem 1, we have that  $\text{supp}(\varphi', \mathcal{D}) \leq \text{supp}(\varphi, \mathcal{D}) < \sigma$  and thus,  $\varphi'$  is not  $\sigma$ -frequent.

**[P3] Prune non-interesting REEs:** Recall that we maintain  $T_k$ , the  $k$ -th highest interestingness in  $\Sigma$ . To prune non-interesting REEs, the intuition is to compute an interestingness upper bound for rules expanded from  $\varphi$ . If the upper bound is less than  $T_k$ , we know that no rules expanded from  $\varphi$  can make a top- $k$  interesting rule and thus, we stop the exact expansion immediately.

In the following, we categorize the interestingness measures into three types: *monotonic*, *anti-monotonic* and *general measures*, based on which we develop the interestingness upper bound.

An interesting measure  $f$  is *monotonic* if  $f(\varphi) \leq f(\varphi')$  as long as  $\varphi \leq \varphi'$ , i.e., adding predicates to  $\varphi$  monotonically increases the interestingness. In contrast, adding predicates in an *anti-monotonic* measure monotonically decreases the interestingness, e.g., support.

If an interesting measure does not have the above properties, it is referred to as *general*, e.g., usually the subjective model.

**Lemma 2:** Given an REE  $\varphi : X \rightarrow p_0$ , let  $\varphi'$  be an REE expanded from  $\varphi$ , and  $F$  and  $G$  be the set of objective/subjective measures. Then we have  $\text{interest}(\varphi') \leq \sum_{f \in F} w_f f_{\text{ub}}(\varphi') + \sum_{g \in G} w_g g_{\text{ub}}(\varphi')$ , where  $f_{\text{ub}}$  and  $g_{\text{ub}}$  are the upper bounds for  $F$  and  $G$ , respectively. Specifically, if  $f$  is monotonic,  $f_{\text{ub}}(\varphi') = f(\mathcal{P}_0 \rightarrow p_0)$ ; if  $f$  is anti-monotonic,  $f_{\text{ub}}(\varphi') = f(\varphi)$ ; and if  $f$  is general,  $f_{\text{ub}}(\varphi') = 1$ ; similarly for  $g$ .  $\square$

**Proof.** We prove it by contradiction. Assume that  $\text{interest}(\varphi') > \sum_{f \in F} w_f f_{\text{ub}}(\varphi') + \sum_{g \in G} w_g g_{\text{ub}}(\varphi')$ . There must exist at least a measure (objective or subjective), say  $f$  or  $g$ , such that  $f(\varphi') > f_{\text{ub}}(\varphi')$  or  $g(\varphi') > g_{\text{ub}}(\varphi')$ . Assume w.l.o.g. that  $f$  is such a measure and it is monotonic. Thus,  $f(\varphi') > f(\mathcal{P}_0 \rightarrow p_0)$ , which contradicts to the definition of monotonic measures that adding more predicates only increases the interestingness; similarly for the other cases.  $\square$

**Example 5:** Assume that the  $k$ -highest interestingness  $T_k$  in  $\Sigma$  is 0.8, and the REE that we are currently processing is  $\varphi$ . If the interestingness upper bound of the REEs expanded from  $\varphi$  is found to be 0.7, we can stop the processing of  $\varphi$  immediately, without actually expanding  $\varphi$ , since it will not yield any top- $k$  rules.  $\square$

**Algorithm.** We now present our algorithm, referred to as Topk-Miner, for top- $k$  interesting REEs discovery on dataset  $\mathcal{D}$ .

As shown in Figure 2, Topk-Miner is a *levelwise search* algorithm. It first initializes a max-heap  $\Sigma$  of maximum size  $k$  (line 1), which is used to store the top- $k$  REEs discovered so far, ordered by their interestingness. Given a consequence  $p_0$  in RHS and its correlated  $\mathcal{P}_0$ , we maintain two predicate sets for discovering new REEs  $\varphi : X \rightarrow p_0$  with  $X \subseteq \mathcal{P}_0$ : (1)  $\mathcal{P}_{\text{sel}}$ , the set of predicates selected to constitute  $X$ ; and (2)  $\mathcal{P}_{\text{re}}$ , the set of remaining predicates in  $\mathcal{P}_0$ . Initially,  $\mathcal{P}_{\text{sel}}$  is empty and  $\mathcal{P}_{\text{re}}$  is  $\mathcal{P}_0$  (line 4). Topk-Miner then traverses the search space level by level by maintaining a queue  $Q$  (line 7), where at the  $i$ -th level, it discovers  $\varphi : X \rightarrow p_0$  with  $|X| = i$ . It iteratively adds predicates from  $\mathcal{P}_{\text{re}}$  to  $\mathcal{P}_{\text{sel}}$  (line 18-19) until one of the following conditions is satisfied: (1)  $\mathcal{P}_{\text{re}}$  is exhaustive; or (2)  $\varphi : \mathcal{P}_{\text{sel}} \rightarrow p_0$  is a minimal REE (line 10-14), since in this case, adding more predicates will not make  $\text{supp}(\varphi, \mathcal{D})$  larger, while it increases the order of  $\varphi$ . We maintain the  $k$ -th highest interestingness  $T_k$  of REEs in  $\Sigma$  (line 10). If  $\text{interest}(\varphi) > T_k$ ,  $\Sigma$  is updated (line 12-13) by adding  $\varphi$  into  $\Sigma$  and removing the least interesting REE from  $\Sigma$  if there are more than  $k$  REEs in  $\Sigma$ . If  $\varphi : \mathcal{P}_{\text{sel}} \rightarrow p_0$  is still not a minimal REE, we expand it (line 18-19); before expansion, we apply the pruning strategies [P1]-[P3] (line 15-17), to check whether we can terminate the current expansion early.

Algorithm Topk-Miner employs the following optimization strategies commonly used in rule discovery. (a) When multiple  $p_0$  in RHS share similar correlated predicates  $\mathcal{P}_0$ , it processes these  $p_0$ 's together (not shown). (b) It pre-computes auxiliary structures (line 2) such as position list indexes (PLI) [54], to efficiently compute supports and confidences for REEs verification.

We also propose the following for top- $k$  REE discovery.

**Correlated predicate learning.** Recall that for each consequence  $p_0$ , we learn a subset  $\mathcal{P}_0$  of its correlated predicates. In addition to logic predicates,  $\mathcal{P}_0$  includes ML predicates as well. To this end, we maintain a pool of pre-trained ML models. Given a schema  $\mathcal{R}$ , we

associate attributes in  $\mathcal{R}$  to compatible ML models in the pool and initialize the ML predicates. Then, to learn correlated  $\mathcal{P}_0$  for a given  $p_0$ , we apply graphical lasso [24, 78] to learn how an attribute is affected by other attributes. Informally, given a predicate  $p$ , either a logic or an ML predicate, if attributes in  $p$  have strong impact on the attributes in  $p_0$ ,  $p$  is correlated to  $p_0$  and is included in  $\mathcal{P}_0$ .

**Handling multiple relation atoms.** To efficiently support multiple relation atoms in Topk-Miner, we incrementally discover multi-variable REEs in rounds at each level (not shown). In the  $j$ -th round, we discover  $j$ -variable REEs based on the results of the  $(j-1)$ -th round. That is, we process each non-minimal  $(j-1)$ -variable REE  $\varphi$  found in the  $(j-1)$ -th round, by constructing  $\mathcal{P}_{\text{re}}$  incrementally, which is the set of remaining predicates that can be used to expand  $\varphi$ , such that the expanded rules contain exactly  $j$  relation atoms. Here  $\mathcal{P}_{\text{re}}$  can be constructed incrementally by enumerating the predicates that contain one new relation atom and at most one existing relation atom used in  $\varphi$ . Then we discover  $j$ -variable rules by expanding  $\varphi$  with the predicates in the newly constructed  $\mathcal{P}_{\text{re}}$ .

**Early termination.** As remarked above, Topk-Miner allows early termination for top- $k$  discovery, i.e., we can terminate the expansion of an REE  $\varphi$  early if one of the following happens (line 13): (1) if there exists an REE  $\varphi'$  in  $\Sigma$  such that  $\varphi' \leq \varphi$ , then there is no need to expand  $\varphi$ , which cannot be minimal [P1]; (2) if  $\text{supp}(\varphi) \leq \sigma$  [P2], further expanding  $\varphi$  will not lead to  $\sigma$ -frequent REEs by the anti-monotonicity of support; or (3) if  $\text{UB} < T_k$  [P3], where  $\text{UB}$  is the interestingness upper bound of the rules expanded from  $\varphi$  (computed by Lemma 2), then no rules expanded from  $\varphi$  are more interesting than any rule in  $\Sigma$ , and we can stop immediately.

In addition, we further optimize Topk-Miner by considering effective processing orders for the predicates in  $\mathcal{P}_{\text{re}}$  (line 18). Intuitively, a proper processing order can speed up Topk-Miner, e.g., by tightening the interestingness bound at a faster rate. Below we present two practical processing orders of  $\mathcal{P}_{\text{re}}$  based on support and interestingness, respectively, where the former guarantees to prune some predicates in  $\mathcal{P}_{\text{re}}$  under certain conditions, while the latter *dynamically* determines the next predicate to be processed so that rules with higher interestingness are likely to be discovered earlier.

**(1) Support-based processing order.** The first strategy is to add the predicates  $p$  from  $\mathcal{P}_{\text{re}}$  to  $\mathcal{P}_{\text{sel}}$  based on  $\text{supp}(\mathcal{P}_{\text{sel}} \wedge p, \mathcal{D})$ , such that predicates with high supports are processed first. This processing order helps us prune those predicates in  $\mathcal{P}_{\text{re}}$  that are clearly useless in generating  $\sigma$ -frequent REEs (in addition to [P2]). The intuition is that if including a predicate  $p$  with high support cannot lead to an  $\sigma$ -frequent REE, it is even more difficult for a predicate  $p'$  with low support to do so. Formally, we have the following lemma.

**Lemma 3:** Given an REE  $\varphi : \mathcal{P}_{\text{sel}} \rightarrow p_0$ , and two predicates  $p$  and  $p'$  in  $\mathcal{P}_{\text{re}}$ , expanding  $\mathcal{P}_{\text{sel}}$  with  $p'$  will not give any  $\sigma$ -frequent REE if  $\mathcal{P}_{\text{sel}} \wedge p \rightarrow p_0$  is not  $\sigma$ -frequent and  $\text{spset}(p', \mathcal{D}) \subseteq \text{spset}(p, \mathcal{D})$ .  $\square$

**Proof.** Let  $\varphi$  and  $\varphi'$  be  $\mathcal{P}_{\text{sel}} \wedge p \rightarrow p_0$  and  $\mathcal{P}_{\text{sel}} \wedge p' \rightarrow p_0$ , respectively. Since  $\text{spset}(p', \mathcal{D}) \subseteq \text{spset}(p, \mathcal{D})$ , we have  $\text{spset}(\mathcal{P}_{\text{sel}} \wedge p' \wedge p_0, \mathcal{D}) \subseteq \text{spset}(\mathcal{P}_{\text{sel}} \wedge p \wedge p_0, \mathcal{D})$ . Then for each  $h'$  in  $\text{spset}(\varphi', \mathcal{D})$ ,  $h'$  must be in  $\text{spset}(\varphi, \mathcal{D})$ . Thus,  $\text{supp}(\varphi', \mathcal{D}) \leq \text{supp}(\varphi, \mathcal{D}) < \sigma$ .  $\square$

Note that  $\text{spset}(p', \mathcal{D}) \subseteq \text{spset}(p, \mathcal{D})$  indicates that  $\text{supp}(\mathcal{P}_{\text{sel}} \wedge$

$p, \mathcal{D}) \geq \text{supp}(\mathcal{P}_{\text{sel}} \wedge p', \mathcal{D}))$ , and thus,  $p$  will be processed before  $p'$ . In other words, if  $\mathcal{P}_{\text{re}}$  is ordered by supports, every time we see a predicate  $p$  and if  $\mathcal{P}_{\text{sel}} \wedge p \rightarrow p_0$  is not  $\sigma$ -frequent, we can prune all  $p'$  in  $\mathcal{P}_{\text{re}}$  ordered after  $p$  with  $\text{spset}(p', \mathcal{D}) \subseteq \text{spset}(p, \mathcal{D})$ .

**Example 6:** Consider the relation in Table 2. Assume that  $\mathcal{P}_{\text{sel}} = \{\text{RentalContract}(t)\}$  and  $p_0$  is  $t.\text{montly\_rent} = 1,500$ . Let  $p$  and  $p'$  be  $t.\text{renter} = \text{"Ada"}$  and  $t.\text{pid} = p_3$ , respectively. Clearly,  $\text{spset}(p', \mathcal{D}) = \{t_7 \mapsto t, t_8 \mapsto t\}$  is a subset of  $\text{spset}(p, \mathcal{D}) = \{t_7 \mapsto t, t_8 \mapsto t, t_9 \mapsto t\}$  and thus,  $p$  is processed before  $p'$ . If we find that  $\mathcal{P}_{\text{sel}} \wedge p \rightarrow p_0$  is not  $\sigma$ -frequent, there is no need to process  $p'$ , since  $\text{supp}(\mathcal{P}_{\text{sel}} \wedge p \rightarrow p_0) = 2 > 1 = \text{supp}(\mathcal{P}_{\text{sel}} \wedge p' \rightarrow p_0)$ .  $\square$

(2) *Interestingness-based processing order.* Alternatively, we can also process the predicates in  $\mathcal{P}_{\text{re}}$  based on their “potential” in discovering rules with high interestingness. Intuitively, if more high interestingness rule are discovered, the interestingness bound  $T_k$  (the  $k$ -th highest interestingness observed so far) is tighter and thus, more rules are likely to be pruned by [P3] at an earlier stage.

More specifically, given REE  $\varphi : \mathcal{P}_{\text{sel}} \rightarrow p_0$  and a predicate  $p$  in  $\mathcal{P}_{\text{re}}$ , we define an indicator  $\Delta_p$ , expressing the best possible interestingness that can be achieved by including  $p$  to  $\mathcal{P}_{\text{sel}}$ :

$$\Delta_p = \sum_{f \in F} w_f f_{\text{ub}}(\varphi') + \sum_{g \in G} w_g g_{\text{ub}}(\varphi'),$$

where  $\varphi'$  is  $\mathcal{P}_{\text{sel}} \wedge p \rightarrow p_0$ , and  $f_{\text{ub}}, g_{\text{ub}}$  are defined in Lemma 2.

Then, the predicate  $p$  in  $\mathcal{P}_{\text{re}}$  with the maximum  $\Delta_p$  is selected as the next predicate to be used to expand  $\mathcal{P}_{\text{sel}}$ .

*Remark.* The support and interestingness based processing orders can be combined together, e.g., by ordering predicates in  $\mathcal{P}_{\text{re}}$  by their supports and breaking ties based on the interestingness indicator.

**Example 7:** We show how  $\varphi_2$  in Example 1 is found. Assume that  $\mathcal{P}_{\text{sel}} = \{\text{Property}(t_a), \text{RentalContract}(s_a), \text{RentalContract}(s_b), \mathcal{M}_{\text{rent}}(t_a, s_a)\}$ ,  $\mathcal{P}_{\text{re}} = \{t_a.\text{pid} = s_a.\text{pid}, s_a.\text{rent\_start\_year} \neq s_b.\text{rent\_start\_year}\}$  and  $p_0$  is  $s_a.\text{monthly\_rent} \neq s_b.\text{monthly\_rent}$ . We expand  $\mathcal{P}_{\text{sel}}$  by adding predicates in  $\mathcal{P}_{\text{re}}$  one by one. In fact, before we perform the exact expansion, we first compute an interestingness upper bound, say UB, by Lemma 2. If UB is less than the  $k$ -th highest interestingness maintained in  $\Sigma$ , we can terminate early since expanding  $\mathcal{P}_{\text{sel}}$  will not result in any top- $k$  interesting REEs [P3]. Otherwise, assume that  $p : t_a.\text{pid} = s_a.\text{pid}$  has a larger  $\text{supp}(\mathcal{P}_{\text{sel}} \wedge p)$ . Then we add  $p$  to  $\mathcal{P}_{\text{sel}}$  first, by the support-based processing order, until we find that  $\varphi_2$  is a minimal REE. If  $\varphi_2$  is more interesting than some rules in  $\Sigma$ ,  $\Sigma$  is updated using  $\varphi_2$ .  $\square$

*Correctness.* The correctness of Topk-Miner is verified as follows.

**Theorem 4:** Topk-Miner correctly discovers top- $k$  REEs.  $\square$

**Proof.** One can verify that Topk-Miner does not miss any rule  $\varphi$  unless there exists an REE  $\varphi'$  with  $\varphi' \leq \varphi$ , such that when processing  $\varphi'$ , we decide to stop expanding  $\varphi'$  according to the pruning strategies [P1]-[P3]. This is correct since by [P1]-[P3],  $\varphi$  cannot be a top- $k$  REE, since it is non-minimal, not  $\sigma$ -frequent, or less interesting than at least  $k$  rules that have been discovered.  $\square$

## 5 ANYTIME DISCOVERY

In this section, we extend Topk-Miner and develop an anytime algorithm Anytime-Miner, such that users can continuously get

the next top- $k$  interesting rules when needed, via lazy evaluation.

A brute-force approach for supporting this is to compute the full ranking of all REEs first. Every time a user wants the next top- $k$  results, we retrieve the corresponding results from the ranked list. Clearly, this method is inefficient. Users are typically only interested in the first few top-ranked results, and should not pay the cost of waiting for discovering the entire set of REEs on a dataset.

We convert Topk-Miner into Anytime-Miner. Denote by  $\Sigma$  (see Figure 2) the heap of discovered top-ranked REEs so far. We expand  $\Sigma$  to lazily discover the next top- $k$  results as follows.

(1) Instead of just maintaining the top- $k$  results in  $\Sigma$ , all minimal rules discovered are kept in  $\Sigma$ , referred to as *complete rules*. In addition, we maintain *partial rules* in  $\Sigma$ , where an REE  $\varphi$  is said to be *partial* if at the time it is processed, its interestingness upper bound is lower than those of at least  $k$  complete rules in  $\Sigma$ . In other words, the pruning strategy [P3] in the original Topk-Miner is revised: instead of directly dropping those rules with relatively low interestingness, we keep them as partial rules in  $\Sigma$ . Intuitively, these partial rules are likely to be expanded and contribute to the top- $k$  ones in later rounds when the user wants more results. For partial rules, their remaining predicates, say  $\mathcal{P}_{\text{re}}$ , are also stored for later expansion, and we use its interestingness upper bound as the key in  $\Sigma$ .

(2) We only return the next top- $k$  *complete rules* in  $\Sigma$ . For each partial rule whose interestingness upper bound is among the next top- $k$ , we resume its levelwise search in order. Since we have stored the remaining predicates  $\mathcal{P}_{\text{re}}$  for each partial rule, the resumption is straightforward. The resumed search updates the rules maintained in  $\Sigma$ ; it continues until the next top- $k$  rules in  $\Sigma$  all become complete.

(3) We ensure that the next top- $k$  results are not “redundant”, i.e., not logical consequences of the rules that have been shown before. Thus, we apply the implication analysis [20] on each newly discovered rules. Formally, we say that a set of REEs  $\Sigma$  entails another REE  $\varphi$  over  $\mathcal{R}$ , denoted by  $\Sigma \models \varphi$ , if for any instance  $\mathcal{D}$  of  $\mathcal{R}$ , if  $\mathcal{D} \models \Sigma$  then  $\mathcal{D} \models \varphi$ . Then, every time a complete rule  $\varphi$  is discovered, we add it to the heap  $\Sigma$  only if  $\Sigma \not\models \varphi$ . While the implication problem is  $\Pi_2^P$ -complete [21], we develop an efficient heuristic for checking.

**Example 8:** Assume that  $k = 3$  and rules in  $\Sigma$  are currently stored in order:  $\varphi_1^c, \varphi_2^p, \varphi_3^p, \varphi_4^c$ , where  $\varphi_i^c$  and  $\varphi_j^p$  denote complete rules and partial rules, respectively. Since there are partial rules in top-3 of  $\Sigma$ , we process them in order. Assume that we first resume the levelwise search for  $\varphi_2^p$  and obtain three new REEs:  $\varphi_5^c, \varphi_6^c, \varphi_7^p$ , and  $\Sigma$  is updated as:  $\varphi_1^c, \varphi_5^c, \varphi_6^c, \varphi_7^p, \varphi_3^p, \varphi_4^c$ . At this point, all top-3 rules in  $\Sigma$  ( $\varphi_1^c, \varphi_5^c, \varphi_6^c$ ) are complete rules, and they are returned to the user.  $\square$

## 6 PARALLEL TOP- $k$ RULE DISCOVERY

In this section we parallelize top- $k$  discovery to scale with large datasets. We first review a criterion for measuring the effectiveness of parallel algorithms (Section 6.1). We then parallelize Topk-Miner, denoted by PTopk-Miner, with the performance guarantees (Section 6.2); Anytime-Miner is parallelized along the same lines.

### 6.1 Parallel Scalability

We revisit the widely adopted notion of parallel scalability [39].

Assume that  $\mathcal{A}$  is a sequential algorithm which, given a dataset



$\mathcal{D}$ , consequences RHS and thresholds  $\sigma$  and  $\delta$  for support and confidence, respectively, computes a set  $\Sigma$  of top- $k$  interesting REEs on  $\mathcal{D}$ . Denote its worst running time as  $t(|\mathcal{D}|, |\text{RHS}|, \sigma, \delta)$ . We say that a parallel algorithm  $\mathcal{A}_p$  is *parallelly scalable relative to  $\mathcal{A}$*  if its running time by using  $n$  processors can be expressed as:

$$T(|\mathcal{D}|, |\text{RHS}|, \sigma, \delta) = \tilde{O}\left(\frac{t(|\mathcal{D}|, |\text{RHS}|, \sigma, \delta)}{n}\right),$$

where the notation  $\tilde{O}()$  hides  $\log(n)$  factors.

Intuitively, parallel scalability guarantees “linear” speedup of  $\mathcal{A}_p$  relative to the “yardstick” algorithm  $\mathcal{A}$ . That is, the more processors are used, the faster  $\mathcal{A}_p$  is. Hence  $\mathcal{A}_p$  can scale with large databases by adding processors and makes REEs discovery feasible in practice.

## 6.2 Parallel Algorithm

We next parallelize Topk-Miner to be PTopk-Miner (Figure 3).

**Setting.** PTopk-Miner runs with one coordinator  $S_c$  and  $n$  workers  $P_1, \dots, P_n$  under the Bulk Synchronous Parallel (BSP) model [71], where the coordinator is responsible for distributing and balancing workloads, and workers discover rules in parallel. The overall computation is divided into supersteps of a fixed duration.

**Overview.** Similar to Topk-Miner, the coordinator maintains a max-heap of maximum size  $k$ , consisting of the top-ranked REEs discovered so far (line 1). Denote the heap at superstep  $i$  by  $\Sigma_i$ , and the  $k$ -th highest interestingness in  $\Sigma_i$  by  $T_k^i$ . The coordinator first distributes the workloads evenly to all the workers (see below; line 2-5). Then, each worker parallelly processes its workload and discovers rules in supersteps (line 6-15). At each superstep, the coordinator informs each worker the latest interestingness bound  $T_k^i$  (line 10), based on which each worker performs the subsequent discovery (line 11) by applying the pruning strategies in Section 4.2. The coordinator  $S_c$  pulls the newly discovered top- $k$  rules from each worker at the end of each superstep (line 12). In addition, it adjusts and balances the workload when needed (line 13-14; see below). Moreover,  $S_c$  extends the heap  $\Sigma_i$  to  $\Sigma_{i+1}$  with the new rules, and updates interestingness bound  $T_k^i$  to  $T_k^{i+1}$  (line 15). The process continues until all workers finish their work, i.e., when no more rules with interestingness above the bound can be found.

**Workload assignment.** Given RHS,  $S_c$  evenly divides it into  $n$  partitions, namely  $\text{RHS}_1, \dots, \text{RHS}_n$ , and constructs a set of *work units* based on each  $\text{RHS}_j$  ( $j \in [1, n]$ ) for the  $j$ -th worker  $P_j$  as follows.

For each consequence  $p_0$  in  $\text{RHS}_j$ , it constructs a work unit, which is a triple  $w = \langle \mathcal{P}_{\text{sel}}, \mathcal{P}_{\text{re}}, p_0 \rangle$ , where  $\mathcal{P}_{\text{sel}}$  denotes the set of predicates that are selected to constitute the rules, and  $\mathcal{P}_{\text{re}}$  denotes the set of remaining predicates. Initially,  $\mathcal{P}_{\text{sel}}$  is empty and  $\mathcal{P}_{\text{re}}$  is  $\mathcal{P}_0$ , which is the set of predicates correlated to  $p_0$ . Then, it sends workload  $\mathcal{W}_j = \{w = \langle \mathcal{P}_{\text{sel}}, \mathcal{P}_{\text{re}}, p_0 \rangle \mid p_0 \in \text{RHS}_j\}$  to worker  $P_j$ .

Upon receiving  $\mathcal{W}_j$ , worker  $P_j$  fetches a subset  $\mathcal{D}_{\mathcal{W}_j}$  of data from  $\mathcal{D}$ , guided by  $\mathcal{W}_j$ , where  $\mathcal{D}_{\mathcal{W}_j} = \{t \in \mathcal{D} \mid \exists s \in \mathcal{D}, p \in \mathcal{P}_0 \text{ s.t. } h\langle t, s \rangle \models p \text{ or } h\langle t, s \rangle \models p_0, \text{ where } p_0 \in \text{RHS}_j\}$ ; it also constructs the corresponding auxiliary structures for performing rule discovery. In this way, the same data will be merged and transmitted to  $P_j$  only once even if it satisfies multiple predicates, reducing the total communication cost when processing multiple predicates.

**Example 9:** Consider relation RentalContract in Table 2. Let  $p$  be  $t.\text{pid} = s.\text{pid}$  and  $p'$  be  $t.\text{renter} = s.\text{renter}$ . Assume that the

### Algorithm PTopk-Miner

*Input:*  $\mathcal{D}$ , RHS,  $k$ ,  $\sigma$ ,  $\gamma$ , a coordinator  $S_c$  and  $n$  workers  $P_1, \dots, P_n$ .

*Output:* A max-heap  $\Sigma$  of top- $k$  REEs on  $\mathcal{D}$ .

*/\* executed at coordinator  $S_c$  \*/*

1.  $i := 0$ ;  $\Sigma_i :=$  an empty max-heap of maximum size  $k$ ;
2. **for each**  $p_0 \in \text{RHS}$  **do**
3.   Construct a work unit  $w = \langle \mathcal{P}_{\text{sel}}, \mathcal{P}_{\text{re}}, p_0 \rangle$ , where  $\mathcal{P}_{\text{sel}} = \emptyset$  and  $\mathcal{P}_{\text{re}} = \mathcal{P}_0$ ;
4. Evenly divide RHS into  $n$  partitions, namely  $\text{RHS}_1, \dots, \text{RHS}_n$ ;
5. Assign workload  $\mathcal{W}_j = \{w = \langle \mathcal{P}_{\text{sel}}, \mathcal{P}_{\text{re}}, p_0 \rangle \mid p_0 \in \text{RHS}_j\}$  to worker  $P_j$ ;
6. */\* run on  $n$  workers in parallel, in supersteps \*/*
7. **for each worker**  $P_j$  **do**
8.   Fetch  $\mathcal{D}_{\mathcal{W}_j} = \{t \in \mathcal{D} \mid \exists s \in \mathcal{D}, p \in \mathcal{P}_0 \text{ s.t. } h\langle t, s \rangle \models p \text{ or } h\langle t, s \rangle \models p_0, \text{ where } p_0 \in \text{RHS}_j\}$  and build the corresponding auxiliary structures;
9. **while** there exists unfinished work **do** */\* superstep  $i$  \*/*
10.   **for each**  $P_j$  with non-empty workload  $\mathcal{W}_j$  **do**
11.      $T_k^i :=$  the  $k$ -th highest interestingness in  $\Sigma_i$  (informed by  $S_c$ );
12.     Run Topk-Miner at  $P_j$  based on  $\mathcal{W}_j$  and  $\mathcal{D}_{\mathcal{W}_j}$  in parallel;
13.      $S_c$  pulls top- $k$  REEs  $\varphi$  newly discovered ( $\text{interest}(\varphi) > T_k^i$ );
14.     **for each**  $P_x$  that has finished the assigned workload **do**
15.       Balance workload between  $P_j$  and  $P_x$ , where  $P_j$  is the heaviest worker;
16.     Upon receiving new REEs from workers,  $S_c$  updates  $\Sigma_i$  to  $\Sigma_{i+1}$ ;
17.     update  $T_k^i$  to  $T_k^{i+1}$ ; broadcast  $T_k^{i+1}$  to all workers;  $i := i + 1$ ;
18. **return**  $\Sigma_i$ ;

Figure 3: Algorithm PTopk-Miner

coordinator  $S_c$  constructs and assigns workload  $\mathcal{W}_j = \{w, w'\}$  to worker  $P_j$ , where  $w = \langle \emptyset, \{p\}, p_0 \rangle$  and  $w' = \langle \emptyset, \{p'\}, p'_0 \rangle$ . It is easy to see that  $\mathcal{D}_{\mathcal{W}_j} = \{t_7, t_8, t_9\}$ . In particular, although  $h\langle t_7, t_8 \rangle \models p$  and  $h\langle t_7, t_8 \rangle \models p'$ ,  $t_7$  and  $t_8$  are only transmitted to  $P_j$  once.  $\square$

**Workload balancing.** At each superstep, if the workloads across workers are “skewed”, i.e., there is an idle worker  $P_x$  that has finished its assigned works, we re-distribute the workload to  $P_x$  from the heaviest worker  $P_j$ , in the following two steps.

- (1) If there are more than one work unit in  $\mathcal{W}_j$ ,  $P_j$  sends half of  $\mathcal{W}_j$  (and the corresponding auxiliary structures) to  $P_x$ .
- (2) If there is only one remaining work unit  $w = \langle \mathcal{P}_{\text{sel}}, \mathcal{P}_{\text{re}}, p_0 \rangle$  in  $\mathcal{W}_j$ , we split this heavy unit into two smaller ones, namely  $w' = \langle \mathcal{P}_{\text{sel}} \cup \{p\}, \mathcal{P}_{\text{re}} \setminus \{p\}, p_0 \rangle$  and  $w'' = \langle \mathcal{P}_{\text{sel}}, \mathcal{P}_{\text{re}} \setminus \{p\}, p_0 \rangle$ , where  $p$  is the predicate in  $\mathcal{P}_{\text{re}}$  with the highest processing order (see the processing order in Section 4.2), and send one of the two to  $P_x$ . Intuitively, it means that we divide the work unit  $w$  into two, i.e., selecting  $p$  into  $\mathcal{P}_{\text{sel}}$  and excluding  $p$  from  $\mathcal{P}_{\text{sel}}$ .

**Parallel scalability.** The parallel scalability is shown as follows.

**Theorem 5:** Algorithm PTopk-Miner is *parallelly scalable relative to the sequential algorithm* Topk-Miner.  $\square$

**Proof.** The time complexity of Topk-Miner is  $t(|\mathcal{D}|, |\text{RHS}|, \sigma, \delta) = O(\sum_{\varphi \in C(\mathcal{P}_0) \times \text{RHS}} |\mathcal{D}|^{|\varphi|})$ , where  $C(\mathcal{P}_0)$  is the power set of  $\mathcal{P}_0$  and  $|\varphi|$  denotes the number of predicates in  $\varphi$ , since Topk-Miner examines the entire  $C(\mathcal{P}_0)$  for each  $p_0 \in \text{RHS}$  in the worst case.

We next show that the parallel runtime of PTopk-Miner is in  $O(\frac{t(|\mathcal{D}|, |\text{RHS}|, \sigma, \delta)}{n})$ . In PTopk-Miner,  $S_c$  conducts workload assignment and heap maintenance, which maintains the global top- $k$  by collecting REEs from each worker. The former takes  $O(|\text{RHS}|)$  time, while the latter takes  $O(nk \log(k))$  time using merge-sort. Both are smaller than the discovery cost, which dominates the complexity.

The cost at each worker is dominated by the following: (a) trans-

mit its top- $k$  rules to the coordinator in time much less than  $O(|\mathcal{D}|)$  since each rule is discovered from  $\mathcal{D}$  and  $k$  is a small number; (b) receive  $T_k$  from  $S_c$  in  $O(1)$  time; (c) balance its workload, where  $O(|\mathcal{D}|)$  data is sent to idle workers; and (d) locally perform discovery in  $O(\frac{t(|\mathcal{D}|, |\text{RHS}|, \sigma, \delta)}{n})$  time, since the workload is evenly distributed in (c). Taken together, the parallel cost of PTopk-Miner is  $O(\frac{t(|\mathcal{D}|, |\text{RHS}|, \sigma, \delta)}{n})$  in the worst-case. In practice, the pruning strategies in Section 4.2 effectively remove useless candidates.  $\square$

## 7 EXPERIMENTAL STUDY

Using real-life and synthetic data, we experimentally evaluated (1) the scalability of PTopk-Miner for top- $k$  discovery and Anytime-Miner for anytime discovery, (2) the effectiveness of the proposed interestingness model, and (3) the effectiveness of top- $k$  discovery.

**Experimental setting.** We start with the experimental setting.

**Datasets.** Following the setting in studies [46, 54], we used seven datasets as shown in Table 3. Adult, Airport, Hospital, Inspection and NCVoter are real-life datasets that are commonly used in the literature. We additionally used an academic dataset DBLP that has multiple relations, and a synthetic dataset Tax that is obtained by first duplicating tuples of the original tax data (1M) [8, 12] 10 times and then modifying their attributes using a program of [17].

**ML models.** We used two ML models as predicates in REEs. (1) For long textual attributes where logic predicates do not work well, we applied SentenceBert [59]. (2) For ER, we adopted ditto [44]. For the interestingness model, we set the embedding size to 512. We used Adam optimizer with a batch-size of 20; the learning rate is 0.001. We trained our model with 50 epochs on Tesla V100 GPU.

**Baselines.** We implemented the following, all in Java. (1) PTopk-Miner, (2) Anytime-Miner; (3) PTopk-Miner<sub>nop</sub>, a variant of PTopk-Miner without using pruning strategy [P3], which enumerates all minimal rules, computes their interestingness and returns the top- $k$  ones; (4) PTopk-Miner<sub>dc</sub>, another variant of PTopk-Miner, which discovers only DCs, *i.e.*, REEs on single relations without ML predicates; (5) DCFinder [54], a state-of-the-art DCs discovery algorithm, which mines all DCs that hold on the dataset; as shown in [54], it outperforms other DC discovery methods [54]. We parallelize it for a fair comparison and extend it to support constant predicates of the form  $t.A = c$ . We compared with PTopk-Miner<sub>nop</sub> to test the effectiveness of pruning strategies, and with PTopk-Miner<sub>dc</sub> and DCFinder for efficiency although they mine a special case of REEs.

We conducted experiments on a cluster of up to 21 virtual machines (one for the coordinator), each powered by 32GB RAM and 2 processors with 3.10 GHz. We ran the experiments 3 times, and report the average here. We do not include the time of loading datasets and constructing auxiliary structure, *i.e.*, PLI for all algorithms.

**Experimental results.** We next report our findings.

**Exp-1: Scalability.** We first evaluated the scalability of PTopk-Miner and Anytime-Miner vs. PTopk-Miner<sub>nop</sub>, PTopk-Miner<sub>dc</sub> and DCFinder, by varying (1) the number  $n$  of machines, (2) the support threshold  $\sigma$ , (3) the confidence threshold  $\delta$ , (4) the parameter  $k$  in top- $k$  discovery, (5) the number of relation atoms in REEs, (6) the rounds of top- $k$  discovery, and (7) the size of synthetic data. Unless stated explicitly, the default setting is  $n = 20$ ,  $\sigma = 10^{-6} \cdot |\mathcal{D}|^2$ ,

Name	Type	#tuples	#attributes	#relations
Adult [38, 46, 54]	real-life	32,561	15	1
Airport [46, 54]	real-life	55,113	18	1
Hospital [8, 12, 46, 54]	real-life	114,919	15	1
Inspection [46, 54, 60]	real-life	170,000	19	1
NCVoter [38, 46, 54]	real-life	1,681,617	12	1
DBLP [70]	real-life	1,799,559	18	3
Tax [8, 12, 17, 46, 54]	synthetic	10,000,000	15	1

**Table 3: Dataset statistic**

$\delta = 0.75$ ,  $k = 10$ , and the default number of relation atoms in REEs is set to 2, for a fair comparison with PTopk-Miner<sub>dc</sub> and DCFinder that are restricted to bi-variable discovery. For large datasets DBLP, NCVoter and Tax, we included 5 predicates in RHS, which is typical in an application. We adopted the combined predicate processing order and used 4 objective measures and 1 subjective measures. For the lack of space, we only show the results on some of the datasets; the results on the other datasets are consistent.

**Varying  $n$ .** We varied the number  $n$  of machines from 4 to 20. As shown in Figures 4(a) and 4(b), (a) PTopk-Miner scales well with the increase of machines: it is 3.7 times faster when  $n$  varies from 4 to 20. (b) PTopk-Miner is feasible in practice. It takes 667s on DBLP when  $n = 20$ , as opposed to 8726s by DCFinder, although PTopk-Miner discovers REEs that carries ML predicates and relation atoms across multiple tables, which are not supported by DCFinder. (c) PTopk-Miner is 2.89 and 11.90 times faster than PTopk-Miner<sub>nop</sub> and DCFinder on average, respectively, up to 3.29 and 13.11 times. This verifies that our pruning strategies for finding top- $k$  interesting rules effectively reduce the execution time. (d) Although PTopk-Miner<sub>dc</sub> is slightly faster than PTopk-Miner, it only discovers DCs, a special case of REEs. (e) Anytime-Miner is slightly slower than PTopk-Miner, since it has to maintain more rules in the heap and conduct implication checking to remove redundant rules (Section 5). Nonetheless, its advantage is evident when the users continuously require the next top- $k$  results, as will be seen shortly.

**Varying  $\sigma$ .** Varying the support threshold  $\sigma$  from  $10^{-1} \cdot |\mathcal{D}|^2$  to  $10^{-8} \cdot |\mathcal{D}|^2$ , we report the results in Figures 4(c) and 4(d). As expected, all algorithms take longer when  $\sigma$  is smaller since they need to examine more candidates, *e.g.*, PTopk-Miner<sub>nop</sub> is 8.2 times faster on Inspection when  $\sigma$  changes from  $10^{-1} \cdot |\mathcal{D}|^2$  to  $10^{-8} \cdot |\mathcal{D}|^2$ . Nevertheless, PTopk-Miner is faster than PTopk-Miner<sub>nop</sub> and DCFinder under all values of  $\sigma$ , which is consistent with Figures 4(a) and 4(b). In particular, PTopk-Miner is less sensitive to  $\sigma$ , since it checks much less REEs than PTopk-Miner<sub>nop</sub> due to its pruning strategies. Anytime-Miner enjoys a similar trend as PTopk-Miner.

**Varying  $\delta$ .** Varying the confidence threshold  $\delta$  from 0.75 to 0.95, we report the results in Figures 4(e) and 4(f). As shown there, (a) almost all algorithms become faster given a smaller  $\delta$ , *e.g.*, PTopk-Miner and Anytime-Miner are 1.10 times and 1.11 times faster, respectively, when  $\delta$  varies from 0.95 to 0.75. This is because a higher  $\delta$  indicates REEs with fewer violations, and this stronger constraint leads to more minimal REEs to be checked. (b) PTopk-Miner still consistently outperforms PTopk-Miner<sub>nop</sub> and DCFinder.

**Varying  $k$ .** Varying  $k$  from 1 to 100, Figure 4(g) shows that PTopk-Miner<sub>nop</sub> and DCFinder are indifferent to the values of  $k$ , since (a) PTopk-Miner<sub>nop</sub> mines all minimal REEs regardless of  $k$ , and (b)

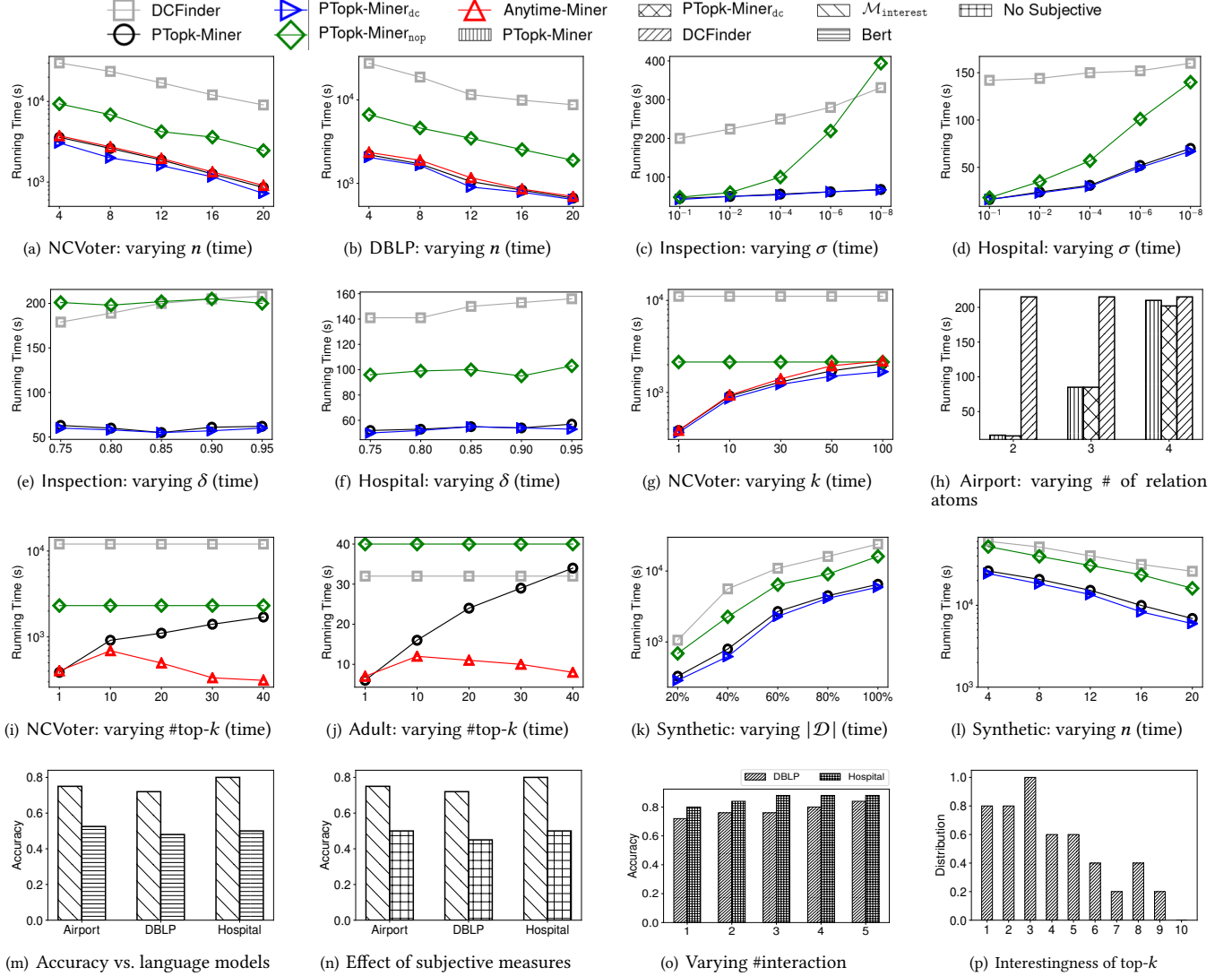


Figure 4: Performance evaluation

DCFinder discovers all minimal DCs subject to  $\sigma$  and  $\delta$ . In contrast, PTopk-Miner takes much less time than the two due to its pruning strategies for top- $k$  discovery; its cost increases slightly when  $k$  gets larger since more REEs have to be checked. Anytime-Miner is also faster than PTopk-Miner<sub>nop</sub> and DCFinder, by adopting a lazy evaluation strategy and skipping unnecessary REEs expansions.

**Varying # of relation atoms.** In Figure 4(h), we varied the number of relation atoms from 2 to 4. We observe that it takes longer to discover REEs with more relation atoms, as expected. Nonetheless, PTopk-Miner is still faster than DCFinder that mines bi-variable DCs, by 5.66 times on average and up to 13.43 times, even when PTopk-Miner mines more complicated REEs with 4 relation atoms.

**Varying #top- $k$ .** Fixing  $k=10$ , we varied the number #top- $k$  of top- $k$  results that users wish to see. Different from the lazy evaluation of Anytime-Miner, when users continue to find the next top- $k$ , PTopk-Miner is executed with an increased value of  $k$  and an increased heap size so that it exactly returns the desired results. For instance, when #top- $k$  = 4, PTopk-Miner discovers top-40 REEs.

Results on NCVoter and Adult are shown in Figures 4(i) and 4(j), respectively, when varying #top- $k$  from 1 and 4. For the first top- $k$  REEs, there is no big difference in the runtime between Anytime-Miner and PTopk-Miner. However, the advantage of Anytime-Miner over PTopk-Miner is more evident when the users want to see more top- $k$  results, e.g., when one asks for the 4th top-10 REEs, Anytime-Miner is 5.46 times and 4.25 times faster than PTopk-Miner on NCVoter and Adult, respectively. Intuitively, this is because Anytime-Miner maintains partial results and is more efficient to resume the discovery when needed. After Anytime-Miner accumulates sufficient partial results, its runtime may even decrease for larger #top- $k$ . This shows the effectiveness of Anytime-Miner when the users continue want to explore the next top- $k$  rules.

Using large Tax synthetic data  $\mathcal{D}$  (15 attributes and 10M tuples), we tested the impact of the size  $|\mathcal{D}|$  and the number  $n$  of machines.

**Varying  $|\mathcal{D}|$  (synthetic).** In Figure 4(k), we varied the scaling factor of  $\mathcal{D}$  from 20% to 100%, i.e., we changed the number of relations and tuples per relation from 2 million to 10 million. As shown there, all

algorithms take longer, as expected. Consistent with the previous test results, PTopk-Miner still outperforms PTopk-Miner<sub>nop</sub> and DCFinder. It takes 1.8h when  $\mathcal{D}$  has 10M tuples, as opposed to 4.4h and 6.5h by PTopk-Miner<sub>nop</sub> and DCFinder, respectively.

*Varying  $n$  (synthetic).* Fixing the data size  $|\mathcal{D}|$  as 10M, we varied the number  $n$  of machines from 4 to 20. The results are reported in Figure 4(l) and are consistent with those in Figures 4(a) and 4(b). PTopk-Miner is 3.8 times faster when  $n$  varies from 4 to 20.

**Exp-2: Effectiveness of the interestingness model.** In this set of experiments, we studied the effectiveness of our interestingness model, by (1) comparing the accuracy of  $\mathcal{M}_{\text{interest}}$  against the state-of-the-art language models, (2) evaluating the effectiveness of subjective measures, and (3) varying the rounds of interaction.

*Accuracy vs. language models.* We compared our interestingness model against Bert [13], which is implemented as the binary classification model with pairs of rules as input. We constructed a test set of rule pairs and asked the user to label the relative interestingness for each rule pair. Then, the model accuracy is measured by the percentage of rule pairs whose relative interestingness is correctly identified by the model. As reported in Fig. 4(m), our model consistently outperforms state-of-the-art language models in accuracy on all datasets, e.g., on Hospital, our accuracy is 0.80, 30% higher than the language model. This is mainly due to unique features such as token segmentation and permutation invariant of logic rules, which are unable to be captured by existing language models.

*Effectiveness of subjective measures.* We next studied the effectiveness of subjective measures by comparing the accuracy of our interestingness model with and without subjective measures. As shown in Fig. 4(n), the results verify that objective measures alone cannot meet the users' need well since different users have diverse preference, e.g., on DBLP, introducing subjective measures improves the accuracy by 27%, which clearly demonstrates the effectiveness of incorporating subjective measures in discovery.

*Varying #interaction.* To verify the usefulness of active learning, we report the accuracy by varying the number  $\#r$  of rounds of interaction in Figure 4(o). With larger  $\#r$ , the accuracy increases, e.g., after 5 rounds, the accuracy changes from 0.72 to 0.84 on DBLP. This said, after 5 rounds, the accuracy gets stable since the model has accumulated enough training data. We find that it typically requires 5 rounds of interactions to achieve a stable accuracy.

*Interestingness of top- $k$  REEs.* In this set of experiment, we mined top-10 REEs using PTopk-Miner and asked the users to label at most 5 of them which they think are interesting. As shown Figure 4(p), a rule ranked higher on the returned list is indeed more favored by the users. More specifically, the rules ranked first to fifth on the list are labeled as interesting by 76% of the users on average, as opposed to 24% for the rules ranked sixth to tenth. The user ranking justifies the semantic of top- $k$  interesting rule discovery.

**Exp-3: Effectiveness of PTopk-Miner.** We manually examined REEs discovered by PTopk-Miner from DBLP and Airport. Below are three REEs with support above 10 and confidence above 0.85.

(1)  $\text{Author}_{\text{DBLP}}(t_0) \wedge \text{Author}_{\text{DBLP}}(t_1) \wedge \mathcal{M}_{\text{Bert}}(t_0.\text{Affiliate}, t_1.\text{Affiliate}) \wedge \mathcal{M}_{\text{ditto}}(t_0, t_1) \rightarrow t_0.\text{Name} = t_1.\text{Name}$ . This REE

states that if two authors have similar affiliations and they semantically match, then they have the same name. The rule employs ML models for similarity checking (Bert) and ER (ditto).

(2)  $\text{Paper}_{\text{DBLP}}(t_0) \wedge \text{Paper}_{\text{DBLP}}(t_1) \wedge \mathcal{M}_{\text{Bert}}(t_0.\text{title}, t_1.\text{title}) \wedge t_0.\text{year} = t_1.\text{year} \wedge t_0.\text{venue} = t_1.\text{venue} \rightarrow \mathcal{M}_{\text{ditto}}(t_0.\bar{A}, t_1.\bar{A})$ , where  $\bar{A}$  is the set of all attributes in relation Paper of DBLP. This REE explores the possibility of interpreting an ML predication for ER by logic predicates and ML models. It says that two papers are predicted to match by an ML model because they have matching venue, year and title attributes. Note that the REE interprets the ML prediction of  $\mathcal{M}_{\text{ditto}}$  with another ML predicate  $\mathcal{M}_{\text{Bert}}$ .

(3)  $\text{Airport}(t_0) \wedge \text{Airport}(t_1) \wedge \text{Airport}(t_2) \wedge t_0.\text{iso\_region} = t_2.\text{iso\_region} \wedge t_1.\text{iata\_code} = t_2.\text{iata\_code} \wedge t_0.\text{municipality} = t_1.\text{municipality} \wedge t_0.\text{longitude\_deg} = t_2.\text{longitude\_deg} \rightarrow t_0.\text{iso\_region} = t_1.\text{iso\_region}$ . This REE collectively checks consistency by using three relation atoms. Current DC discovery algorithms do not support rules with more than two relation atoms.

(4)  $\text{Paper}_{\text{DBLP}}(t_0) \wedge \text{Paper}_{\text{DBLP}}(t_1) \wedge t_0.\text{year} \neq t_1.\text{year} \rightarrow t_0.\text{id} \neq t_1.\text{id}$ . It distinguishes papers that are published in different years.

**Summary.** We find the following. (1) Top- $k$  REEs discovery speeds up PTopk-Miner<sub>dc</sub> and DCFinder by 2.6 and 8.4 times on average, respective, up to 4.2 and 13.1 times. When  $n = 20$ , it takes less than 670s to mine top-10 REEs from DBLP that has 3 relations, 18 attributes and 1.8M tuples, as opposed to 3043s and 8726s by PTopk-Miner<sub>nop</sub> and DCFinder, respectively. (2) PTopk-Miner scales well with various parameters; e.g., when the support threshold  $\sigma$  decreases on Inspection, PTopk-Miner and Anytime-Miner take only slightly longer. (3) PTopk-Miner is parallelly scalable: on average, it is 3.6 times faster when the number  $n$  of machines varies from 4 to 20. (4) The lazy evaluation strategy of Anytime-Miner is effective when the users continuously want the next top- $k$  results: Anytime-Miner is 4.85 times faster than PTopk-Miner, when the users want the 4th top-10 REEs. (5) Our pruning strategies are effective, e.g., reducing the the runtime of PTopk-Miner by 2.89 times on the Tax data of 10M tuples. (6) Our interestingness model is on average 25.5% more accurate than the state-of-the-art language models. In particular, the subjective measures introduced in the model improves the accuracy from 0.48 to 0.76. (7) PTopk-Miner is capable of finding truly interesting REEs from real-life data.

## 8 CONCLUSION

We have studied discovery of top- $k$  interesting rules. The novelty of the work consists of the following: (1) an interestingness model with both objective and subjective measures; (2) an active-learning method to learn the subjective model and weights of various measures; (3) a top- $k$  algorithm for discovering REEs, which subsume CFDs, DCs and MDs as special cases; (4) an anytime algorithm to continuously mine the next top- $k$  rules via lazy evaluation, and (5) parallelization of the algorithms with the parallel scalability. Our experimental study has shown that the method is promising.

One topic for future work is to study incremental top- $k$  rule discovery in response to updates to both users interest and datasets. Another topic is to integrate top- $k$  algorithm and data sampling, to further speed up the discovery process, with accuracy guarantees.

## REFERENCES

- [1] Ziawasch Abedjan, Patrick Schulze, and Felix Naumann. 2014. DFD: Efficient functional dependency discovery. In *CIKM*. 949–958.
- [2] Serge Abiteboul, Richard Hull, and Victor Vianu. 1995. *Foundations of Databases*. Addison-Wesley.
- [3] Marcelo Arenas, Leopoldo Bertossi, and Jan Chomicki. 1999. Consistent Query Answers in Inconsistent Databases. In *PODS*.
- [4] Zeinab Bahmani and Leopoldo E. Bertossi. 2017. Enforcing Relational Matching Dependencies with Datalog for Entity Resolution. In *FLAIRS*.
- [5] Laure Berti-Équille, Hazar Harmouch, Felix Naumann, Noël Novelli, and Saravanan Thirumuruganathan. 2018. Discovery of Genuine Functional Dependencies from Relational Data with Missing Values. *PVLDB* 11, 8 (2018), 880–892.
- [6] Leopoldo E. Bertossi, Solmaz Kolahi, and Laks V. S. Lakshmanan. 2013. Data Cleaning and Query Answering with Matching Dependencies and Matching Functions. *Theory Comput. Syst.* (2013).
- [7] Indrajit Bhattacharya and Lise Getoor. 2007. Collective entity resolution in relational data. *TKDD* (2007).
- [8] Tobias Bleifuß, Sebastian Kruse, and Felix Naumann. 2017. Efficient Denial Constraint Discovery with Hydra. *PVLDB* 11, 3 (2017), 311–323.
- [9] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165* (2020).
- [10] Fei Chiang and Renee Miller. 2008. Discovering Data Quality Rules. In *VLDB*.
- [11] Vassilis Christophides, Vasilis Efthymiou, Themis Palpanas, George Papadakis, and Kostas Stefanidis. 2021. An Overview of End-to-End Entity Resolution for Big Data. *ACM Comput. Surv.* 53, 6 (2021), 127:1–127:42.
- [12] Xu Chu, Ihab F. Ilyas, and Paolo Papotti. 2013. Discovering Denial Constraints. *PVLDB* (2013).
- [13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *NAACL-HLT*. 4171–4186.
- [14] Honghua Dong, Jiayuan Mao, Tian Lin, Chong Wang, Lihong Li, and Denny Zhou. 2019. Neural Logic Machines. In *ICLR*.
- [15] Muhammad Ebraheem, Saravanan Thirumuruganathan, Shafiq R. Joty, Mourad Ouzzani, and Nan Tang. 2018. Distributed Representations of Tuples for Entity Resolution. *PVLDB* (2018).
- [16] Wenfei Fan, Hong Gao, Xibei Jia, Jianzhong Li, and Shuai Ma. 2011. Dynamic constraints for record matching. *VLDB J.* 20, 4 (2011), 495–520.
- [17] Wenfei Fan, Floris Geerts, Xibei Jia, and Anastasios Kementsitsidis. 2008. Conditional functional dependencies for capturing data inconsistencies. *TODS* 33, 2 (2008), 6:1–6:48.
- [18] Wenfei Fan, Floris Geerts, Jianzhong Li, and Ming Xiong. 2011. Discovering conditional functional dependencies. *TKDE* 23, 5 (2011), 683–698.
- [19] Wenfei Fan, Floris Geerts, Nan Tang, and Wenyuan Yu. 2014. Conflict resolution with data currency and consistency. *J. Data and Information Quality* 5, 1-2 (2014), 6:1–6:37.
- [20] Wenfei Fan, Ping Lu, and Chao Tian. 2020. Unifying Logic Rules and Machine Learning for Entity Enhancing. *Science China Information Sciences* (2020).
- [21] Wenfei Fan, Chao Tian, Yanghao Wang, and Qiang Yin. 2021. Discrepancy Detection and Incremental Detection. *PVLDB* (2021).
- [22] Wenfei Fan, Xin Wang, Yinghui Wu, and Jingbo Xu. 2015. Association Rules with Graph Patterns. *PVLDB* 8, 12 (2015), 1502–1513.
- [23] Peter A Flach and Iztok Savnik. 1999. Database dependency discovery: a machine learning approach. *AI communications* 12, 3 (1999), 139–160.
- [24] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. 2007. Sparse inverse covariance estimation with the graphical lasso. *Biostatistics* 9, 3 (2007), 432–441.
- [25] Cheng Fu, Xianpei Han, Le Sun, Bo Chen, Wei Zhang, Suhui Wu, and Hao Kong. 2019. End-to-end multi-perspective matching for entity resolution. In *IJCAI*. 4961–4967.
- [26] Eve Garnaud, Nicolas Hanusse, Sofian Maabout, and Noël Novelli. 2014. Parallel mining of dependencies. In *HPCS*. IEEE, 491–498.
- [27] Gartner. 2018. How to create a business case for data quality improvement. <https://www.gartner.com/smarterwithgartner/how-to-create-a-business-case-for-data-quality-improvement/>.
- [28] Chang Ge, Ihab F. Ilyas, and Florian Kerschbaum. 2019. Secure Multi-Party Functional Dependency Discovery. *PVLDB* 13, 2 (2019), 184–196.
- [29] Liqiang Geng and Howard J. Hamilton. 2006. Interestingness measures for data mining: A survey. *ACM Comput. Surv.* 38, 3 (2006), 9.
- [30] Lise Getoor and Ashwin Machanavajjhala. 2012. Entity resolution: tutorial. *VLDB* (2012).
- [31] Lukasz Golab, Howard Karloff, Flip Korn, Divesh Srivastava, and Bei Yu. 2008. On generating near-optimal tableaux for conditional functional dependencies. *VLDB* (2008).
- [32] Alireza Heidari, Joshua McGrath, Ihab F. Ilyas, and Theodoros Rekatsinas. 2019. HoloDetect: Few-Shot Learning for Error Detection. In *SIGMOD*.
- [33] Ykä Huhtala, Juha Kärrkäinen, Pasi Porkka, and Hannu Toivonen. 1999. TANE: An efficient algorithm for discovering functional and approximate dependencies. *The computer journal* (1999).
- [34] Ykä Huhtala, Juha Kärrkäinen, Pasi Porkka, and Hannu Toivonen. 1999. TANE: An Efficient Algorithm for Discovering Functional and Approximate Dependencies. *Comput. J.* (1999).
- [35] Hanna Köpcke, Andreas Thor, and Erhard Rahm. 2010. Evaluation of entity resolution approaches on real-world match problems. *PVLDB* 3, 1-2 (2010), 484–493.
- [36] Ioannis Koumarelas, Thorsten Papenbrock, and Felix Naumann. 2020. MDedup: Duplicate detection with matching dependencies. *PVLDB* 13, 5 (2020), 712–725.
- [37] Sebastian Kruse and Felix Naumann. 2018. Efficient Discovery of Approximate Dependencies. *PVLDB* 11, 7 (2018), 759–772.
- [38] Sebastian Kruse and Felix Naumann. 2018. Efficient discovery of approximate dependencies. *VLDB* (2018).
- [39] Clyde P Kruskal, Larry Rudolph, and Marc Snir. 1990. A complexity theory of efficient parallel algorithms. *TCS* (1990).
- [40] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2019. Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942* (2019).
- [41] Bing Li, Wei Wang, Yifang Sun, Linhan Zhang, Muhammad Asif Ali, and Yi Wang. 2020. GraphER: Token-Centric Entity Resolution with Graph Convolutional Neural Networks. In *AAAI*. 8172–8179.
- [42] Weibang Li, Zhanhuai Li, Qun Chen, Tao Jiang, and Hailong Liu. 2015. Discovering functional dependencies in vertically distributed big data. In *WISE*. 199–207.
- [43] Weibang Li, Zhanhuai Li, Qun Chen, Tao Jiang, and Zhilei Yin. 2016. Discovering approximate functional dependencies from distributed big data. In *APWeb*. 289–301.
- [44] Yuliang Li, Jinfeng Li, Yoshihiko Suhara, AnHai Doan, and Wang-Chiew Tan. 2020. Deep Entity Matching with Pre-Trained Language Models. *arXiv preprint arXiv:2004.00584* (2020).
- [45] Yinhan Liu, Mylène Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).
- [46] Ester Livshits, Alireza Heidari, Ihab F. Ilyas, and Benny Kimelfeld. 2020. Approximate Denial Constraints. *PVLDB* 13, 10 (2020), 1682–1695.
- [47] Stéphane Lopes, Jean-Marc Petit, and Lotfi Lakhal. 2000. Efficient discovery of functional dependencies and Armstrong relations. In *EDBT*. Springer, 350–364.
- [48] Mohammad Mahdavi, Ziawasch Abedjan, Raul Castro Fernandez, Samuel Madden, Mourad Ouzzani, Michael Stonebraker, and Nan Tang. 2019. Raha: A Configuration-Free Error Detection System. In *SIGMOD*. 865–882.
- [49] Sidharth Mudgal, Han Li, Theodoros Rekatsinas, AnHai Doan, Youngchoon Park, Ganesh Krishnan, Rohit Deep, Esteban Arcaute, and Vijay Raghavendra. 2018. Deep learning for entity matching: A design space exploration. In *SIGMOD*.
- [50] Noel Novelli and Rosine Cicchetti. 2001. Fun: An efficient algorithm for mining functional and embedded dependencies. In *ICDT*. Springer, 189–203.
- [51] Thorsten Papenbrock, Jens Ehrlich, Jannik Marten, Tommy Neubert, Jan-Peer Rudolph, Martin Schönberg, Jakob Zwiener, and Felix Naumann. 2015. Functional dependency discovery: An experimental evaluation of seven algorithms. *PVLDB* 8, 10 (2015), 1082–1093.
- [52] Thorsten Papenbrock and Felix Naumann. 2016. A Hybrid Approach to Functional Dependency Discovery. In *SIGMOD*.
- [53] Eduardo H. M. Pena, Eduardo Cunha de Almeida, and Felix Naumann. 2019. Discovery of Approximate (and Exact) Denial Constraints. *PVLDB* 13, 3 (2019), 266–278.
- [54] Eduardo H. M. Pena, Eduardo Cunha de Almeida, and Felix Naumann. 2019. Discovery of Approximate (and Exact) Denial Constraints. *PVLDB* 13, 3 (2019), 266–278.
- [55] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proc. of NAACL*.
- [56] Meng Qu, Junkun Chen, Louis-Pascal A. C. Xhonneux, Yoshua Bengio, and Jian Tang. 2020. RNNLogic: Learning Logic Rules for Reasoning on Knowledge Graphs. *CoRR abs/2010.04029* (2020).
- [57] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI blog* 1, 8 (2019), 9.
- [58] Thomas C. Redman. 2016. Bad Data Costs the U.S. \$3 Trillion Per Year. Harvard Business Review. <https://hbr.org/2016/09/bad-data-costs-the-u-s-3-trillion-per-year>.
- [59] Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. In *EMNLP-IJCNLP*. 3980–3990.
- [60] Theodoros Rekatsinas, Xu Chu, Ihab F. Ilyas, and Christopher Ré. 2017. HoloClean: Holistic Data Repairs with Probabilistic Inference. *PVLDB* (2017).
- [61] Theodoros Rekatsinas, Xu Chu, Ihab F. Ilyas, and Christopher Ré. 2017. HoloClean: Holistic data repairs with probabilistic inference. *PVLDB* (2017).
- [62] Theodoros Rekatsinas, Manas Joglekar, Hector Garcia-Molina, Aditya G. Parameswaran, and Christopher Ré. 2017. SLIMFast: Guaranteed Results for

- Data Fusion and Source Reliability. In *SIGMOD*. 1399–1414.
- [63] Hemant Saxena, Lukasz Golab, and Ihab F Ilyas. 2019. Distributed discovery of functional dependencies. In *ICDE*. IEEE, 1590–1593.
  - [64] Hemant Saxena, Lukasz Golab, and Ihab F Ilyas. 2019. Distributed implementations of dependency discovery algorithms. *PVLDB* 12, 11 (2019), 1624–1636.
  - [65] Philipp Schirmer, Thorsten Papenbrock, Ioannis Koumarelas, and Felix Naumann. 2020. Efficient Discovery of Matching Dependencies. *RODS* 45, 3 (2020), 1–33.
  - [66] Philipp Schirmer, Thorsten Papenbrock, Sebastian Kruse, Felix Naumann, Dennis Hempfing, Torben Mayer, and Daniel Neuschäfer-Rube. 2019. DynFD: Functional Dependency Discovery in Dynamic Datasets. In *EDBT*.
  - [67] Shaoyun Shi, Hanxiong Chen, Weizhi Ma, Jiaxin Mao, Min Zhang, and Yongfeng Zhang. 2020. Neural Logic Reasoning. In *CIKM*. 1365–1374.
  - [68] Rohit Singh, Venkata Vamsikrishna Meduri, Ahmed K. Elmagarmid, Samuel Madden, Paolo Papotti, Jorge-Arnulfo Quiané-Ruiz, Armando Solar-Lezama, and Nan Tang. 2017. Synthesizing Entity Matching Rules by Examples. *PVLDB* 11, 2 (2017), 189–202.
  - [69] Shaoxu Song and Lei Chen. 2009. Discovering matching dependencies. In *CIKM*.
  - [70] Jie Tang, Jing Zhang, Limin Yao, Juanzi Li, Li Zhang, and Zhong Su. 2008. Arnet-Miner: Extraction and Mining of Academic Social Networks. In *KDD'08*. 990–998.
  - [71] Leslie G. Valiant. 1990. A Bridging Model for Parallel Computation. *Commun. ACM* 33, 8 (1990), 103–111.
  - [72] G. I. Webb and S. Zhang. 2005. k-Optimal Rule Discovery. *Data Mining and Knowledge Discovery* 10, 1 (2005), 39–79.
  - [73] Renzhi Wu, Sanya Chaba, Saurabh Sawlani, Xu Chu, and Saravanan Thirumuganathan. 2020. ZeroER: Entity Resolution using Zero Labeled Examples. In *SIGMOD*. 1149–1164.
  - [74] Catharine M. Wyss, Chris Giannella, and Edward L. Robertson. 2001. FastFDs: A Heuristic-Driven, Depth-First Algorithm for Mining Functional Dependencies from Relation Instances - Extended Abstract. In *DaWaK*.
  - [75] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. In *NeurIPS*. 5753–5763.
  - [76] H Yao, H Hamilton, and C Butz. 2002. Fd\_mine: discovering functional dependencies in a database using equivalences, Canada. In *IEEE ICDM*. 1–15.
  - [77] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Ruslan R Salakhutdinov, and Alexander J Smola. 2017. Deep Sets. In *Advances in Neural Information Processing Systems*. 3391–3401.
  - [78] Yunjia Zhang, Zhihan Guo, and Theodoros Rekatsinas. 2020. A Statistical Perspective on Discovering Functional Dependencies in Noisy Data. In *SIGMOD*. 861–876.
  - [79] Yunjia Zhang, Zhihan Guo, and Theodoros Rekatsinas. 2020. A Statistical Perspective on Discovering Functional Dependencies in Noisy Data. In *SIGMOD*. 861–876.
  - [80] Chen Zhao and Yeye He. 2019. Auto-EM: End-to-end Fuzzy Entity-Matching using Pre-trained Deep Models and Transfer Learning. In *WWW*. 2413–2424.