

## Probabilistic Constraint Logic Theories

Inductive Constraint Logic (De Raedt and Van Laer): learning Constraint Logic Theories (sets of integrity constraints) from interpretations rather than from entailment.

In this work: *Probabilistic* Constraint Logic Theories (PCLTs) and a system (PASCAL) to perform discriminative learning of their structure and parameters from interpretations.

**Language** A PCLT is a set of probabilistic integrity constraints (ICs) of the form  $p_i :: L_1, \dots, L_b \rightarrow \exists(P_1); \dots; \exists(P_n); \forall \neg(N_1); \dots; \forall \neg(N_m)$

**Semantics** A PCLT  $T$  defines a probability distribution on the set  $W$  of ground constraint logic theories called *possible theories*; for each grounding of the body of each IC, we include the IC in a possible theory with probability  $p_i$  and we assume all groundings to be independent.

An IC's probability is the sum of the probabilities of the possible theories where a grounding of the constraint is present.

In each possible theory, An IC  $C$  is *true in an interpretation*  $I$  ( $I \models C$ ) if and only if, for each substitution  $\theta$  such that each literal in  $Body(C)\theta$  is ground and true in  $I$ , at least one disjunct in  $Head(C)\theta$  is true in  $I$ .

The probability  $P(\oplus|I)$  of the positive class given an interpretation  $I$ , the probability of a PCLT  $T$  satisfying  $I$ , is

$$P(\oplus|I) = \sum_{w \in W} P(\oplus, w|I) = \sum_{w \in W} P(\oplus|w, I)P(w|I) = \sum_{w \in W, M(\mathbf{BG} \cup I) \models w} P(w)$$

**Example** (from Bongard)

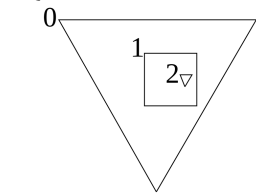
Background knowledge:

$in(A, B) \leftarrow inside(A, B)$

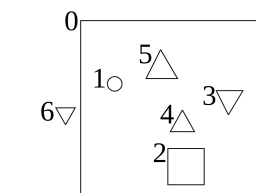
$in(A, D) \leftarrow inside(A, C), in(C, D)$

PCLT:

$\{C_1 = 0.5 :: triangle(T), square(S), in(T, S) \rightarrow false\}$



The body of  $C_1$  is true for the single substitution  $T/2$  and  $S/1$  thus  $m_1 = 1$  and  $P(\oplus|I) = 0.5$ .



The body of  $C_1$  is true for three pairs (triangle, square) thus  $m_1 = 3$  and  $P(\oplus|I) = 0.125$ .

## PCLT parameter learning

**Problem** Given

- a PCLT theory  $T$
- a set  $\mathcal{I}^+ = \{I_1, \dots, I_Q\}$  of positive interpretations
- a set  $\mathcal{I}^- = \{I_{Q+1}, \dots, I_R\}$  of negative interpretations
- a normal logic program  $\mathbf{BG}$  (background knowledge)

find the parameters of  $T$  such that the likelihood

$$L = \prod_{q=1}^Q P(\oplus|I_q) \prod_{r=Q+1}^R P(\ominus|I_r)$$

is maximized. The likelihood is given by the probability that the example labels are observed for each example.

The equation  $\frac{\partial L}{\partial p_i} = 0$  does not admit a closed form solution, so we must use optimization to find the maximum of  $L$ .

We can optimize the likelihood with gradient descent, where weights are updated using the formula

$$\mathbf{p}_{n+1} = \mathbf{p}_n - \epsilon \nabla_{\mathbf{p}} L(\mathbf{p}) = \mathbf{p}_n - \epsilon \nabla_{\mathbf{p}} \frac{\partial L}{\partial \mathbf{p}}$$

where  $\epsilon$  is the learning rate defining the size of the step done by gradient descent along the gradient and  $\mathbf{p}$  is the vector containing the parameters  $p_i$ , or with a second order method such as Limited-memory BFGS (L-BFGS).

Experiments show that gradient descent outperforms L-BFGS in most cases in terms of area under the PR and ROC curves, and execution time.

## PASCAL: PCLT structure learning

**Problem** Given

- a set  $\mathcal{I}^+ = \{I_1, \dots, I_Q\}$  of positive interpretations
- a set  $\mathcal{I}^- = \{I_{Q+1}, \dots, I_R\}$  of negative interpretations
- a normal logic program  $\mathbf{BG}$  (background knowledge)
- a language bias

find a PCLT  $T$  that maximizes the likelihood

$$L = \prod_{q=1}^Q P(\oplus|I_q) \prod_{r=Q+1}^R P(\ominus|I_r)$$

PASCAL solves this problem by first identifying good candidate ICs and then searching for a theory guided by the log likelihood (LL) of the data.

Parameters:

- $MLB$ , the maximum number of literals in the body of ICs;
- $MD$ , the maximum number of disjuncts in the head of ICs;
- $MLP$  and  $MLN$ , the maximum number of literals allowed in a  $P$  disjunct and a  $N$  disjunct respectively.

```

1: function PASCAL( $\mathcal{I}^+, \mathcal{I}^-, \mathbf{BG}, NC, MLB, MD, MLP, MLN, BeamSize, MaxSteps$ )
2:   Steps = 1
3:   Beam  $\leftarrow$  {false  $\leftarrow$  true,  $-\infty$ } ▷ Empty IC
4:   repeat
5:     NewBeam = []
6:     while Beam is not empty do ▷ ICs search
7:       Remove the first IC ( $C, LL$ ) from Beam
8:       Ref  $\leftarrow$  all refinements of  $C$  respecting  $MLB, MD, MLP$ , and  $MLN$ 
9:       for all  $C' \in Ref$  do
10:        ( $\{C''\}, LL''$ )  $\leftarrow$  LEARNPARAMS( $\{C'\}, \mathcal{I}^+, \mathcal{I}^-, \mathbf{BG}$ ) ▷ gradient descent
11:        NewBeam  $\leftarrow$  INSERT( $\{C'', LL''\}, NewBeam$ )
12:        if size(NewBeam) > BeamSize then
13:          Remove the last element of NewBeam
14:        end if
15:      end for
16:    end while
17:    Beam  $\leftarrow$  NewBeam
18:    Steps = Steps + 1
19:  until Steps > MaxSteps
20:  T  $\leftarrow$   $\emptyset$ , LL  $\leftarrow$   $-\infty$  ▷ Theory search
21:  repeat
22:    Remove the first couple ( $C, LL$ ) from Beam
23:    ( $T', LL'$ )  $\leftarrow$  LEARNPARAMS( $T \cup \{C\}, \mathcal{I}^+, \mathcal{I}^-, \mathbf{BG}$ )
24:    if  $LL' > LL$  then
25:      T  $\leftarrow$  T', LL  $\leftarrow$  LL'
26:    end if
27:  until Beam is empty or T contains NC ICs
28:  return T
29: end function

```

## Experimental results

We compared PASCAL with the PLP algorithms LIFTCOVER, SLIPCOVER and LEMUR, the MLNs algorithms MLN-BC/MLN-BT, the (probabilistic) relational classifier TILDE.

Average AUC-PR

Dataset	SLIPCOVER	LIFT-EM	LIFT-LBFGS	PASCAL	TILDE
Bupa	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	0.420
Carcinogenesis	0.745	0.672	0.561	<b>0.770</b>	0.707
Financial	0.173	0.126	0.187	<b>0.317</b>	0.123
Mondial	<b>0.776</b>	0.763	0.723	0.652	0.650
Mutagen.	0.920	<b>0.971</b>	0.725	0.902	0.851
Pyrimidine	0.956	<b>1</b>	0.819	0.990	0.769
Sisya	<b>0.708</b>	0.706	0.706	0.622	0.621
Sisyb	<b>0.287</b>	0.286	0.286	0.286	0.286
Triazine	0.560	0.734	0.760	<b>0.855</b>	0.685
Yeast	0.428	0.502	0.448	0.469	<b>0.588</b>
Bongard	0.899	0.966	<b>0.970</b>	0.635	0.300

Dataset	MLN-BC			MLN-BT		PASCAL
	LEMUR	MLN-BC	samp.	MLN-BT	samp.	
Carcinogenesis	0.691	0.619	0.633	0.503	0.494	<b>0.770</b>
Mondial	<b>0.864</b>	0.585	0.742	0.735	0.781	0.652
Mutagenesis	<b>0.952</b>	0.690	0.831	0.872	-	0.902

Average AUC-ROC

Dataset	SLIPCOVER	LIFT-EM	LIFT-LBFGS	PASCAL	TILDE
Bupa	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>	0.500
Carcinogenesis	0.695	<b>0.766</b>	0.472	0.763	0.667
Financial	0.568	0.432	0.535	<b>0.745</b>	0.478
Mondial	0.630	<b>0.663</b>	0.643	0.495	0.500
Mutagen.	0.826	<b>0.931</b>	0.649	0.806	0.778
Pyrimidine	0.925	<b>1</b>	0.850	0.993	0.815
Sisya	0.719	0.372	<b>0.721</b>	0.502	0.499
Sisyb	<b>0.500</b>	<b>0.500</b>	<b>0.500</b>	<b>0.500</b>	<b>0.500</b>
Triazine	0.544	0.713	0.760	<b>0.803</b>	0.600
Yeast	0.733	0.786	0.721	<b>0.794</b>	0.718
Bongard	0.944	0.975	<b>0.987</b>	0.749	0.500

Dataset	MLN-BC			MLN-BT		PASCAL
	LEMUR	MLN-BC	samp.	MLN-BT	samp.	
Carcinogenesis	0.691	0.619	0.633	0.503	0.494	<b>0.770</b>
Mondial	<b>0.864</b>	0.585	0.742	0.735	0.781	0.652
Mutagenesis	<b>0.952</b>	0.690	0.831	0.872	-	0.902

The three datasets where PASCAL performs well - Triazine, Financial and Carcinogenesis - have a small number of examples but a large number of different predicates, possibly indicating that the expressive language bias is beneficial when the dataset is not very big but has a rich structure.