

PRACTICAL USES OF EXISTENTIAL RULES IN KNOWLEDGE REPRESENTATION

Part 4: Practical Applications of Rules

David Carral,¹ Markus Krötzsch,¹ and Jacopo Urbani²

1. TU Dresden

2. Vrije Universiteit Amsterdam

ECAI, September 4, 2020

Outline

Goal

Show some example where either **rules** or **related ideas** were crucial to achieve the state of the art

- PLP
- Data integration
- Stream reasoning

Take-home message

1. Rules can be used also in scenarios where not everything is definite
2. A **declarative** approach is (often) intuitive and **decreases** the development time
3. Developing robust tools is fundamental

1st Scenario: Probabilistic Logic Programming

PLP

How can we perform logic-based reasoning in an uncertain domain?

PLP

Probabilistic Logic Programming (PLP): Formalisms to combine logic and probability for reasoning in uncertain domains.

Basic idea: Reason over facts which may be true with a certain probability

State of the art

Several PLP formalisms exist. **ProbLog** (Raedt, Kimmig, and Toivonen 2007) is one of the most popular ones

ProbLog

Definition

A ProbLog program \mathcal{P} is a triple $(\mathcal{R}, \mathcal{F}, \pi)$ where \mathcal{R} is set of (function-free) rules, \mathcal{F} is a set of facts and $\pi : \mathcal{F} \rightarrow [0, 1]$ is the function that labels facts with probabilities.

Key problem

Given \mathcal{P} and query q as input, what is $Pr(q)$ (the probability of q)?

General Approach

It has been shown that computing $Pr(q)$ can be expressed using Weighted Model Counting (WMC) over weighted logical formulas (Vlasselaer et al. 2016)

The Grounding Problem

ProbLog2, a state-of-the-art engine, proceeds as follows:

1. Find relevant **ground** program for q with backward chaining
2. Execute a custom implementation of fixpoint operator $T_{\mathcal{P}}$:
 - $T_{\mathcal{P}}$ proceeds bottom-up, akin to chase procedures
 - $T_{\mathcal{P}}$ incrementally computes, for each inferred fact f , a propositional formula λ_f which “remembers” all the possible ways f can be inferred
3. After $T_{\mathcal{P}}$ has finished, it computes WMC for λ_q

Problem

Grounding can be a major performance bottleneck with large knowledge bases

Datalog to the rescue

Some ideas developed for Datalog are useful here (Tsamoura, Gutiérrez-Basulto, and Kimmig 2020)

First idea

Don't ground \mathcal{P} with backward chaining. Rewrite it with **magic sets** (Bancilhon et al. 1985)

Second idea

Apply **semi-naïve evaluation** (Abiteboul, Hull, and Vianu 1995) while computing $T_{\mathcal{P}}$ to reduce the number of duplicates

Magic sets

Consider database I and program P . Our goal is to answer query Q

Idea

The main idea is to rewrite P into P' where additional **magic** relations restrict the derivations to facts relevant for answering Q

Magic sets

Consider database I and program P . Our goal is to answer query Q

Example 1

Consider the rules below and assume we want to answer $Q = \text{lives}(\text{linda}, X)$

$$\text{married}(X, Y), \text{lives}(X, Z) \rightarrow \text{lives}(Y, Z) \quad (r_1)$$

$$\text{married}(X, Y) \rightarrow \text{married}(Y, X) \quad (r_2)$$

The rewriting procedure produces the program

$$\text{mgc}_1(Y), \text{married}(X, Y), \text{lives}(X, Z) \rightarrow \text{lives}(Y, Z) \quad (r_3)$$

$$\text{mgc}_1(X) \rightarrow \text{mgc}_2(X) \quad (r_4)$$

$$\text{mgc}_2(Y), \text{married}(X, Y) \rightarrow \text{married}(Y, X) \quad (r_5)$$

Then, we can reason on $I \cup \{\text{mgc}_1(\text{linda})\}$

Semi naïve evaluation

Semi naïve evaluation is a well-known technique to avoid the recomputation of duplicate derivation during the materialization

Naïve Evaluation

Input: Facts I , program P

```
1 while true do  
2    $J := I$ ;  
3   for  $r \in P$  do  
4     Let  $r$  be  $B \rightarrow H$   
5      $J := J \cup \{H\sigma \mid B\sigma \subseteq I\}$ ;  
6   if  $J = I$  then return  $J$  ;
```

Semi Naïve Evaluation

Input: Facts I , program P

```
1  $\Delta := I$ ;  
2 while true do  
3    $J := I$ ;  
4   for  $r \in P$  do  
5     Let  $r$  be  $B \rightarrow H$ ;  
6      $J := J \cup \{H\sigma \mid B\sigma \subseteq I \wedge B\sigma \cap \Delta \neq \emptyset\}$ ;  
7   if  $J = I$  then return  $J$ ;  
8    $\Delta := J \setminus I$ ;
```

New approach

Tsamoura et al. (2020) proposed a new procedure:

1. ~~Find relevant **ground** program for q with backward chaining.~~ Use Magic Set to obtain a **non-ground** program
2. ~~Execute a custom implementation of fixpoint operator $T_{\mathcal{P}}$~~ Offload the computation to a chase engine (VLog):
 - Leverage semi-naïve evaluation
 - Introduce some rules to compute formulas (called λ -transformation)
3. After $T_{\mathcal{P}}$ has finished, compute WMC for λ_q

Impact

The new procedure removes the need for grounding, which was a performance bottleneck

Performance improvement

Some key results from (Tsamoura, Gutiérrez-Basulto, and Kimmig 2020)

- The runtime of query answering was two order of magnitude and 25% faster than ProbLog2 in the best and worst cases, respectively
- VLog enabled the computation on much larger DBs than what was possible before

Lesson learned

Well-known ideas developed for rule-based query answering can be re-used as-is for other problems as well

2nd Scenario: Entity Resolution

Problem

Scientific advancement requires an extensive analysis of prior knowledge in the literature, but this is **time consuming**

AI can help!

Long-term vision: Develop an accurate and large-scale KB of scientific knowledge

A KB of Scientific Knowledge

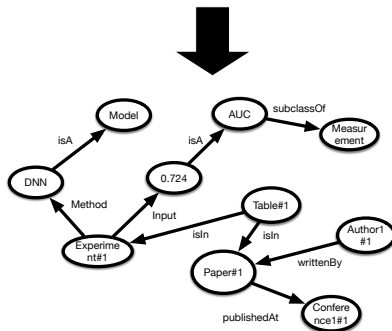
valuable experimental knowledge

l_1 - l_2	#S	# l_1 -W	# l_2 -W	# l_1 -V	# l_2 -V	Type	Example Words
en-de	1.9M	55M	52M	40k	50k	Offensive	disgusting, filthy, nasty, rude, horrible, terrible, awful, worst, idiotic, stupid, dumb, ugly, etc.
en-fr	2.0M	50M				Non-offensive	help, love, respect, believe, congrats, hi, like, great, fun, nice, neat, happy, good, best, etc.
en-es	1.9M	49M					

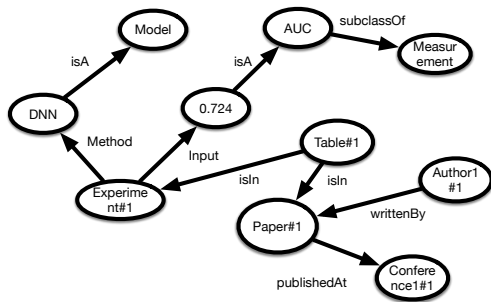
Distribution Parameters		Gaussian $\mu \in \mathbf{R}, \sigma^2 > 0$	
Models	AUC	RIG	$[0, 1]$
DNN	0.724	0.094	0
miDNN	0.747	0.119	$\mathbf{R}, b > 0$
miRNN	0.765	0.141	$b \in \mathbf{R}$
miRNN+attention	0.774	0.156	

Model	AUC	RIG	Latency (ms)
DNN	0.724	0.094	100
miDNN	0.747	0.119	150
miRNN	0.765	0.141	200
miRNN+attention	0.774	0.156	250

Figure 6: The word latency increase with respect to base size.
 Table 2: The AUC increase in A/B test.
 Table 3: The RIG increase in A/B test.



Advantages



Potential use cases:

- Retrieve experimental results with entity-based search
- Exploit co-authorship networks
- Identify potential inconsistencies across papers

Tab2Know: General pipeline

Tab2Know is a recent work to construct a KB from tables in scientific papers (Kruit, He, and Urbani 2020)

Key features:

- Heuristic-based methods to recognize and extract tables from PDFs
- Machine learning models to predict the type of tables and columns
- **Weak supervision** with SPARQL queries to counter the problem of lack of training data
- **(Focus of today)** logic-based reasoning for **entity resolution**

Tab2Know: General pipeline

From (Kruit, He, and Urbani 2020)

TABLE I. RANKING OF SUBMITTED METHODS TO TASK 1.1

Method Name	Recall (%)	Precision (%)	F-score
USTB_TexStar	82.38	93.83	87.74
TH-TextLoc	75.85	86.82	80.96
I2R_NUS_FAR	71.42	84.17	77.27
Baseline	69.21	84.94	76.27
Text Detection [15], [16]	73.18	78.62	75.81
I2R_NUS	67.52	85.19	75.34
BDTD_CASIA	67.05	78.98	72.53
OTCYMIST [7]	74.85	67.69	71.09
Inkam	52.21	58.12	55.00

Method Name	Recall (%)	Precision (%)	F-score
USTB_TexStar	82.38	93.83	87.74
TH-TextLoc	75.85	86.82	80.96
I2R_NUS_FAR	71.42	84.17	77.27
Baseline	69.21	84.94	76.27
Text Detection [15], [16]	73.18	78.62	75.81
I2R_NUS	67.52	85.19	75.34
BDTD_CASIA	67.05	78.98	72.53
OTCYMIST [7]	74.85	67.69	71.09
Inkam	52.21	58.12	55.00

Input: PDF Figure

APIs



Semantic Scholar

Ontology



SPARQL Queries



SPARQL Query 1
SPARQL Query 2
SPARQL Query 3
...

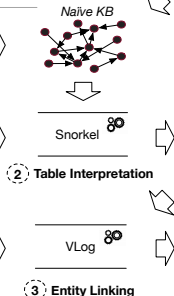
Rules



Rule 1
Rule 2
Rule 3
...

Assets

1 Table Extraction



2 Table Interpretation

3 Entity Linking

Table type classification

Method Name	Recall (%)	Precision (%)	F-score
USTB_TexStar	82.38	93.83	87.74
TH-TextLoc	75.85	86.82	80.96
I2R_NUS_FAR	71.42	84.17	77.27
Baseline	69.21	84.94	76.27
Text Detection [15],[16]	73.18	78.62	75.81
I2R_NUS	67.52	85.19	75.34
BDTD_CASIA	67.05	78.98	72.53
OTCYMIST [7]	74.85	67.69	71.09
Inkam	52.21	58.12	55.00

Header detection

Column type classification

Output: KB (with linked entities)

Entity Resolution

Entity resolution is the task of recognizing and linking entities across different tables. It is a well-known task in database literature (96+ papers between 2009-2014, see (Papadakis, Ioannou, and Palpanas 2020))

- Magellan (Konda et al. 2016)
- Deep Learning (Mudgal et al. 2018)
- Crowd-sourcing (Das et al. 2017)
- Embeddings (Cappuzzo, Papotti, and Thirumuruganathan 2020)
- ...

A declarative approach

Tab2Know's approach: Use (existential) rules!

TGDs

Used to create new entities from the cells

EGDs

Used to infer equality among the entities

Output

After reasoning is completed, entities are used to populate a KB

A declarative approach: TGDs

Two TGDs are used:

$$\text{type}(X, \text{Column}) \rightarrow \exists Y. \text{colEntity}(X, Y) \quad (r_1)$$

$$\text{type}(X, \text{Cell}) \rightarrow \exists Y. \text{cellEntity}(X, Y) \quad (r_2)$$

- Two types of entities are introduced. One describes columns, the other describes cells;
- Every cell is assigned to a entity; it is likely that the same entity is represented with multiple labeled nulls!

A declarative approach: EGDs

EGDs determines whether multiple cells refer to the same entity

$$ceNoTypLabel(X, L) \wedge ceNoTypLabel(Y, L) \rightarrow X \approx Y \quad (r_3)$$

$$eNoTypLabel(X, C, L), eNoTypLabel(Y, C, L) \rightarrow X \approx Y \quad (r_4)$$

$$eTableLabel(X, T, L), eTableLabel(Y, T, L) \rightarrow X \approx Y \quad (r_5)$$

$$eTypLabel(X, S, L), eTypLabel(Y, S, M), STR_EQ(L, M) \rightarrow X \approx Y \quad (r_6)$$

$$eAuthLabel(X, A, L), eAuthLabel(Y, A, M), STR_EQ(L, M) \rightarrow X \approx Y \quad (r_7)$$

- Special built-in predicates (*STR_EQ*) encode string similarities
- Other predicates include authors of the paper
- Program can be easily extended with other rules \rightarrow rapid KB construction

Preliminary results

Input

Approach was tested on a collection with 142k CS open-access papers and 73k tables (IJCAI, ECAI, etc.)

Key results

- Table interpretation superior than previous state-of-the-art approach (Yu et al. 2020)
- EGDs reduced number of “column” entities of 65% and of “cell” entities of 55%
- Every rule contributed by linking some entities
- On a sample of 541 entities, average precision was 97%

Lessons learned

1. A declarative approach is ideal for non-CS domain experts
2. Rules can be easily changed or adapted depending on the performance
3. VLog was scalable enough to perform rapid prototyping with large KGs
4. Support to built-in predicates was crucial

3rd Scenario: Stream Reasoning

A few of slides are a modified version of Harald Beck's ISWC17 presentation, used with permission

Motivation

Stream reasoning: add reasoning on top of stream processing. Central question: “**What is true now?**” (Margara et al. 2014)

- E.g. public transport: What are the current expected arrival times?
- Is there currently a good connection between two lines?

Semantic Web: RDF Stream Processing - SPARQL extensions: C-SPARQL, CQELS, SPARQL_{Stream}, . . . Typical: **Window operators** select snapshots of recent data

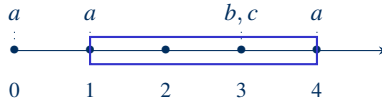
- Window examples: **[RANGE 3m], [TRIPLES 2]**

Goals & Challenges

- Goal: expressive stream reasoning solutions
 - (1) based on model-based semantics
 - (2) high performance
- Central challenge: **throughput vs. expressiveness**

LARS: A Logic for Analytic Reasoning over Streams

LARS (Beck, Dao-Tran, and Eiter 2018) is a logic-based framework to reason on streams



- Stream $S = (T, \nu)$
 - **Timeline** T closed interval in \mathbb{N} , $t \in T$ **time point**
 - **Evaluation** function $\nu : T \rightarrow 2^{\mathcal{A}}$ (sets of atoms)
- Window function w yields window $w(S, t) \subseteq S$
- Formulas ψ : evaluated on S at t

ψ	holds in S at t iff φ holds ...	Ex.: $S, 4 \models \psi$?
$\boxplus^w \varphi$	in $w(S, t)$ at t	
$\diamond \varphi$	at some time point $t' \in T$	$\boxplus^3 \diamond a$ ✓
$\square \varphi$	at all time points $t' \in T$	$\boxplus^3 \square a$ ✗
$@_{t'} \varphi$	at time point t' and $t' \in T$	$\boxplus^3 @_1 a$ ✓

Plain LARS

Observations

- Many practical problems do not need a multiple model semantics
- **Time-based** and **tuple-based windows** often suffice
- **Sliding windows** can be exploited for incremental reasoning

Plain LARS (Bazoobandi, Beck, and Urbani 2017)

Focus on **positive LARS programs** where for each rule $\alpha \leftarrow \beta_1, \dots, \beta_n$ we have:

- head α : atom a or $@_t a$
- body elements: $\beta_i ::= a \mid @_t a \mid \boxplus^w @_t a \mid \boxplus^w \diamond a \mid \boxplus^w \square a$

Consider **non-ground programs**, using substitutions due to available ground atoms, as usual

From LARS to Datalog

Observation

LARS rules can be rewritten into Datalog rules

- How do we represent time?
 - Increase arity of the relations, e.g., $P(X) \rightarrow P(X, T)$
- How can we translate LARS rules?
 - $@_S P(X)$ as $P(X, S)$
 - $\boxplus^2 \diamond P(X) \rightarrow Q(X)$ as $P(X, T) \rightarrow Q(X)$ and $P(X, T - 1) \rightarrow Q(X)$

Semi-naïve evaluation (SNE)

One key novelty of (Bazoobandi, Beck, and Urbani 2017) is to show how to replicate SNE in a stream

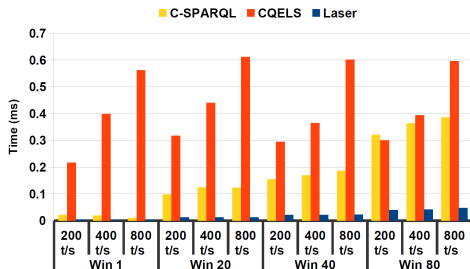
From LARS to Datalog

- For formula $\varphi = \alpha, \beta_i$ in any rule $\alpha \leftarrow \beta_1, \dots, \beta_n$, consider **annotated ground formulas** $\varphi\sigma_{[c,h]}$, where
 - $\varphi\sigma$ is the **ground instance** of φ due to **substitution** σ
 - $[c, h]$ is an **annotation** stating that $\varphi\sigma$ holds from **consideration time** c to **horizon time** h
- Horizon time can be extended in the future, e.g., at time point t , $\boxplus^3 \diamond p(a)$ can be annotated as $\boxplus^3 \diamond p(a)_{[t,t+3]}$
- When computing substitution σ for instantiating rule $B_1 \wedge B_2 \wedge \dots \wedge B_n \rightarrow H$ at time point t , at least one $B_i\sigma_{[c,h]}$ has $c = t$, i.e., has been derived at the current time point

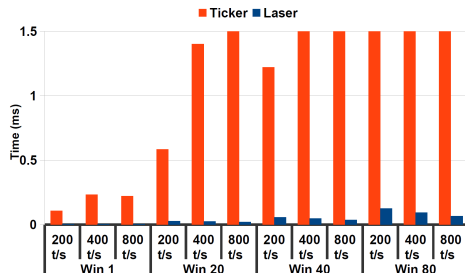
Laser: Implementation & Evaluation

Evaluation: Time per triple

- Compare to C-SPARQL, CQELS, and Ticker
- Micro benchmarks to test **(1)** $q(A, B) \leftarrow \boxplus^n \diamond p(A, B)$ (resp. \square); elementary data join; multiple rules; **(2)** small show case example requiring LARS features.
- Window sizes: 1s to 80s; stream rate: 200 to 800 triples/second



(1)



(2)

Lesson learned






- A good idea remains a good idea (even if is old)
- ... but it might need to be properly implemented

To conclude

We have described cases where rules turned out to be very useful

- In some scenarios, existential quantification was necessary (data integration). In others, Datalog rules were enough (PLP, stream reasoning)
- Sometimes, the tools could be directly used (data integration). In other cases, some modifications are required (PLP)
- Finally, we have seen how sometimes **ideas** rather than technology can make the difference

References I

-  Abiteboul, Serge, Richard Hull, and Victor Vianu (1995). **Foundations of databases**. Vol. 8. Addison-Wesley Reading.
-  Bancilhon, Francois, David Maier, Yehoshua Sagiv, and Jeffrey D. Ullman (1985). “Magic sets and other strange ways to implement logic programs”. In: **Proceedings of the fifth ACM SIGACT-SIGMOD symposium on Principles of database systems**. ACM, pp. 1–15. (Visited on 02/25/2015).
-  Bazoobandi, Hamid R., Harald Beck, and Jacopo Urbani (2017). “Expressive Stream Reasoning with Laser”. In: **ISWC**, pp. 87–103.
-  Beck, Harald, Minh Dao-Tran, and Thomas Eiter (2018). “LARS: A Logic-based framework for Analytic Reasoning over Streams”. In: **Artificial Intelligence 261**, pp. 16–70. ISSN: 0004-3702.
-  Cappuzzo, Riccardo, Paolo Papotti, and Saravanan Thirumuruganathan (2020). “Creating Embeddings of Heterogeneous Relational Datasets for Data Integration Tasks”. In: **SIGMOD**, pp. 1335–1349.

References II

-  Das, Sanjib, Paul Suganthan G.C., AnHai Doan, Jeffrey F. Naughton, Ganesh Krishnan, Rohit Deep, Esteban Arcaute, Vijay Raghavendra, and Youngchoon Park (2017). “Falcon: Scaling Up Hands-Off Crowdsourced Entity Matching to Build Cloud Services”. In: **SIGMOD**, pp. 1431–1446.
-  Konda, Pradap, Sanjib Das, Paul Suganthan G. C., AnHai Doan, Adel Ardalan, Jeffrey R. Ballard, Han Li, Fatemah Panahi, Haojun Zhang, Jeff Naughton, Shishir Prasad, Ganesh Krishnan, Rohit Deep, and Vijay Raghavendra (2016). “Magellan: toward building entity matching management systems”. In: **PVLDB** 9.12, pp. 1197–1208.
-  Kruit, Benno, Hongu He, and Jacopo Urbani (2020). “Tab2Know: Building a Knowledge Base from Tables in Scientific Papers”. In: **To appear at ISWC 2020**, pp. xxx–xxx.
-  Margara, Alessandro, Jacopo Urbani, Frank Van Harmelen, and Henri Bal (2014). “Streaming the web: Reasoning over dynamic data”. In: **Web Semantics: Science, Services and Agents on the World Wide Web** 25, pp. 24–44. (Visited on 04/30/2017).

References III

-  Mudgal, Sidharth, 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**, pp. 19–34.
-  Papadakis, George, Ekaterini Ioannou, and Themis Palpanas (2020). “Entity Resolution: Past, Present and Yet-to-Come.”. In: **EDBT**, pp. 647–650.
-  Raedt, Luc De, Angelika Kimmig, and Hannu Toivonen (2007). “ProbLog: A Probabilistic Prolog and Its Application in Link Discovery”. In: **IJCAI**, pp. 2462–2467.
-  Tsamoura, Efthymia, Víctor Gutiérrez-Basulto, and Angelika Kimmig (2020). “Beyond the Grounding Bottleneck: Datalog Techniques for Inference in Probabilistic Logic Programs”. In: **AAAI**, pp. 10284–10291.
-  Vlasselaer, Jonas, Guy Van den Broeck, Angelika Kimmig, Wannes Meert, and Luc De Raedt (2016). “TP-Compilation for inference in probabilistic logic programs”. In: **International Journal of Approximate Reasoning** 78, pp. 15–32.

References IV



Yu, Wenhao, Wei Peng, Yu Shu, Qingkai Zeng, and Meng Jiang (2020).
“Experimental Evidence Extraction System in Data Science with Hybrid Table
Features and Ensemble Learning”. In: **WWW**, pp. 951–961.