

Generating Contrastive Explanations for Inductive Logic Programming

Based on a Near Miss Approach

Johannes Rabold, Michael Siebers, Ute Schmid

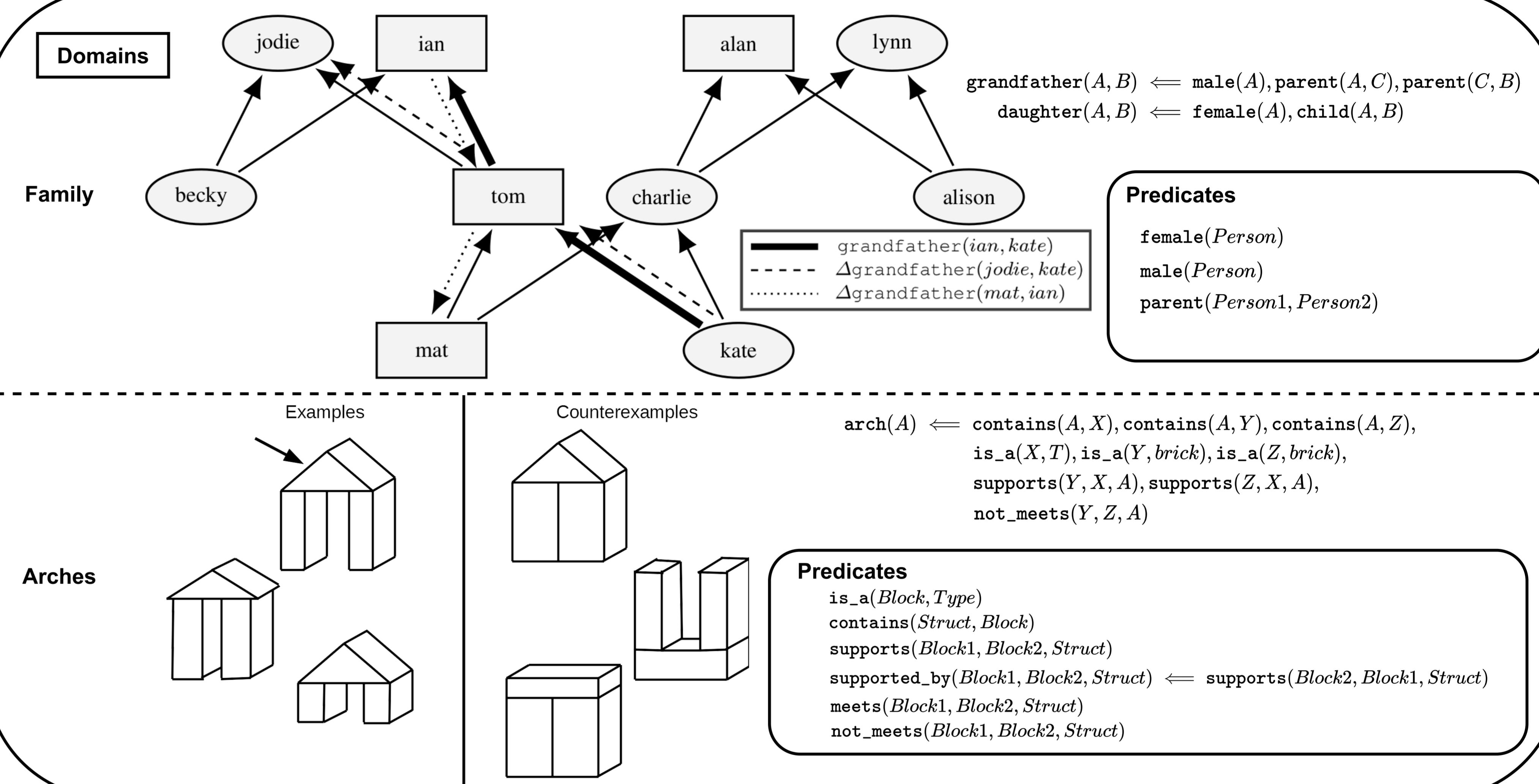
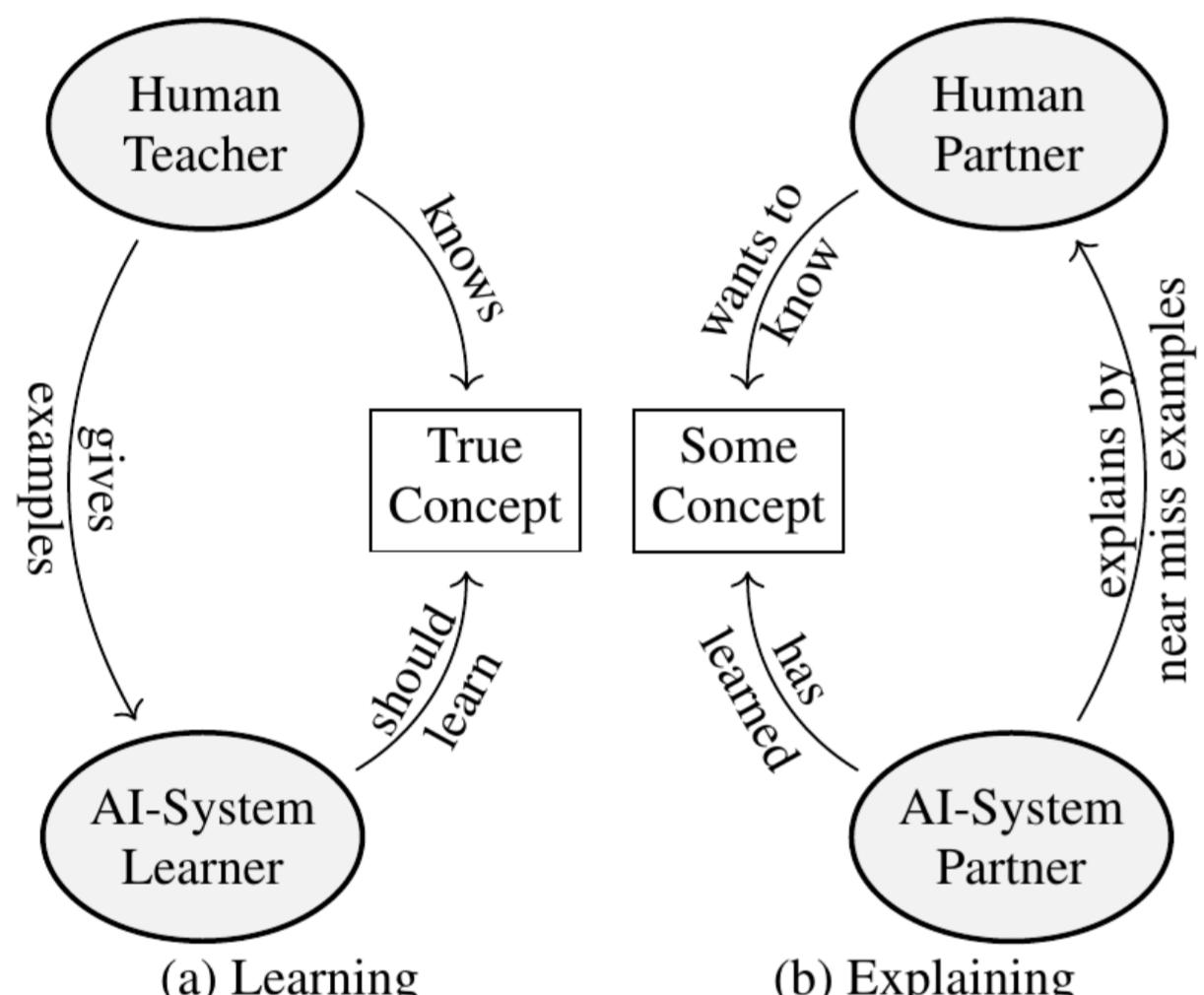
International Joint Conference on Learning and Reasoning

Introduction

Counterfactuals are explanations in the form of negative examples for a concept that only differ minimally from a given positive example.

You were denied a loan because your annual income was £30,000. If your income had been £45,000, you would have been offered a loan. (Wachter et al. 2017)

Such contrastive explanations were successfully applied in the image domain (e.g. CEM, Dhurandhar et al. 2018; MMD-critic, Kim et al. 2016) but to our knowledge not yet in relational domains. Winston showed in 1970 that presenting near miss examples in a relational domain results in faster learning. We propose what is effective for learning is also effective for explaining a learned model. In this work, we present GeNME, an algorithm that generates helpful near miss explanations for relational models.



GeNME: Generation of Near Miss Explanations

Require: Theory T , Finite set of rewriting filters O , Positive Example P

```

1: Initialize family of result sets  $(\mathcal{E}_i)_{i \in \mathbb{N}}$  to empty sets
2: Initialize the set of all near miss candidates  $\mathcal{N}$ 
3: for all local explanations  $C\theta$  for  $P$  where  $C \in T$  do
4:   for all  $\forall p \rightarrow q \in O$  do
5:     for all  $\mathcal{L} \in \forall p \rightarrow q(C)$  do
6:       create  $C'$  from  $C$  by replacing  $p$  with  $q$  in every literal in the body which is in  $\mathcal{L}$ 
7:       for all  $N \in \mathcal{N}$  do
8:          $d \leftarrow 0$ 
9:          $\mathcal{E} \leftarrow \{\}$ 
10:        while  $\mathcal{E} = \{\}$  and  $d < |\theta|$  do
11:           $d \leftarrow d + 1$ 
12:          for all partitions of  $\theta$  into  $\theta_1$  and  $\theta_2$  such that  $|\theta_2| = d$  do
13:            for all  $\theta'_2 = \{x_i \mapsto t'_i \mid x_i \mapsto t_i \in \theta_2\}$  such that no  $t'_i = t_i$  do
14:               $E \leftarrow C'(\theta_1 \cup \theta'_2)$ 
15:               $\mathcal{E} \leftarrow \mathcal{E} \cup \{E \mid T \models \text{body}(E) \text{ and head}(E) = N\}$ 
16:            end for
17:          end for
18:        end while
19:         $\mathcal{E}_d \leftarrow \mathcal{E}_d \cup \mathcal{E}$ 
20:      end for
21:    end for
22:  end for
23: end for
24: return  $(\mathcal{E}_i)_{i \in \mathbb{N}}$ 

```

	gf(ian,kate) $ \mathcal{N} = 96$	dt(becky,jodie) $ \mathcal{N} = 92$
male ↔ female		
$ \mathcal{E}_1 $	1	1
$ \mathcal{E}_2 $	2	3
$ \mathcal{E}_3 $	1	0
parent ↔ child		
$ \mathcal{E}_1 $	0	0
$ \mathcal{E}_2 $	2	6
$ \mathcal{E}_3 $	2	0
arch(struct1) $ \mathcal{N} = 3$		
meets ↔ not_meets		
$ \mathcal{E}_1 $	1	
$ \mathcal{E}_2 $	0	
$ \mathcal{E}_3 $	1	
supports ↔ supported_by		
$ \mathcal{E}_1 $	0	
$ \mathcal{E}_2 $	0	
$ \mathcal{E}_3 $	1	

Empirical Study

Types of explanations

General rule (R): global explanation of the concept a specific instance belongs to
Example (E): example-based explanation in form of a specific instance belonging to the concept
Near Miss (N): contrastive (negative) example with high degree of structural similarity to specific instance
Far Miss (F): contrastive (negative) example with low degree of structural similarity to specific instance

Information needs

Understanding general concept
Understanding particular example instance for concept
Understanding what is not in the concept (exclusion)

	(a) Family			
	(R)ule	(E)xample	(N)eir Miss	(F)ar Miss
general	4.97	4.52	2.93	2.19
example	4.14	4.70	2.49	2.37
exclusion	2.95	2.62	4.30	3.67

	(b) Arches			
	(R)ule	(E)xample	(N)eir Miss	(F)ar Miss
general	4.45	4.70	2.70	2.36
example	4.56	4.27	2.74	2.38
exclusion	3.25	2.73	3.95	3.82