

Learning First-Order Rules with Differentiable Logic Program Semantics

Kun Gao¹, Katsumi Inoue², Yongzhi Cao¹, Hanpin Wang^{3, 1}

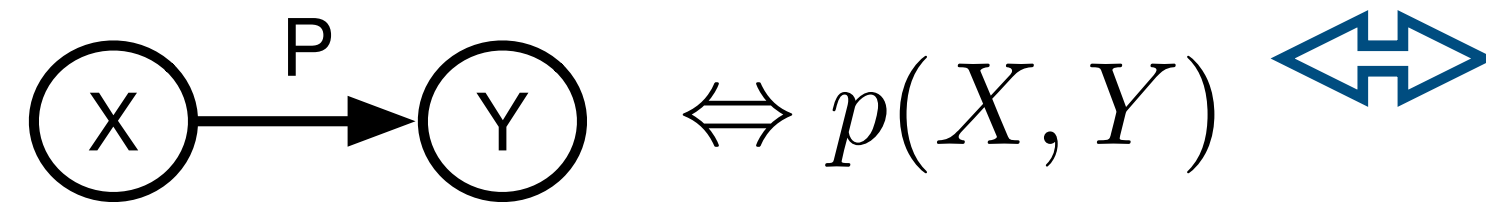
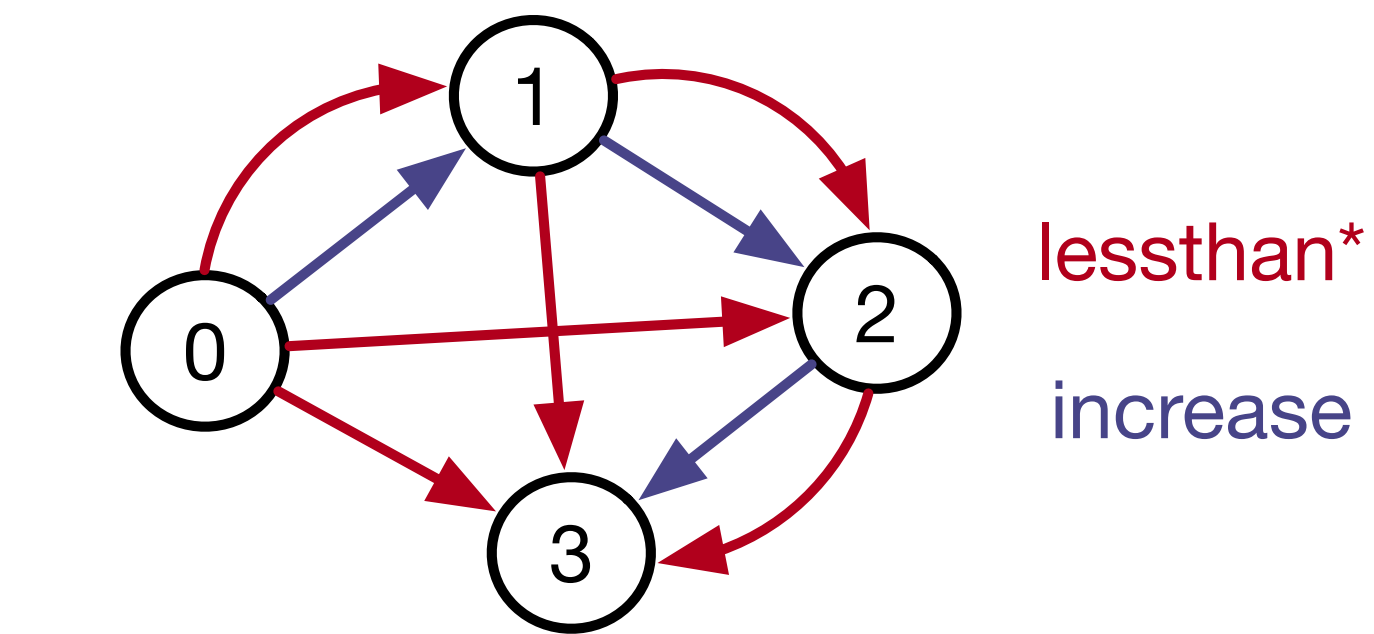
¹ Peking University

² National Institute of Informatics

³ Guangzhou University



Motivation



lessthan*
increase

Background Knowledge

increase(0,1) increase(1,2) increase(2,3)

B

Positive Example

lessthan(0,1) lessthan(1,2) lessthan(2,3)
lessthan(0,3) lessthan(0,2)...

P

ILP



```
lt(X,Y) :- inc(X,Y).
lt(X,Y) :- inc(X,Z) &
            inc(Z,Y).
```

Solution

Negative Example

lessthan(1,0) lessthan(2,0) lessthan(3,1)...

N

0, 1, 2, 3 ↔ constant

Lessthan, Increase ↔ predicate

X, Y ↔ variable

Preliminaries

Differential logic program semantics

- **Matrix representations for logic program**

- P with the head atom $p(X, Y)$, 5 possible body atoms, and 3 different rules, and one of matrix embedding $M_P \in [0,1]^{3 \times 5}$ is:

$$\begin{aligned}
 p(X, Y) &\leftarrow p(Y, Z) \wedge p(Z, X). \\
 p(X, Y) &\leftarrow p(X, Z) \wedge p(Z, Y). \\
 p(X, Y) &\leftarrow p(Y, X).
 \end{aligned}$$

$$\mathbf{M}_P = \begin{matrix} & \begin{matrix} p(X, Z) & p(Y, X) & p(Y, Z) & p(Z, X) & p(Z, Y) \end{matrix} \\ \begin{bmatrix} 0 & 0 & \frac{1}{3} & \frac{2}{3} & 0 \\ \frac{1}{2} & 0 & 0 & 0 & \frac{1}{2} \\ 0 & 1 & 0 & 0 & 0 \end{bmatrix} & \begin{matrix} \Rightarrow \\ \text{sum} \end{matrix} & 1 \end{matrix}$$

Preliminaries

Differential logic program semantics

- Vector interpretations:**

- An interpretation vector $v \in \{0,1\}^m$ represents an interpretation in a logic program. If the Boolean value of α_k is True, then $a_k = 1$; otherwise, $a_k = 0$

$p(X, Y) \leftarrow p(Y, Z) \wedge p(Z, X).$
 $p(X, Y) \leftarrow p(X, Z) \wedge p(Z, Y).$
 $p(X, Y) \leftarrow p(Y, X).$

Interpretation 1

$p(Y, Z), p(Z, X)$

Interpretation 2

$p(X, Z), p(Z, X)$

$F_{body} = \{p(X, Z), p(Y, X), p(Y, Z), p(Z, X), p(Z, Y)\}$

$\mathbf{v}_1 = [0, 0, 1, 1, 0]$

$\mathbf{v}_2 = [1, 0, 0, 0, 1]$

Preliminaries

Differential logic program semantics

Sakama et al., *ASPOCP*, 2018; Gao et al., *Mach. Learn.*, 2022;

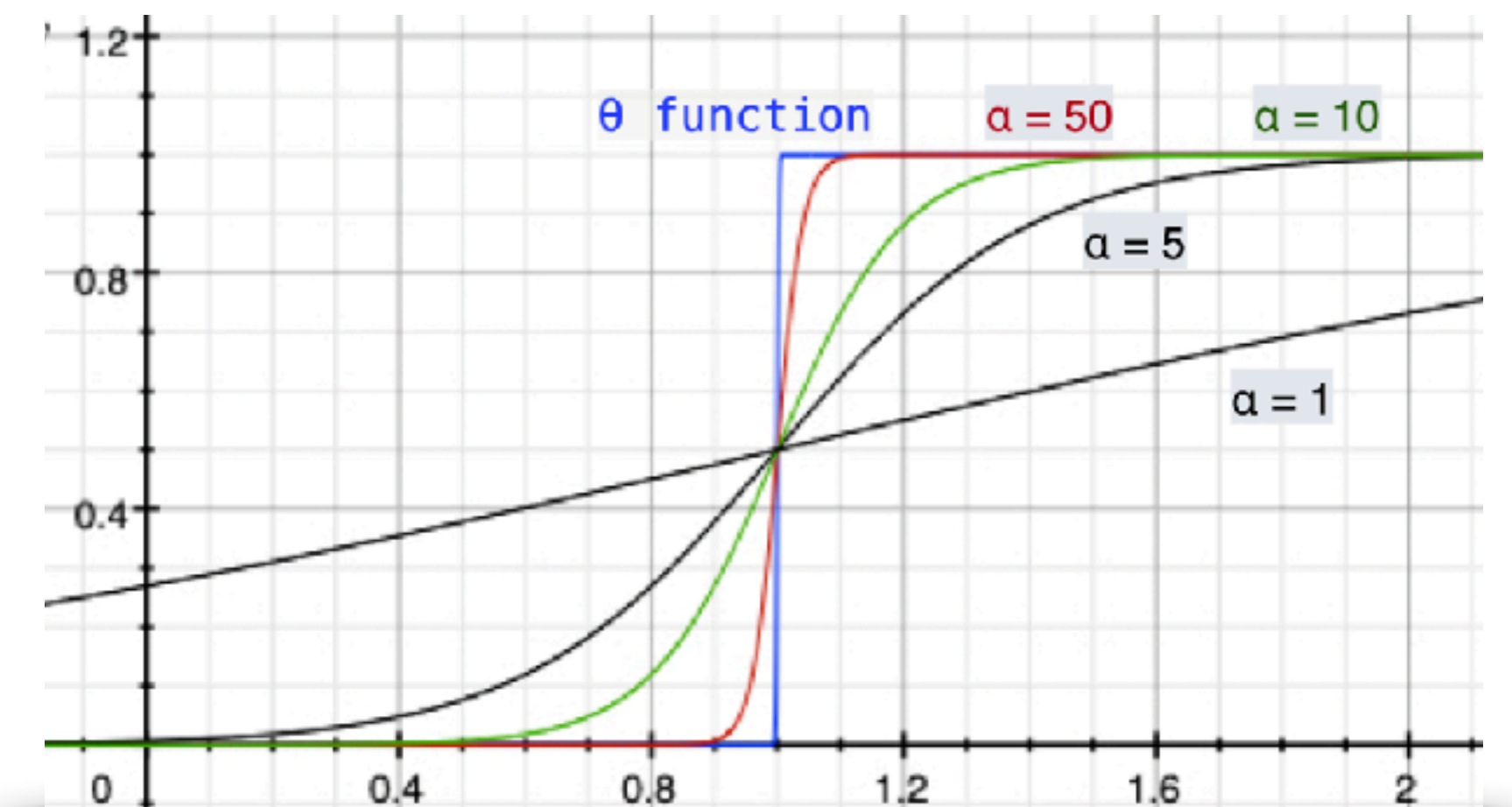
- Deduction reasoning:**

$$\mathbf{v}_o = \bigvee_{k=1}^m \theta(\mathbf{M}_P[k, \cdot] \times \mathbf{v}_i^T) \quad \begin{matrix} p(X, Z), p(Y, X), p(Y, Z), p(Z, X), p(Z, Y) \\ \begin{bmatrix} 0 & 0 & \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & 0 & 0 & 0 & \frac{1}{2} \\ 0 & 1 & 0 & 0 & 0 \end{bmatrix} \end{matrix} \times [0, 0, 1, 1, 0]^T (\mathbf{v}_i) = [1, 0, 0]^T (\mathbf{v}_o)$$

- Differentiable deduction:**

- Differentiable threshold function
- Differentiable fuzzy logic (product-t norm)

$$\theta \Rightarrow \phi = \frac{1}{1 + e^{-\alpha(x-1)}} \quad x \vee y \Rightarrow 1 - (1 - x)(1 - y)$$

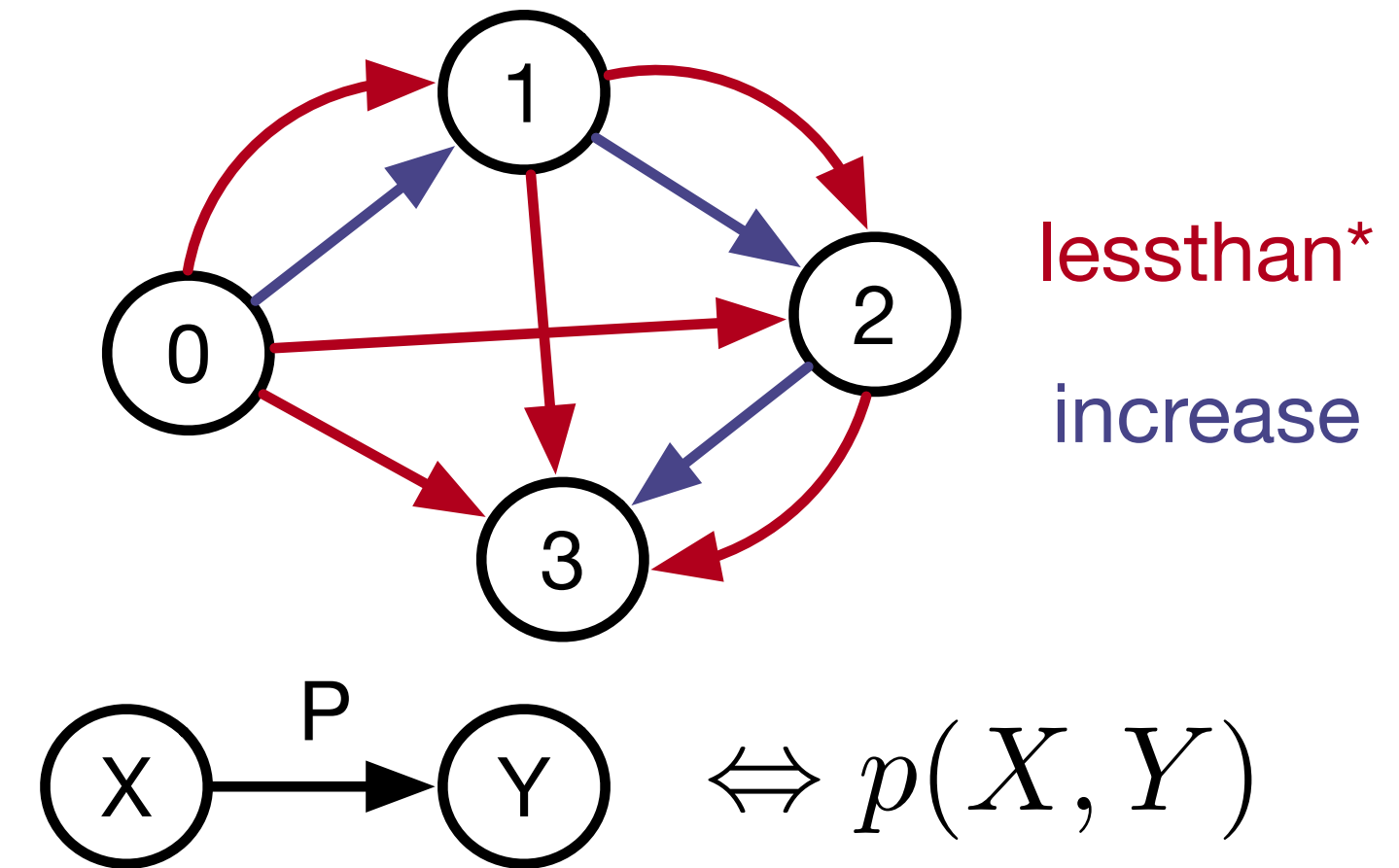


Methods

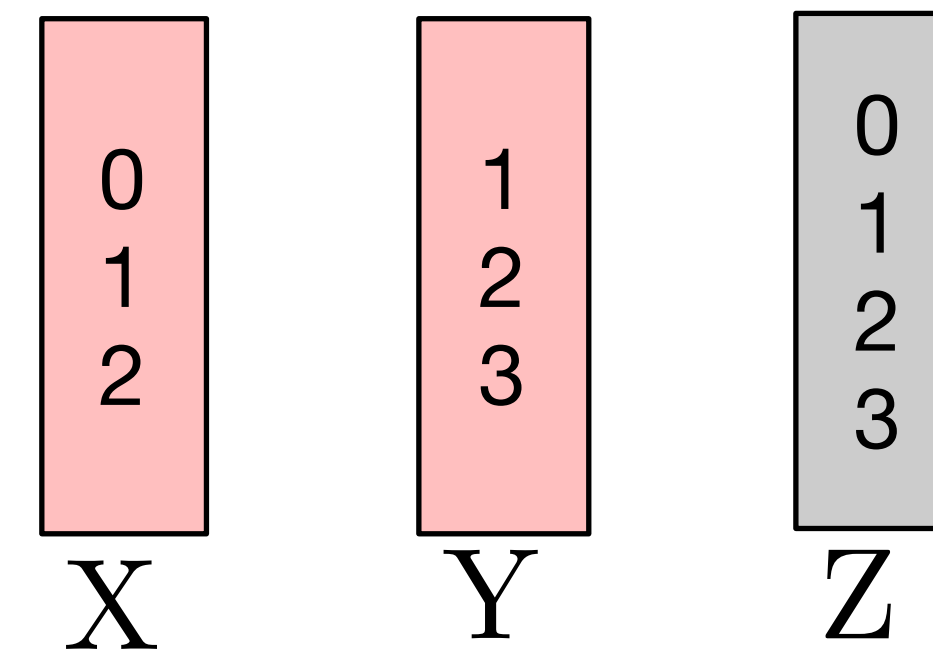
- **Propositionalization method**
 - Data preprocessing for relational facts
- **Differentiable inductive learning**
 - Learn parameters in matrix embeddings through neural networks
- **Logic rules extraction**
 - Interpret symbolic logic programs from matrix embeddings

Propositionalization

- Transfer the relational data into neural network readable attribute-valued data;
- **Input:** The training relational facts;
- **Output:** The trainable pairs of interpretation vectors;



$$d = 1, V = \{X, Y, Z\}$$



1. **Preparation:** Allocate constants to variables

Propositionalization

$$X \times Y \times Z = S$$

$$S = \{(0, 1, 0), (0, 1, 1), \dots, (2, 3, 3)\}$$

$$P_F = \{lt(X, Z), lt(Y, X), lt(Y, Z), lt(Z, X), lt(Z, Y), inc(X, Y), inc(X, Z), inc(Y, X), inc(Y, Z), inc(Z, X), inc(Z, Y)\}$$

$$v_i = [0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 1]$$

$$v_o = [1]$$

⋮

2. **Generation:** Compute all substitutions and generate all input and output vector interpretation

$$v_i = [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0] \quad v_o = [-] \quad \times$$

$$v_i = [0, 0, 1, 0, 1, 1, 0, 0, 0, 0, 1]$$

$$v_i = [1, 0, 0, 0, 1, 0, 1, 0, 1, 0, 1]$$

⋮

$$v_i = [0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 1]$$

$\times \times$

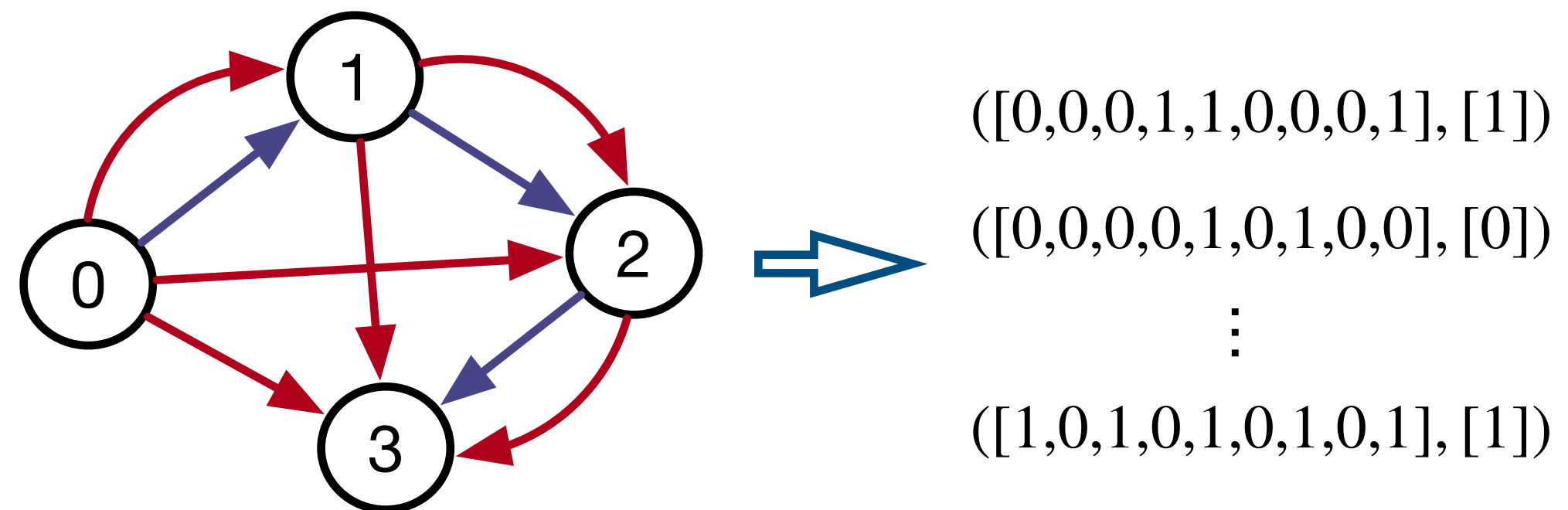
\times

$$P_F = \{lt(X, Z), \boxed{lt(Y, X)}, lt(Y, Z), \boxed{lt(Z, X)}, lt(Z, Y), inc(X, Y), inc(X, Z), inc(Y, X), inc(Y, Z), \boxed{inc(Z, X)}, inc(Z, Y)\}$$

3. **Examination:** Delete the noisy data and features

Propositionalization

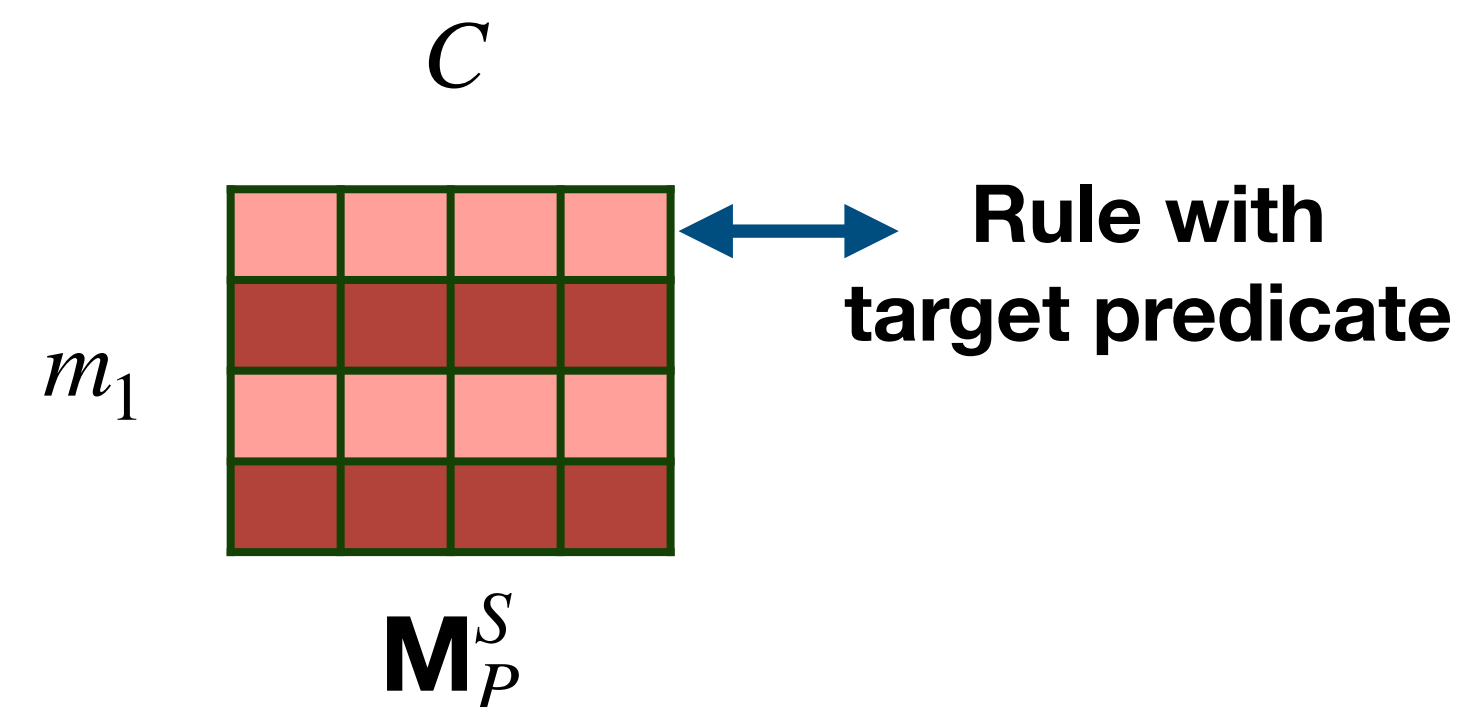
- number of **valid** first-order features: $|v_i| = C$;
- $I_o = T_P(I_i)$ iff $v_o = D_P(v_i)$, where D_P is the differentiable deduction reasoning of a logic program P .



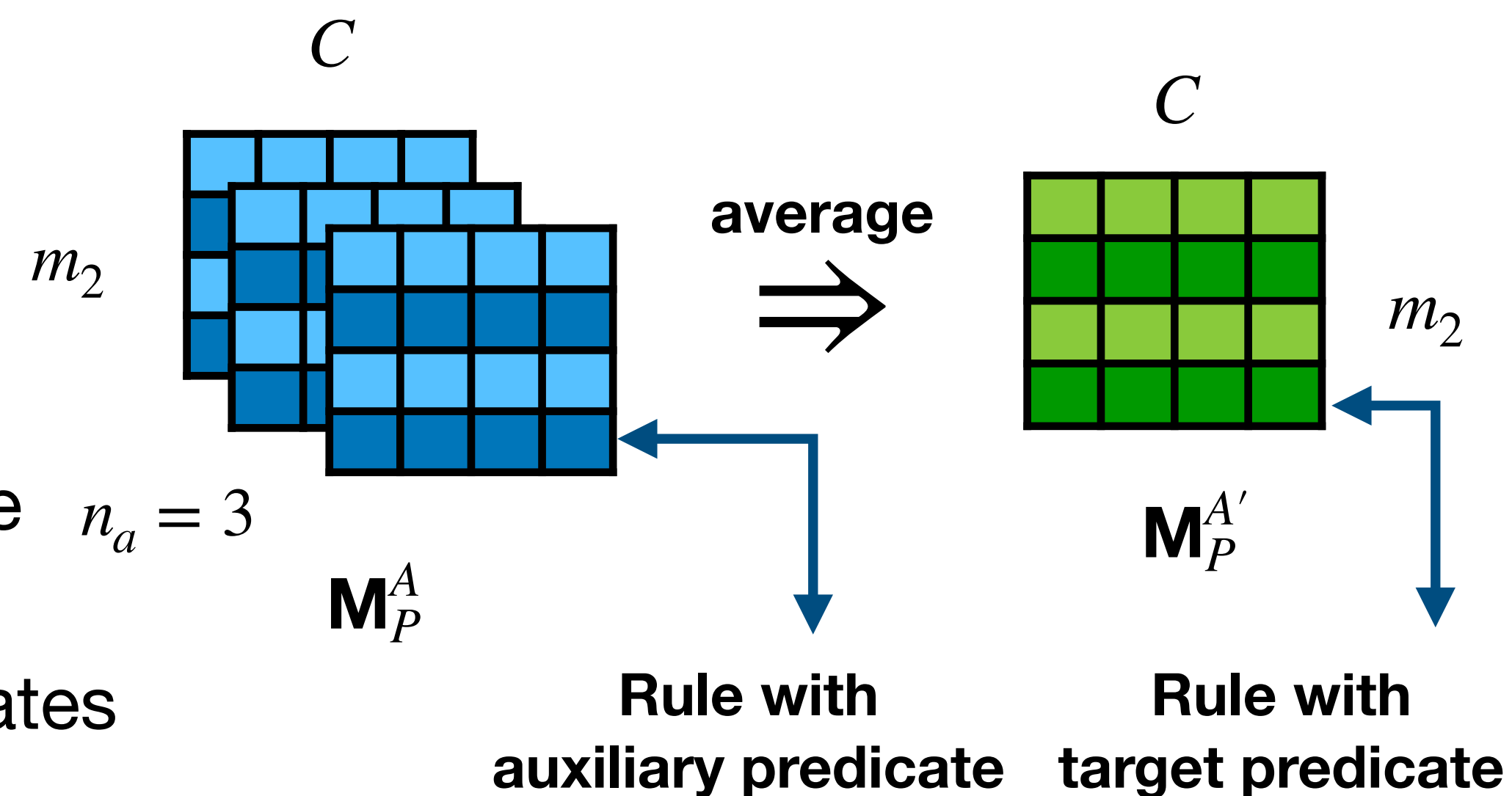
Differentiable Inductive Learning

Matrix Representation

- Trainable matrix $M_P^S \in [0,1]^{m_1 \times C}$
 - m_1 - the number of rules headed by target predicate
 - C - valid first-order features



- Trainable matrix $M_P^A \in [0,1]^{m_2 \times n_a \times C}$
 - m_2 - number of rules headed by target predicate
 - n_a - number of rules headed by auxiliary predicates



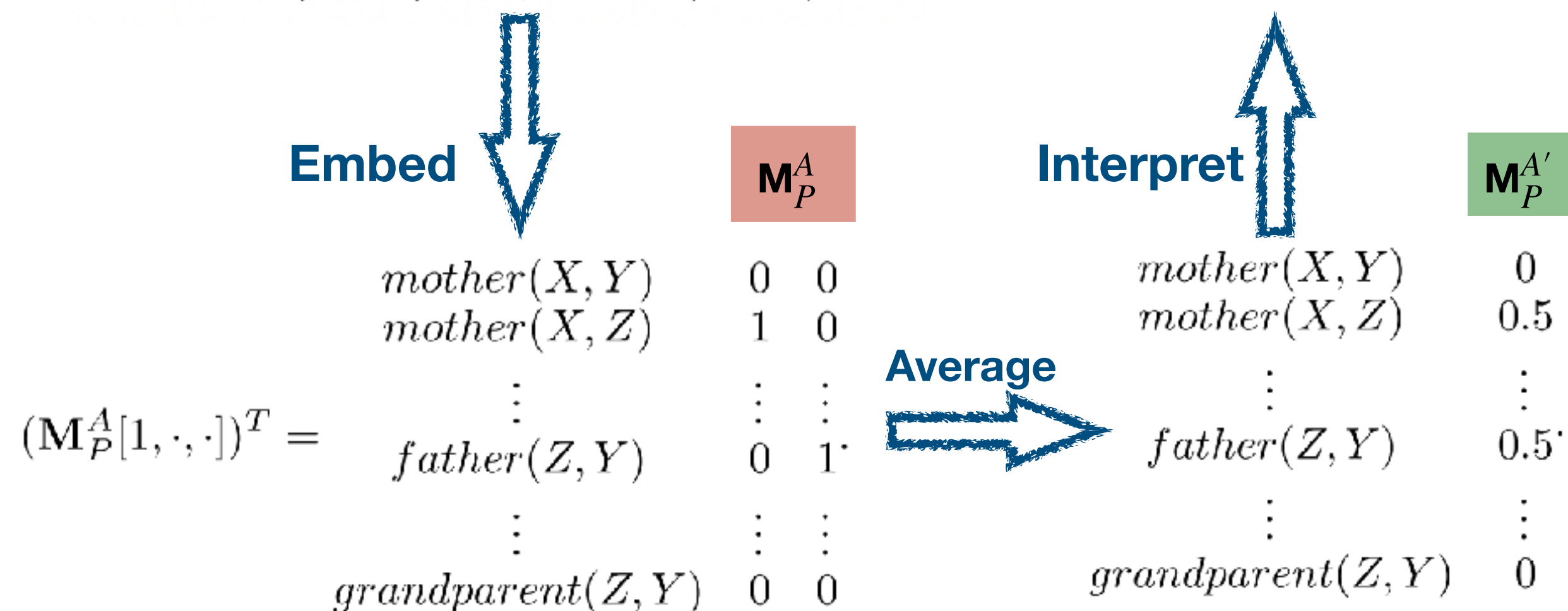
Differentiable Inductive Learning

Matrix Representation

- The reason for applying average operation:
 - Keep the non-zero values index information from M_P^A to $M_P^{A'}$
 - Increase the trainable parameters

$parent(X, Z) \leftarrow mother(X, Z),$
 $parent(Z, Y) \leftarrow father(Z, Y).$

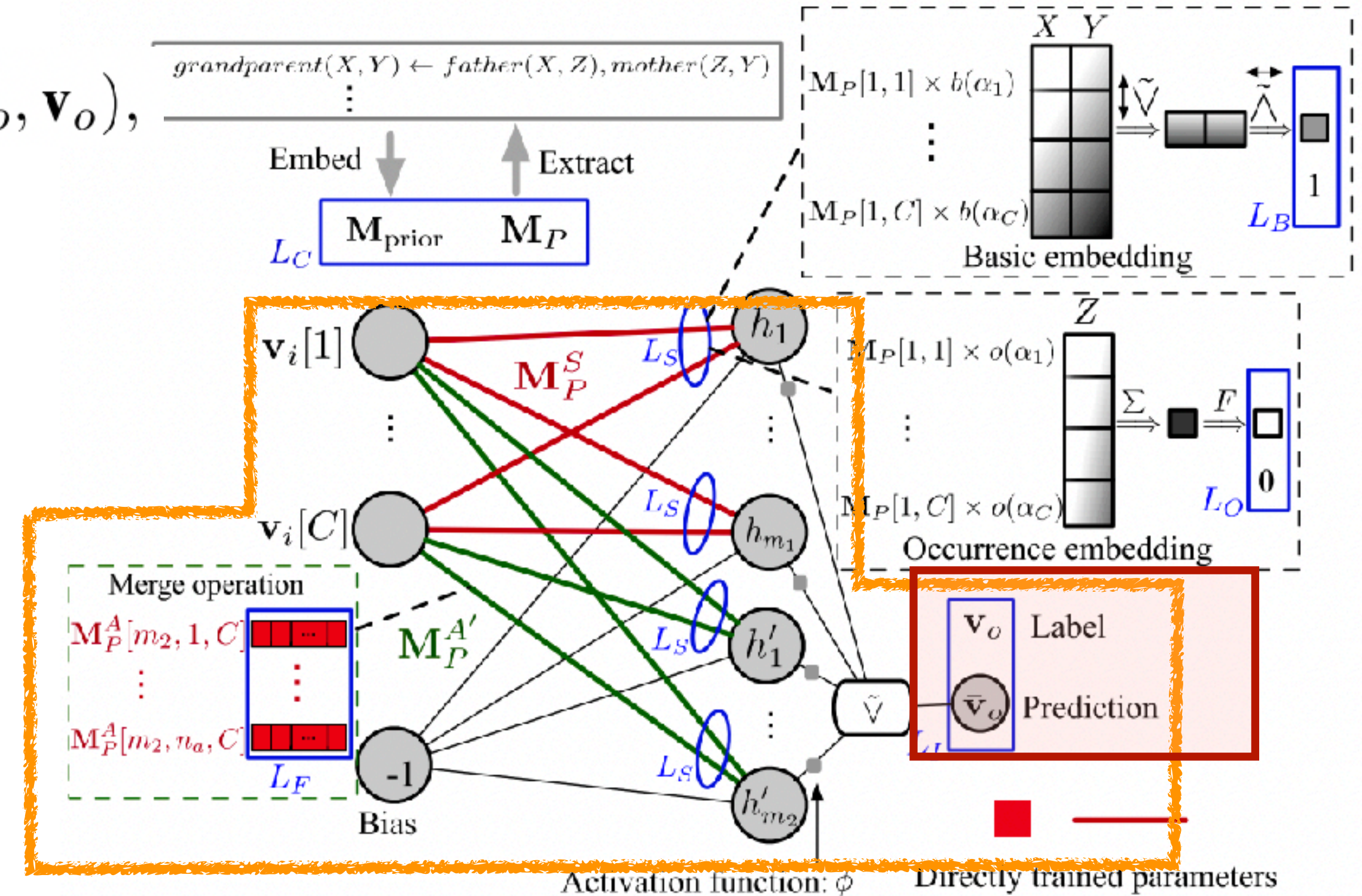
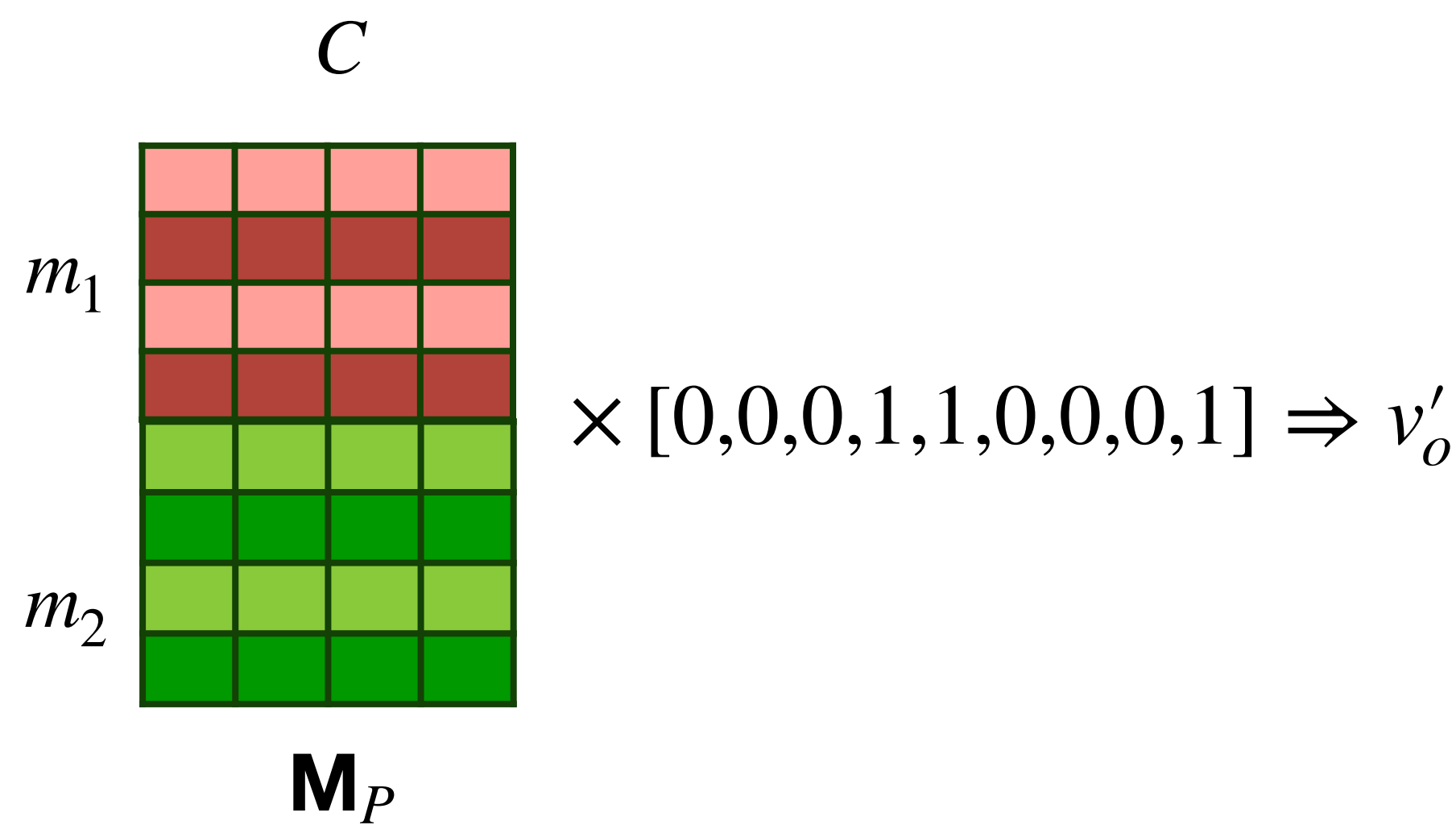
$grandparent(X, Y) \leftarrow mother(X, Z), father(Z, Y).$



Differentiable Inductive Learning Constraint Functions

- Inference loss:

$$\bar{\mathbf{v}}_o = \tilde{\mathbf{V}}_{k=1}^m (\phi(\mathbf{M}_P[k, \cdot] \times \mathbf{v}_i^T - 1)), L_I = H(\bar{\mathbf{v}}_o, \mathbf{v}_o),$$



Differentiable Inductive Learning

Forward-chained formats

- Generate rules with the forward-chained format:

$$p_t(X, Y) \leftarrow p_1(X, Z_1) \wedge p_2(Z_1, Z_2) \wedge \dots \wedge p_{n+1}(Z_n, Y)$$

Head variables

Body variables

- **Basic constraint:** The body of each rule in an LP contains all variables that appear in the head atom.
- **Occurrence constraint:** In the body of each rule, the number of occurrences of each variable that does not appear in the header atom is not one.


Differentiable Inductive Learning

Forward-chained formats

- Generate rules with the forward-chained format:

$$p_t(X, Y) \leftarrow p_1(X, Z_1) \wedge p_2(Z_1, Z_2) \wedge \dots \wedge p_{n+1}(Z_n, Y)$$

$n(Z_1) \geq 2$



- **Basic constraint:** The body of each rule in an LP contains all variables that appear in the head atom.
- **Occurrence constraint:** In the body of each rule, the number of occurrences of each variable that does not appear in the header atom is not one.

Differentiable Inductive Learning Constraint Functions

- Number of variables in head predicate a ;
Number of variable depth d .

- **Basic embedding:** Atom $\rightarrow \{0,1\}^a$

- **Occurrence embedding:** Atom $\rightarrow \{0,1\}^d$

$$p_t(X, Y) \Rightarrow [1,1]$$

$$p_t(X, Y) \Rightarrow [0]$$

$$p_t(X, Z) \Rightarrow [1,0]$$

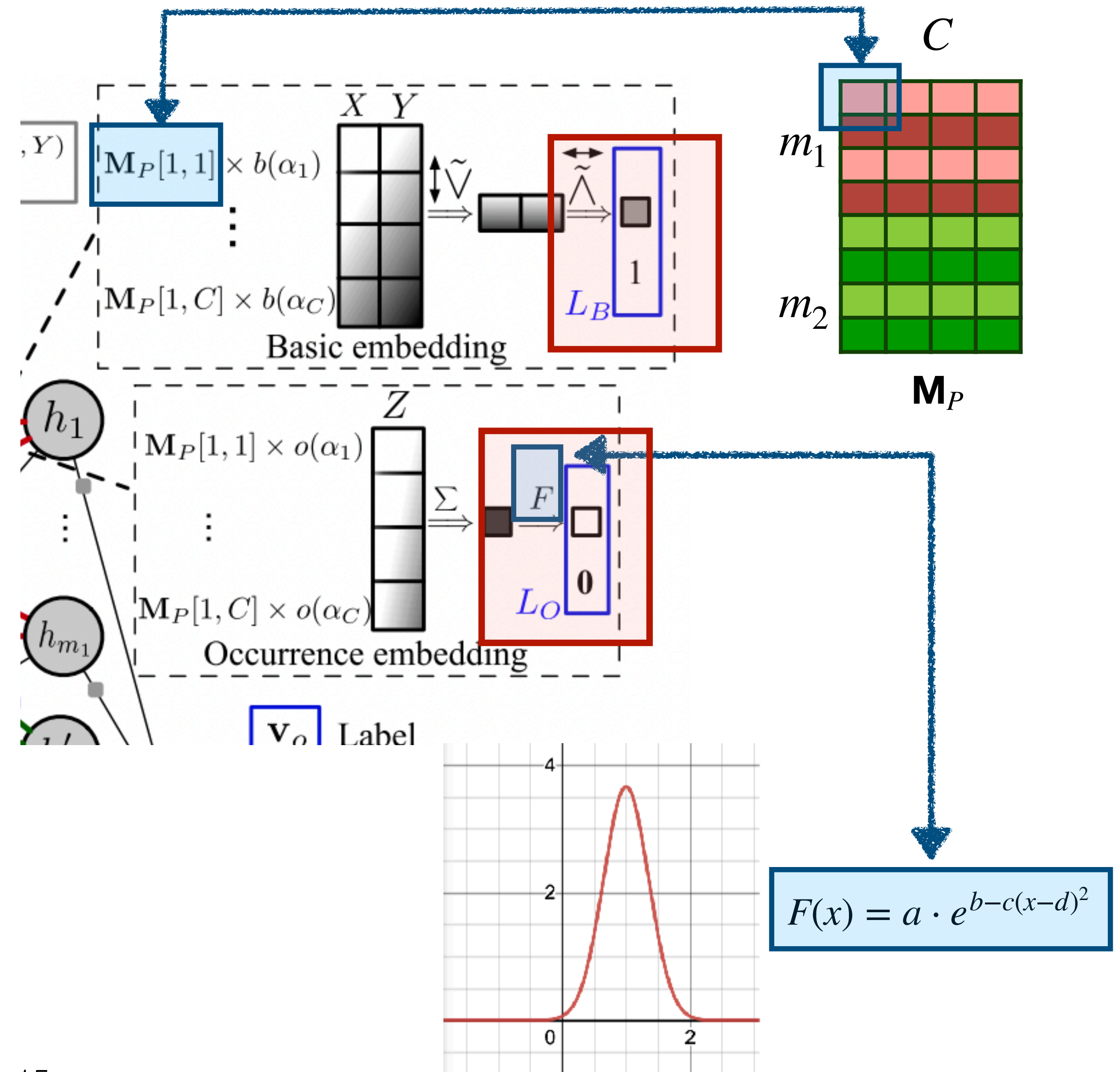
$$p_t(X, Z) \Rightarrow [1]$$

$$p_t(Z, Y) \Rightarrow [0,1]$$

$$p_t(Z, Y) \Rightarrow [1]$$

Basic embeddings Occurrence embeddings

- **Basic Loss, Occurrence Loss**



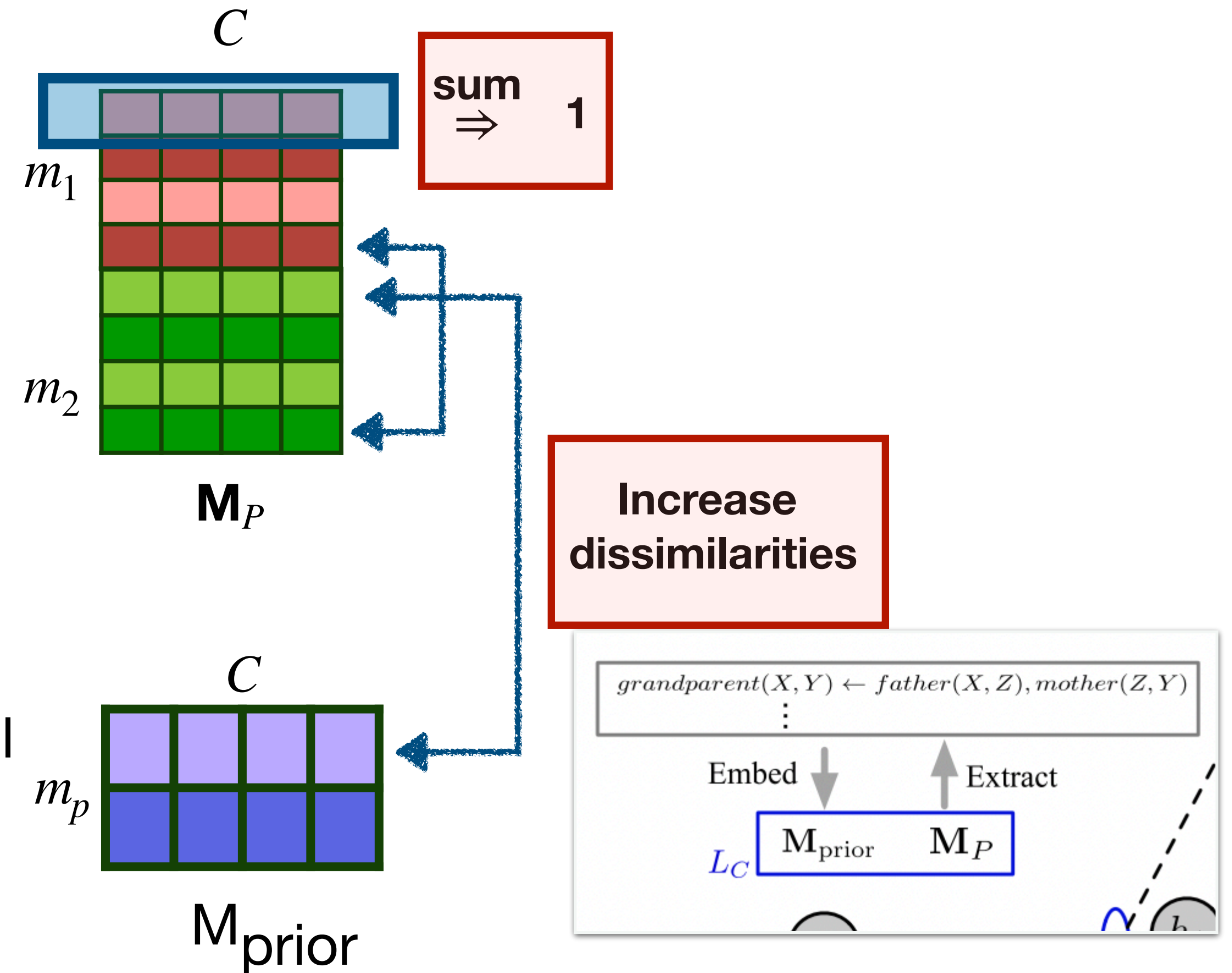
Differentiable Inductive Learning Constraint Functions

- Sum loss:**

$$L_S = \sum_{k=1}^m MSE\left(\sum_i^C M_P[k, i], 1\right)$$

- Similar loss:**

- Each row in M_P have more dissimilarities.
- The rows in trainable matrix M_P and trained-well correct matrix M_{prior} have more dissimilarities.



Logical Rules Extraction

- Use a series of thresholds called τ_f to extract rules from M_P , and use Datalog to check the **precision** $\frac{n_r}{n_b}$ of a rule:
 - n_r - the number of substitutions meet both **head** and **body** of a rule
 - n_b - the number of substitutions meet only **body** of a rule
- Use τ_s to select the rules which precision values are larger than it.
- Use **accuracy (recall)** to indicate the ratio of the covered positive test examples.

$$M_P = \begin{matrix} & \begin{matrix} p(X, Z) & p(Y, X) & p(Y, Z) & p(Z, X) & p(Z, Y) \end{matrix} \\ \begin{bmatrix} 0 & 0 & \frac{1}{3} & \frac{2}{3} & 0 \\ \frac{1}{2} & 0 & 0 & 0 & \frac{1}{2} \\ 0 & 1 & 0 & 0 & 0 \end{bmatrix} \end{matrix}$$

$$\Downarrow \tau_f = \{0, 0.2, \dots, 0.9\}$$

$$\begin{aligned} p(X, Y) &\leftarrow p(Y, Z) \wedge p(Z, X). && \mathbf{1.0} \\ p(X, Y) &\leftarrow p(X, Z) \wedge p(Z, Y). && \mathbf{1.0} \\ p(X, Y) &\leftarrow p(Y, X). && \mathbf{1.0} \\ p(X, Y) &\leftarrow p(Z, X) && \mathbf{0.2} \end{aligned}$$

$$\Downarrow \tau_s = 0.5$$

$$\begin{aligned} p(X, Y) &\leftarrow p(Y, Z) \wedge p(Z, X). && \mathbf{1.0} \\ p(X, Y) &\leftarrow p(X, Z) \wedge p(Z, Y). && \mathbf{1.0} \\ p(X, Y) &\leftarrow p(Y, X). && \mathbf{1.0} \end{aligned}$$

Experimental Results

- From ILP datasets, the generated rules have:

- $\tau_s = 1$, so the precision values are 1.
- Accuracy value is 100%.

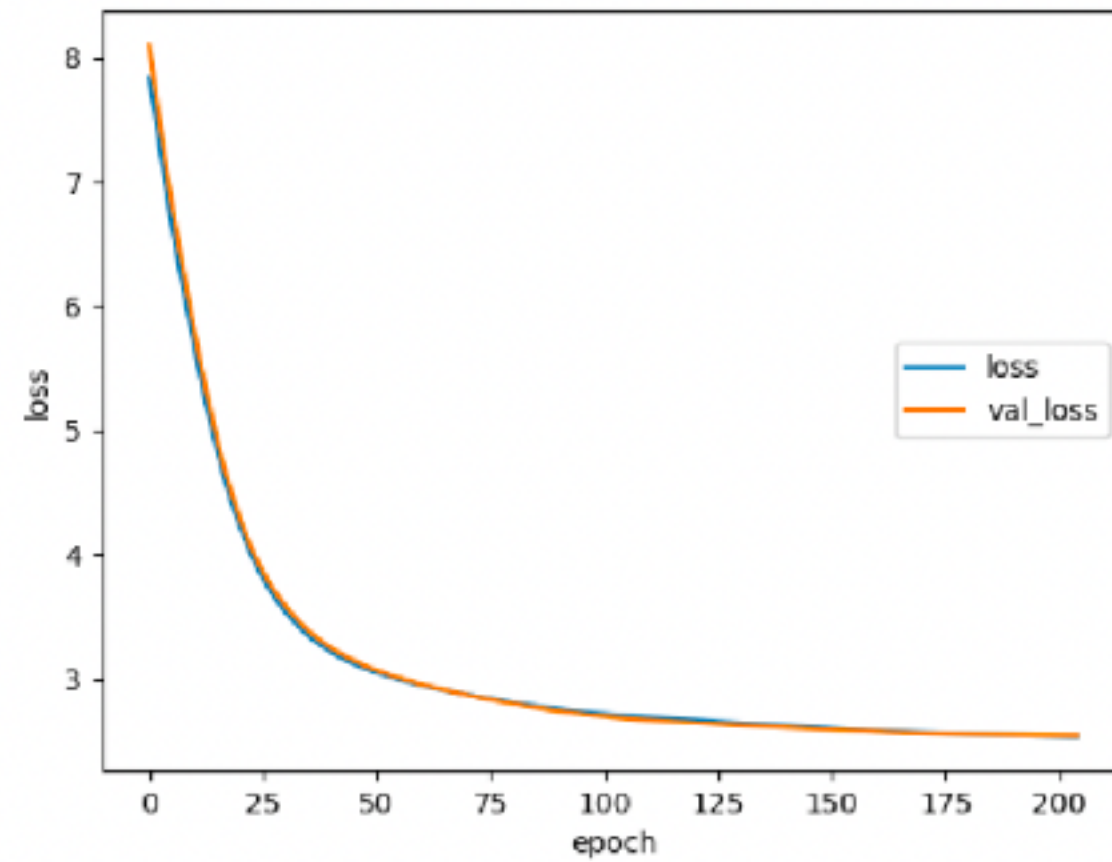


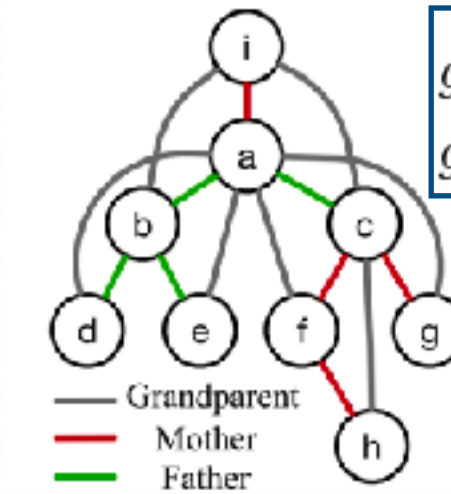
Figure 6: Learning curve in the *less-than* task.

1) $\text{even}(X) \text{ :- zero}(X).$
 2) $\text{even}(X) \text{ :- succ}(X,Z) \& \text{succ}(Z,Y) \& \text{even}(Y).$

The program in **Even** dataset

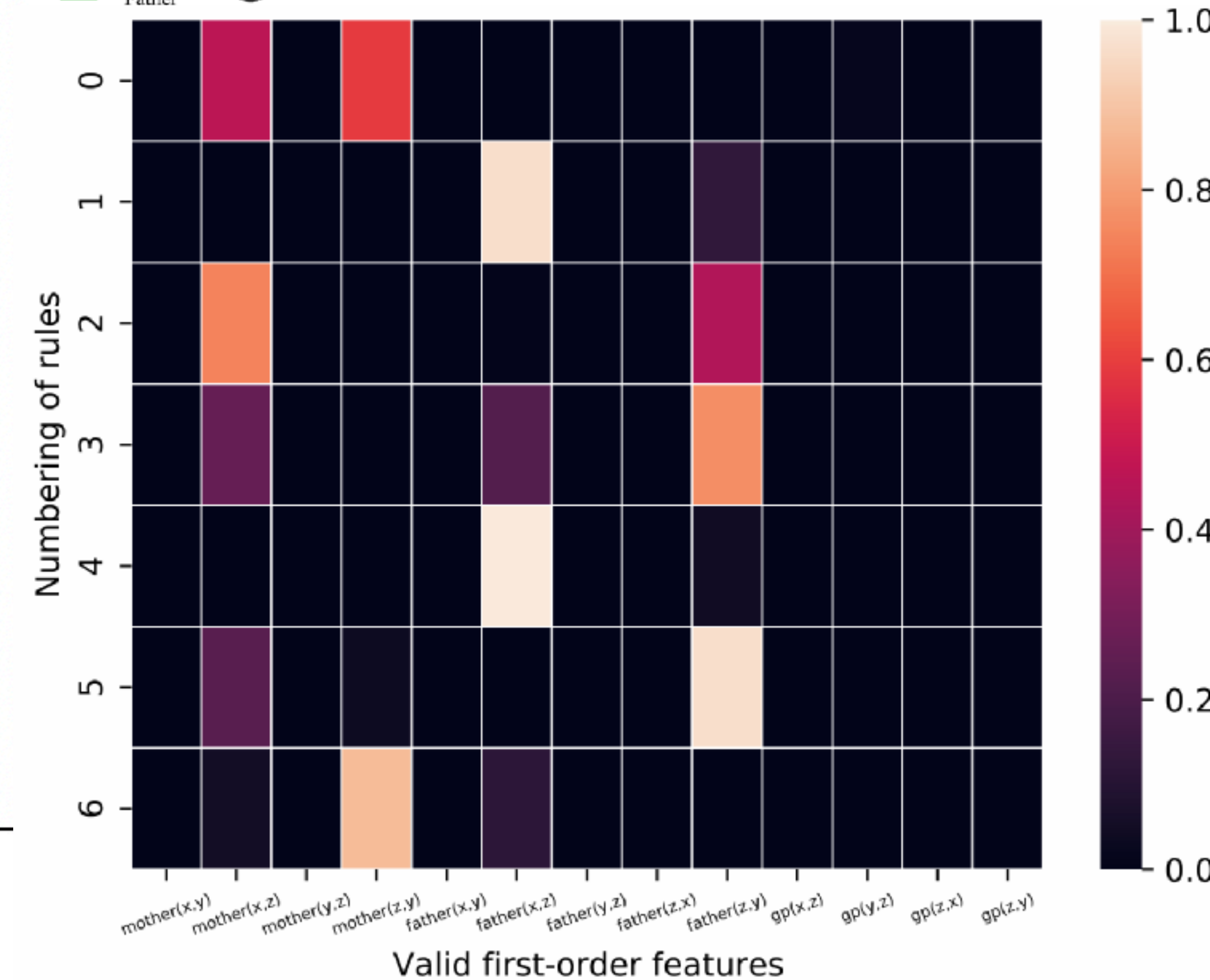
Table 4: The results on ILP datasets. The symbols \checkmark , $*$, and $-$ indicate that the accuracy of the generated logic program is equal to 100%, less than 100%, and equal to 0%, respectively.

Domain	Task	∂ ILP	NeuralLP	DFOL
Arithmetic	Predecessor	\checkmark	\checkmark	\checkmark
	Odd	\checkmark	$-$	\checkmark
	Even / Succ2 (10)	\checkmark	$-$	\checkmark
	Even / Succ2 (20)	$-$	$-$	\checkmark
	Lessthan	\checkmark	\checkmark	\checkmark
	Fizz	\checkmark	$-$	\checkmark
	Buzz	\checkmark	$-$	\checkmark
Lists	Member	\checkmark	$*$	\checkmark
	Length	\checkmark	$-$	\checkmark
	Family Tree			
Family Tree	Son	\checkmark	$*$	\checkmark
	Grandparent	\checkmark	\checkmark	\checkmark
	Husband	\checkmark	$-$	$-$
	Uncle	\checkmark	$-$	$-$
	Relatedness	\checkmark	$*$	\checkmark
	Father	\checkmark	$-$	\checkmark
	Graphs	Directed Edge	\checkmark	$*$
Adjacent to Red		\checkmark	$-$	\checkmark
Two Children		\checkmark	$-$	\checkmark
Graph Coloring (6)		\checkmark	$-$	\checkmark
Graph Coloring (10)		$-$	$-$	\checkmark
Connectedness		\checkmark	$*$	\checkmark
Cyclic		\checkmark	$-$	\checkmark



$g(X, Y) \leftarrow m(X, Z) \wedge m(Z, Y), g(X, Y) \leftarrow f(X, Z) \wedge f(Z, Y)$
 $g(X, Y) \leftarrow m(X, Z) \wedge f(Z, Y), g(X, Y) \leftarrow f(X, Z) \wedge m(Z, Y)$

The program in **Grandparents** dataset



The learned matrix M_P

Experimental Results

- From **fuzzy data**, we let $\varepsilon \sim N(0, \sigma^2)$
 - $p_i^+ = \min(1 - \varepsilon, 1)$
 - $p_i^- = \max(\varepsilon, 0)$
- For the **mislabeled data**, both positive and negative training examples are mislabeled:
 - μ ranging from 0.05 to 1 with step 0.05.

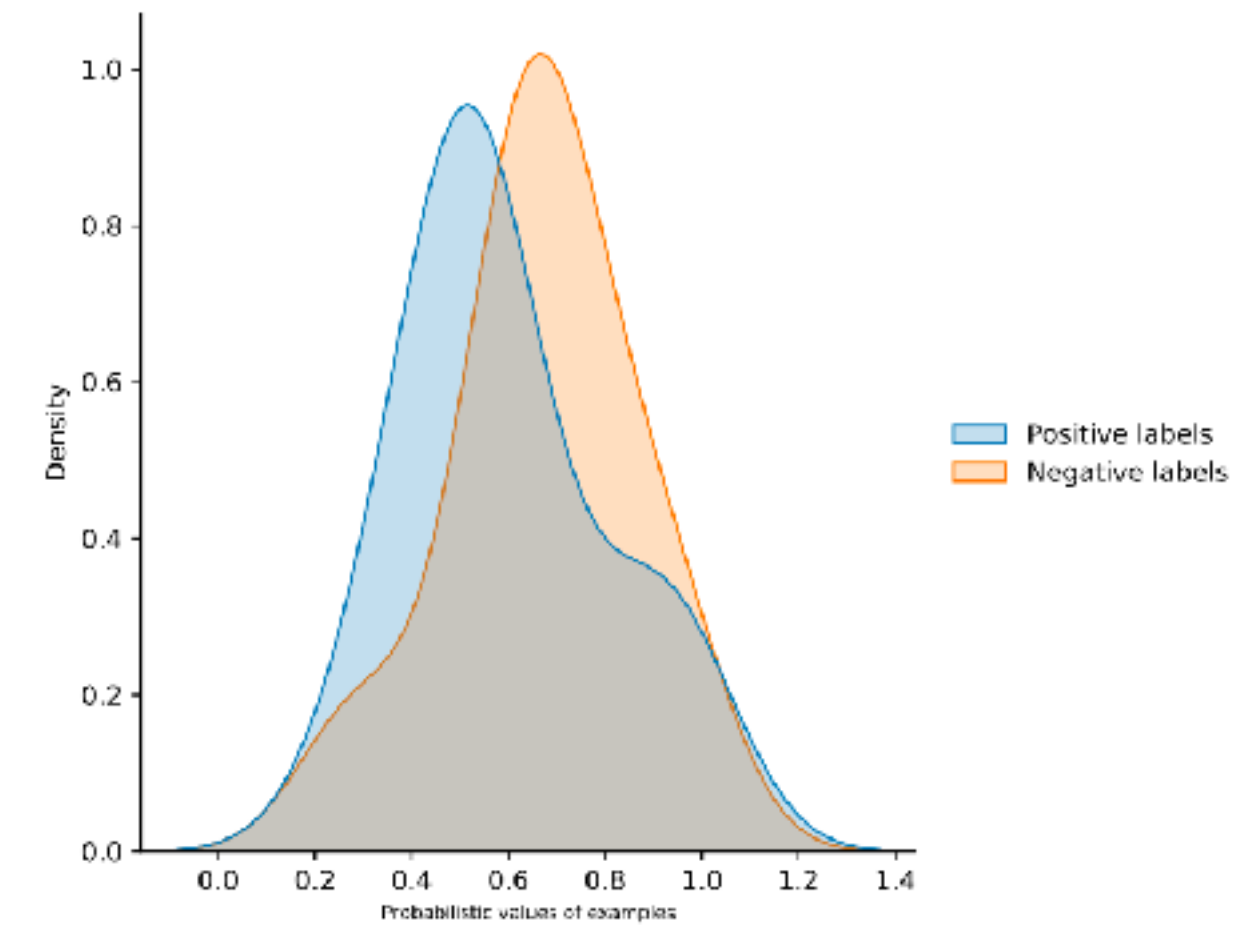


Figure 4: The distribution of examples in the *lessthan* task when $\sigma = 3$.

Table 1: The results on ambiguous datasets. The notations *Con* and *DE* represent *connectedness* and *directed edge* tasks, respectively.

	<i>Lt</i>	<i>Pre</i>	<i>Member</i>	<i>Son</i>	<i>Con</i>	<i>DE</i>
σ	3	3	3	3	2	3
μ	0.95	0.95	0.90	0.95	0.95	0.95

When $\tau_s = 1$ and accuracy is 100%, the largest standard deviation σ and mislabeling rate μ

Experimental Results

- From knowledge graph datasets,
 - Set $\tau_s = 0.3$;
 - Check the HIT@n and MRR on three knowledge graph datasets.

Table 3: The descriptions of the knowledge bases.

Dataset	#Object	#Relation	#Fact
Countries	252	2	1158
Nations	14	56	2565
UMLS	135	49	6529

```

1) locatedIn(X,Y) :- locatedIn(Z,Y)& neighborOf(X,Z)& neighborOf(Z,X). #(0.68, 714, 1053)
2) locatedIn(X,Y) :- locatedIn(Z,Y)& neighborOf(X,Z). #(0.68, 723, 1067)
3) locatedIn(X,Y) :- locatedIn(X,Z)& locatedIn(Z,Y). #(0.81, 187, 231)
4) locatedIn(X,Y) :- locatedIn(Z,Y)& neighborOf(Z,X). #(0.68, 723, 1068)
  
```

The program in **Country-S1** dataset

Table 2: Comparison on knowledge bases. The results in bold indicate the highest accuracy on the corresponding test datasets. The ACC@Sn represent the accuracy of the generated LP on the Sn subset of Countries. The ACC@h represent the accuracy of LP with the head predicate h. Besides, blo, int, neg, and intw denote the relations of blockpositionindex, intergovorgs3, negativecomm, and interacts_with.

Dataset	Metrics	NTP λ	NeuralLP	DFOL
Countries	ACC@S1	100.00	100.00	100.00
	ACC@S2	100.00	100.00	100.00
	ACC@S3	100.00	—	—
Nations	MRR	41.79	56.49	78.88
	HITS@1	41.79	52.49	73.88
	HITS@3	41.79	60.95	84.58
	HITS@10	41.79	61.19	85.07
	ACC@blo	100.00	50.00	100.00
	ACC@int	84.62	84.62	84.62
	ACC@neg	37.50	75.00	75.00
UMLS	MRR	30.03	66.69	74.96
	HITS@1	29.95	61.27	71.41
	HITS@3	30.11	72.31	78.82
	HITS@10	30.11	72.31	78.97
	ACC@isa	65.96	63.83	91.48
	ACC@intw	83.67	86.67	100.00

Conclusion and Future Work

- Conclusion
 - DFOL generates first-order logic programs usually with variable depth 1 or 2;
 - DFOL is a fast, precise, robust, scalable rule learner.
- Future Work
 - Larger variable depth
 - Support negation and function in logic programs

Thanks for your attention!

