

Knowledge Refinement via Rule Selection

Phokion G. Kolaitis^{1,2}, Lucian Popa², Kun Qian²

¹UC Santa Cruz, ²IBM Research – Almaden

AAAI-19



Introduction

- ▶ Rules are ubiquitous in AI and computer science
- ▶ Rules as typically specified by **Horn** formulas
$$\forall x, y, z(\text{PARENT}(x, z) \wedge \text{PARENT}(y, z) \rightarrow \text{SIBLING}(x, y))$$
- ▶ Rules are extensively used in different areas, including:
 - ▶ Logic Programming
$$\text{SIBLING}(x, y) : - \text{PARENT}(x, z), \text{PARENT}(y, z)$$
 - ▶ Data Mining (Mining Association Rules)
 - ▶ Data Exchange and Data Integration
 - ▶ Entity Resolution.
- ▶ **Central Problem:** From a large number of possible rules, select a subset that **best captures** the task at hand.

Concrete Scenario - Refining Entity Resolution Rules

- Discover a rule-based model that identifies the same paper across datasets

DBLP: (id, title, authors, venue, year), ACM: (id, title, authors, venue, year)

- There is a large number of possible entity resolution rules:

r1: DBLP(p, t , a , v , y1)	\wedge ACM(q, t , a , v , y2)	\rightarrow Same(p, q)
r2: DBLP(p, t , a1, v , y)	\wedge ACM(q, t , a2, v , y)	\rightarrow Same(p, q)
r3: DBLP(p, t1, a1, v , y)	\wedge ACM(q, t2, a2, v , y)	\rightarrow Same(p, q)
...

- Rules behave differently.
 - r1 has high precision but low coverage.
 - r3 has high coverage but low precision.
- **Goal:** Find an optimal subset of the given rules that achieves high precision and high coverage.

Concrete Scenario - Refining Entity Resolution Rules

- Given a set of available rules and ground truth data:

$r1: \text{DBLP}(p, \mathbf{t}, \mathbf{a}, \mathbf{v}, y1) \wedge \text{ACM}(q, \mathbf{t}, \mathbf{a}, \mathbf{v}, y2) \rightarrow \text{Same}(p, q)$
 $r2: \text{DBLP}(p, \mathbf{t}, a1, \mathbf{v}, \mathbf{y}) \wedge \text{ACM}(q, \mathbf{t}, a2, \mathbf{v}, \mathbf{y}) \rightarrow \text{Same}(p, q)$
 $r3: \text{DBLP}(p, t1, a1, \mathbf{v}, y1) \wedge \text{ACM}(q, t2, a2, \mathbf{v}, y2) \rightarrow \text{Same}(p, q)$
...

DBLP records	ACM records
(d ₁ , "rule selection", "Jone Doe", AAAI, 2019)	(a, "rule selection", "J. Doe", AAAI, 2019)
(d ₂ , "Invited Talk", "Qiang Yang", AAAI, 2019)	(b, "Invited Talk", "Yu Zheng", AAAI, 2019)
(d ₃ , "Keynote", "Jane Doe", SIGMOD, 2018)	(c, "Keynote", "Jane Doe", SIGMOD, null)

Same
(d ₁ ,a)
(d ₃ ,c)

- Which model minimizes the error over the ground truth?
 - $\{r1\}$: 0 false positives, 1 false negative.
 - $\{r2\}$: 1 false positives, 1 false negative.
 - $\{r3\}$: 2 false positive, 0 false negatives.
 - Both $\{r1\}$ and $\{r1, r2\}$ achieve minimum error.

Basic Concepts: Rules, Data Examples, Rule Evaluation

- ▶ We consider rules of the form $\forall \mathbf{x}(\psi(\mathbf{x}) \rightarrow P(\mathbf{x}))$, where
 - ▶ the **premise** $\psi(\mathbf{x})$ is a conjunction of atoms over a schema **S**;
 - ▶ the **conclusion** $P(\mathbf{x})$ is a atom over a schema **T** disjoint from **S**.
- ▶ A data example is a pair (I, J) , where I is a premise instance I over **S** and a J is a conclusion instance over **T**.
- ▶ If \mathcal{C} is a set of rules, then $\text{Eval}(\mathcal{C}, I) = \bigcup_{r \in \mathcal{C}} \text{Eval}(r, I)$ (i.e., evaluate the premise of each rule on I and populate the conclusion).

Premise Instance I

DBLP records	ACM records
(d ₁ , "rule selection", "Jone Doe", AAAI, 2019)	(a, "rule selection", "J. Doe", AAAI, 2019)
(d ₂ , "Invited Talk", "Qiang Yang", AAAI, 2019)	(b, "Invited Talk", "Yu Zheng", AAAI, 2019)
(d ₃ , "Keynote", "Jane Doe", SIGMOD, 2018)	(c, "Keynote", "Jane Doe", SIGMOD, null)

Conclusion Instance J

Same
(d ₁ , a)
(d ₃ , c)

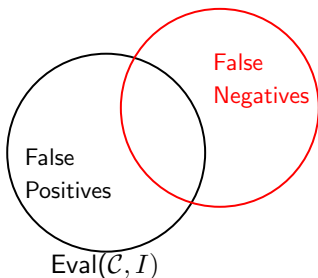
$$\text{Eval}(\{r_1\}, I) = \{(d_3, c)\}, \text{Eval}(\{r_1, r_2\}, I) = \{(d_1, a), (d_2, b), (d_3, c)\}$$

False Positive and False Negative Errors

Definition (Errors)

Given a set \mathcal{C} of rules and a data example (I, J) , we consider the following two types of errors of \mathcal{C} w.r.t. (I, J) :

- ▶ False positive errors: $\text{Eval}(\mathcal{C}, I) \setminus J$;
- ▶ False negative errors: $J \setminus \text{Eval}(\mathcal{C}, I)$.



The Min Rule-Selection Problem

Definition ($\text{MIN RULE-SELECT}_{\text{FP}+\text{FN}}$)

Input: A set \mathcal{C} of rules and a data example (I, J) .

Goal: Find a subset $\mathcal{C}^* \subseteq \mathcal{C}$ such that the sum of the number of the false positive errors and the number of false negative errors of \mathcal{C}^* with respect to (I, J) is minimized.

Definition ($\text{MIN RULE-SELECT}_{\text{FP}}$)

Input: A set \mathcal{C} of rules and a data example (I, J) with $J \subseteq \text{Eval}(\mathcal{C}, I)$.

Goal: Find a subset $\mathcal{C}^* \subseteq \mathcal{C}$ such that the number of false negative errors is zero and the number of false positive errors of \mathcal{C}^* with respect to (I, J) is minimized.

Note: $\text{MIN RULE-SELECT}_{\text{FN}}$ has a trivial solution (the entire set of rules).

Complexity of MIN RULE-SELECT Problems

Theorem

Both $\text{MIN RULE-SELECT}_{\text{FP}+\text{FN}}$ and $\text{MIN RULE-SELECT}_{\text{FP}}$ are NP-hard optimization problems.

Hint of Proof: Reduction from SET COVER.

Question:

Are there polynomial-time approximation algorithms for $\text{MIN RULE-SELECT}_{\text{FP}+\text{FN}}$ and $\text{MIN RULE-SELECT}_{\text{FP}}$?

Approximation Properties of MIN RULE-SELECT_{FP}

Theorem

- ▶ MIN RULE-SELECT_{FP} is approximable within a factor of $2\sqrt{|\mathcal{C}| \log |J|}$, where $|\mathcal{C}|$ is the number of input rules and $|J|$ is the size of the conclusion instance J .
- ▶ Unless $P=NP$, for every $\epsilon > 0$, there is no polynomial time algorithm that approximates MIN RULE-SELECT_{FP} within a factor of $2^{\log^{1-\epsilon}(|\mathcal{C}|)}$, where $|\mathcal{C}|$ is as above.

Hint of Proof: Give approximation-preserving reductions between MIN RULE-SELECT_{FP} and the RED-BLUE SET COVER problem, which was studied by [Peleg, 2007].

Approximation Properties of MIN RULE-SELECT_{FP+FN}

Theorem

- ▶ MIN RULE-SELECT_{FP+FN} is approximable within a factor of $2\sqrt{(|\mathcal{C}| + |J|) \log |J|}$, where $|\mathcal{C}|$ is the number of input rules and $|J|$ is the size of the conclusion instance J .
- ▶ Unless $P=NP$, for every $\epsilon > 0$, there is no polynomial time algorithm that approximates MIN RULE-SELECT_{FP+FN} within a factor of $2^{\log^{1-\epsilon}(|J|)}$, where $|J|$ is as above.

Hint of Proof: Give approximation-preserving reductions between MIN RULE-SELECT_{FP+FN} and the POSITIVE-NEGATIVE PARTIAL SET COVER problem, which was studied by [Miettinen, 2008].

Bi-objective Optimization Problems

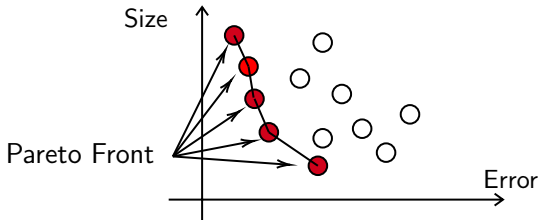
Question: What if we want to optimize both the **errors** and the **size** of rules simultaneously?

- Note that **error** and **size** are qualitatively incomparable quantities.

Answer: Consider **bi-objective** optimization problems.

Definition (Pareto Optimality)

Given a set \mathcal{C} of rules and a data example (I, J) , a subset $\mathcal{C}^* \subseteq \mathcal{C}$ is a **Pareto optimal solution** if there is no $\mathcal{C}' \subseteq \mathcal{C}$ with $\text{size}(\mathcal{C}') < \text{size}(\mathcal{C}^*)$ and $\text{error}(\mathcal{C}') < \text{error}(\mathcal{C}^*)$.



Complexity Results for Bi-objective Optimization

Theorem

- ▶ The following problem is coNP-complete: given an instance $K = \{\mathcal{C}, (I, J)\}$ of MIN RULE-SELECT_{FP} and a subset $\mathcal{C}^* \subseteq \mathcal{C}$, is \mathcal{C}^* a Pareto optimal solution?
- ▶ The following problem is DP-complete: given an instance $K = \{\mathcal{C}, (I, J)\}$ of MIN RULE-SELECT_{FP} and pair (s, e) of integers, is (s, e) on the Pareto front of K ?

Similar results hold for MIN RULE-SELECT_{FP+FN}.

Corollary

Unless $P = NP$, there is no polynomial-time algorithm for constructing the Pareto front.

Concluding Remarks

This work establishes foundations and complexity results for the rule-selection problem, which is fundamental in many areas of AI.

A summary of our results (more details can be found in the paper).

		FP	FP+FN
Rule-Select		NP-complete	NP-complete
Exact Rule-Select		DP-complete	DP-complete
Min Rule-Select	approx. upper	$2\sqrt{ \mathcal{C} \log J }$	$2\sqrt{(\mathcal{C} + J) \log J }$
	approx. lower	$2^{\log^{1-\epsilon}(\mathcal{C})}, \forall \epsilon > 0$	$2^{\log^{1-\epsilon}(J)}, \forall \epsilon > 0$
Pareto Optimal Solution		coNP-complete	coNP-complete
Pareto Front Membership		DP-complete	DP-complete
Bi-level Optimal Solution		coNP-complete	coNP-complete
Bi-level Optimal Value		DP-complete	DP-complete

Directions for future work

- Approximating the Pareto front of the rule selection problem.
- Experimental evaluation.

Thank you

References I



Miettinen, P. (2008).

On the positive-negative partial set cover problem.

Information Processing Letters, 108(4):219 – 221.



Peleg, D. (2007).

Approximation algorithms for the label-covermax and red-blue set cover problems.

Journal of Discrete Algorithms, 5(1):55 – 64.