



# Learning from Interpretation Transitions Using Differentiable Logic Programming Semantics

Kun Gao<sup>1</sup>, Hanpin Wang<sup>1, 2</sup>, Yongzhi Cao<sup>1</sup>, Katsumi Inoue<sup>3</sup>

<sup>1</sup> Peking University

<sup>2</sup> Guangzhou University

<sup>3</sup> National Institute of Informatics

2021/10/25



## Background knowledge

- ❖ A normal logic program is a finite set of rules that satisfies the following form:

$$h \leftarrow l_1 \wedge l_2 \wedge \cdots \wedge l_n,$$

where the  $l_i$ 's are **literals** and  $h$  is the head atom. Given a logic program  $P$ , the set of all atoms (ground atoms) in  $P$  is called a **Herbrand base** and is denoted as  $BP$ .

- ❖ An interpretation of a logic program is a subset of  $BP$ . Given a normal logic program  $P$  and  $BP = \{p_1, \dots, p_n\}$ , an interpretation vector is denoted as  $\mathbf{v} = (v_1, \dots, v_n)^T \in \{0, 1\}^n$ .



- ❖ Given a propositional logic program  $P$ , the immediate consequence operator  $TP: 2^{B_P} \rightarrow 2^{B_P}$  is defined as follows:

$$TP(I) = \{head(r) \mid r \in P, body^+(r) \subseteq I, body^-(r) \cap I = \emptyset\}$$

## The definition of the learning task

- ❖ Input: the pairs of interpretation transitions  $(I, J)$ , where  $J = TP(I)$ .
- ❖ Output: the logic program  $P$ .





## Normal matrices

Given a normal logic program, we use **normal (NOR) matrix** to represent it. The normal matrix of a logic program  $\mathbf{M}_P^{NOR}$  and the pairs of interpretation vectors meet:  $\mathbf{v}_o = \theta(\mathbf{M}_P^{NOR} \mathbf{v}_i)$ , where the threshold is defined as follows:

$$\theta(x) = \begin{cases} 1 & \text{if } x \geq 1; \\ 0 & \text{otherwise.} \end{cases}$$

The denominator of the fraction, corresponding the literal  $p$ , is the number of literals in the conjunctive clause with the least literals that appear in the corresponding rule

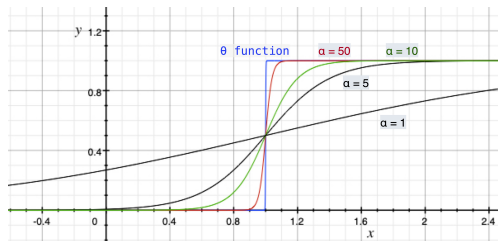
$$\begin{array}{l}
 p \leftarrow (p \wedge q \wedge \neg r) \vee (p \wedge \neg q) \\
 P : q \leftarrow p \wedge r \\
 r \leftarrow
 \end{array}
 \Rightarrow \mathbf{M}_P^{NOR} = \begin{array}{c}
 (p) \\
 (q) \\
 (r)
 \end{array} \begin{array}{cccccc}
 p & q & r & \neg p & \neg q & \neg r \\
 \frac{1}{2} & \frac{1}{3} & 0 & 0 & \frac{1}{2} & \frac{1}{3} \\
 \frac{1}{2} & 0 & \frac{1}{2} & 0 & 0 & 0 \\
 0 & 0 & 1 & 0 & 0 & 0
 \end{array}$$

**NOR matrix**



We use a differentiable function  $\phi$  to replace the function  $\theta$ :

$$\phi = \frac{1}{1 + e^{-ax}}$$



**Figure:** The curves of the function  $\theta$  and  $\phi(x - 1)$



- ❖ We define a differentiable formula to represent the immediate consequence function  $TP$  and the loss function:

$$\bar{\mathbf{v}}_o = \phi(\bar{\mathbf{M}}_P^{NOR} \mathbf{v}_i - \mathbf{1}),$$

$$\text{loss} = \lambda_1 \cdot H(\mathbf{v}, \bar{\mathbf{v}}) + \lambda_2 \cdot \sum_{c=1}^m \sum_{b=1}^{2n} \bar{\mathbf{M}}_P^{NOR}[c, b],$$

penalty term

$H$  is the cross-entropy function



## Meta-info learner

- ❖ Extract the number of non zero elements in every line, denoted as  $l_{p_k}$ , from the generated normal matrices.

$$\mathbf{M}_p^{NOR} = \begin{matrix} & p & q & r & \neg p & \neg q & \neg r \\ (p) & \frac{1}{2} & \frac{1}{3} & 0 & 0 & \frac{1}{2} & \frac{1}{3} \\ (q) & \frac{1}{2} & 0 & \frac{1}{2} & 0 & 0 & 0 \\ (r) & 0 & 0 & 1 & 0 & 0 & 0 \end{matrix} \begin{matrix} \Rightarrow \\ \Rightarrow \\ \Rightarrow \end{matrix}$$

NOR matrix

The meta-information:  
 $l_p = 4, l_q = 2, l_r = 1.$

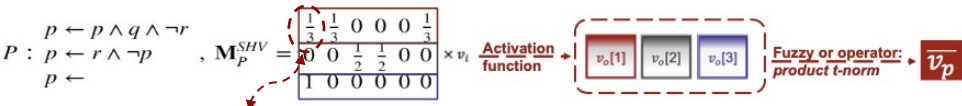
## Interpretation learner

- ❖ The number of rules  $m$  in a same head variable logic program is no larger than  $C\left(l_{p_k}, \left\lfloor \frac{l_{p_k}}{2} \right\rfloor\right)$ .
- ❖ For all  $k \in [1, C\left(l_{p_k}, \left\lfloor \frac{l_{p_k}}{2} \right\rfloor\right)]$ , learn the **same head variable (SHV)** matrix representing the  $k$ -th rule in the disjunctive normal form logic program.





- ❖ We set multiple trainable matrices  $M_P^{SHV}$  to represent the different disjunctive normal form rules in the logic program  $P$ .



The denominator of the fraction, corresponding the literal  $p$ , is the number of literals in the corresponding rule

$$\bar{v}_o[k] = 1 - \prod_{c=1}^{C(L_{p_k}, \lfloor \frac{L_{p_k}}{2} \rfloor)} (1 - \phi(\sum_{b=1}^{2n} \bar{M}_{P_k}^{SHV}[c, b] \cdot v_i[b] - 1)), \quad 1 \leq k \leq m.$$

product t-norm



The loss function is:

$$\text{loss} = \sum_{k=1}^m (\lambda_1 \cdot H(\mathbf{v}[k], \bar{\mathbf{v}}[k]) + \lambda_2 \cdot \sum_{c=1}^{C(L_{P_k}, \lfloor \frac{L_{P_k}}{2} \rfloor)} (1 - \sum_{b=1}^{2n} \bar{\mathbf{M}}_{P_k}^{SHV}[c, b])).$$

H is the cross-entropy function

Constrain the sum of each line of the SHV matrices to one.

## Example:

Logic program:  $p \leftarrow q \vee \neg p$

$q \leftarrow \neg p$

(a) A logic program  $P$

0	1	0	0
0	0	1	0

 $\times \mathbf{v}_i \xrightarrow{\text{or}} \bar{\mathbf{v}}_o[0]$

0	0	1	0
---	---	---	---

 $\times \mathbf{v}_i \xrightarrow{\text{or}} \bar{\mathbf{v}}_o[1]$

(b) Two SHV matrices to represent the logic program  $P$



## Transferring steps:

- ❖ Generate bottom clauses according to the relational database using bottom clause positionalization algorithm<sup>1</sup>;
- ❖ Regard the bottom clauses as propositional logic programs, and make the first-order features as propositional variables;
- ❖ Generate pairs of interpretation transitions.

---

<sup>1</sup> França, M. V. M., Zaverucha, G., & DAvila Garcez, A. S. (2014). Fast relational learning using bottom clause propositionalization with artificial neural networks. *Machine Learning*, 94(1), 81–104.



## Example

- ❖  $F = \{\text{mother}(\text{mom1}, \text{daughter1}), \text{wife}(\text{daughter1}, \text{husband1}), \text{wife}(\text{daughter2}, \text{husband2})\}$ ,
- ❖  $P = \{\text{motherInLaw}(\text{mom1}, \text{husband1})\}$
- ❖  $N = \{\text{motherInLaw}(\text{daughter1}, \text{husband2})\}$ .

Generated bottom clause set (the depth of the variable is set to 1):

Bottom  
clause

$$E_{\perp} = \{\text{motherInLaw}(A, B) : \neg \text{mother}(A, C), \text{wife}(C, B); \\ \sim \text{motherInLaw}(A, B) : \neg \text{wife}(A, C)\}.$$

Pairs of  
interpretation

$I_1: \text{mother}(A, C), \text{wife}(C, B)$

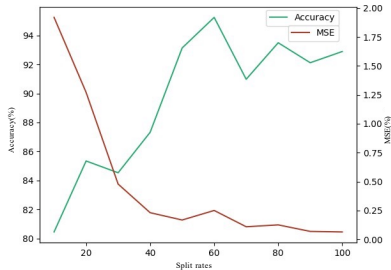
$J_1: \text{motherInLaw}(A, B)$

$I_2: \text{wife}(A, C)$

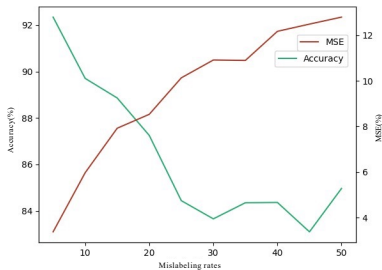
$J_2:$



Test the model on the incomplete and mislabelled datasets:



(a) Incomplete data



(b) Mislabeled data

**Figure:** Mean accuracy of the logic program and the MSE of the predicted Boolean value with respect to different split rates and mislabeling rates of the fission dataset.



**Table 1** Comparison of MSE (%) and accuracy (%) on partial datasets with different split rates

Datasets (Complexity)	Model	The split rates				
		10%	20%	50%	100%	
Fission (10, 2 <sup>16</sup> )	NN-LFIT	<b>1.2, 98.80</b>	<b>0.2, 99.80</b>	<b>0.00, 99.92</b>	<b>0.00, 99.87</b>	
	D-LFIT	1.91, 80.45	1.27, 85.33	0.17, 93.13	0.06, 92.89	
	LFIT	-, 76.85	-, 77.03	-, 77.15	-, <b>100</b>	
	JRip	-, 79.14	-, 78.05	-, 80.04	-, 78.47	
	Mammalian (10, 2 <sup>16</sup> )	NN-LFIT	<b>1.3, 96.6</b>	<b>0.4, 94.35</b>	<b>0.00, 99.89</b>	<b>0.00, 99.91</b>
Mammalian (10, 2 <sup>16</sup> )	D-LFIT	1.73, 71.67	1.21, 75.9	0.79, 80.09	0.52, 82.84	
	LFIT	-, 76.01	-, 76.48	-, 76.73	91.56	
	JRip	-, 77.84	-, 75.44	-, 76.41	-, 74.66	
	Budding (12, 2 <sup>12</sup> )	NN-LFIT	ROT	ROT	ROT	ROT
	D-LFIT	1.03, <b>71.96</b>	0.43, <b>71.39</b>	0.15, <b>70.5</b>	0.09, <b>76.52</b>	
Arabidopsis (15, 2 <sup>15</sup> )	LFIT	ROT	ROT	ROT	ROT	
	JRip	-, 67.97	-, 68.55	-, 67.91	-, 68.32	
	Mutagenesis (1111, 188)	NN-LFIT	ROT	ROT	ROT	ROT
	D-LFIT	0.57, <b>84.35</b>	0.51, <b>86.83</b>	0.48, <b>88.56</b>	0.45, <b>89.70</b>	
	LFIT	ROT	ROT	ROT	ROT	
UW-CSE (601, 1614)	JRip	-, 68.84	-, 69.00	-, 68.79	-, 68.67	
	D-LFIT	<b>77.78</b>	<b>94.44</b>	<b>88.88</b>	<b>83.33</b>	
	JRip	58.37	59.46	65.97	66.45	
	Alzheimers-amine (1084, 686)	D-LFIT	<b>75.00</b>	<b>78.12</b>	<b>78.44</b>	<b>79.44</b>
	JRip	70.18	70.81	73.85	74.35	
Alzheimers-amine (1084, 686)	D-LFIT	58.42	<b>60.29</b>	<b>63.24</b>	<b>67.65</b>	
	JRip	<b>58.74</b>	57.66	57.14	54.81	

**Table 3** Comparison of the MSE (%) and accuracy (%) on the mislabeled data with different mislabeling rates

Datasets	Model name	The rates of mislabeled data				
		5%	20%	35%	50%	
Fission	NN-LFIT	3.25, <b>96.75</b>	10.56, <b>89.43</b>	15.14, <b>84.86</b>	17.74, 82.25	
	D-LFIT	<b>2.23</b> , 92.34	<b>8.53</b> , 87.25	<b>10.89</b> , 84.34	<b>12.08</b> , <b>84.96</b>	
	LFIT	-, 77.25	-, 77.16	-, 77.38	-, 77.32	
	JRip	-, 78.91	-, 78.54	-, 78.55	-, 78.15	
	Mammalian	NN-LFIT	4.77, 79.72	16.78, 11	20.98, 79.01	23.20, 74.49
Mammalian	D-LFIT	<b>3.45</b> , <b>80.00</b>	<b>11.53</b> , <b>78.82</b>	<b>15.74</b> , <b>80.29</b>	<b>16.35</b> , <b>86.27</b>	
	LFIT	-, 76.72	-, 76.73	-, 76.62	-, 77.32	
	JRip	-, 74.21	-, 74.45	-, 74.43	-, 74.00	
	Budding	NN-LFIT	ROT	ROT	ROT	ROT
	D-LFIT	4.9, <b>76.42</b>	13.3, <b>74.28</b>	16.53, <b>73.41</b>	18.21, <b>74.71</b>	
Arabidopsis	LFIT	ROT	ROT	ROT	ROT	
	JRip	-, 67.99	-, 67.15	-, 66.80	-, 66.41	
	Mutagenesis	NN-LFIT	ROT	ROT	ROT	ROT
	D-LFIT	4.8, <b>81.46</b>	11.83, <b>76.59</b>	15.4, <b>70.28</b>	17.35, 64.90	
	LFIT	ROT	ROT	ROT	ROT	
UW-CSE	JRip	-, 68.27	-, 67.34	-, 66.94	-, <b>66.65</b>	
	D-LFIT	<b>88.89</b>	<b>83.33</b>	<b>88.89</b>	<b>88.89</b>	
	JRip	66.49	62.23	62.77	62.23	
	UW-CSE	D-LFIT	<b>73.49</b>	<b>72.67</b>	<b>73.29</b>	<b>73.29</b>
	JRip	72.80	67.35	66.04	66.23	
Alzheimers-amine	D-LFIT	<b>63.24</b>	<b>63.24</b>	<b>60.29</b>	<b>58.83</b>	
	JRip	48.35	49.42	49.27	50.87	

## Comparisons on the incomplete data

**Table 2** Comparison of the accuracy (%) of rules obtained by CILP++

Model	Datasets		
	Mutagenesis	UW_CSE	Alzheimer-amine
CILP++	77.72	<b>81.98</b>	<b>78.70</b>
D-LFIT	<b>83.33</b>	79.44	67.75

## Comparisons on the mislabeled data

## Comparisons on the relational databases



- ❖ D-LFIT is an inductive logic programming learner, which can learn propositional logic programs from **mislabeled** data or **incomplete** data.
- ❖ Through adopting the BCP algorithm, we can learn **first-order logic programs** from relational databases.
- ❖ D-LFIT is a robust, fast learner, which can **curriculum learning strategy** to learn knowledge from data.
- ❖ We will devise more constraints to make the generated logic programs meet destined formats.
- ❖ We will apply the D-LFIT on other dataset such like knowledge graph.



**Thanks for your attention!**

