

Learning the Parameters of Deep Probabilistic Logic Programs

Arnaud Nguembang Fadja Fabrizio Riguzzi Evelina Lamma

Dipartimento di Ingegneria – University of Ferrara

Dipartimento di Matematica e Informatica – University of Ferrara
[arnaud.nguembafadja, fabrizio.riguzzi, evelina.lamma]@unife.it



Introduction

- Probabilistic logic programming is a powerful tool for reasoning with uncertain relational models
- Learning probabilistic logic programs is expensive due to the high cost of inference.
- We consider a restriction of the language of Logic Programs with Annotated Disjunctions called **hierarchical PLP** in which clauses and predicates are hierarchically organized.
- Inference is then much cheaper.

PLP under the Distribution Semantics

- A PLP language under the distribution semantics with a general syntax is **Logic Programs with Annotated Disjunctions** (LPADs)
- Heads of clauses are disjunctions in which each atom is annotated with a probability.
- LPAD T with n clauses: $T = \{C_1, \dots, C_n\}$.
- Each clause C_i takes the form:

$$h_{i1} : \pi_{i1}; \dots; h_{iv_i} : \pi_{iv_i} :- b_{i1}, \dots, b_{iu_i}$$

- Each grounding $C_i\theta_j$ of a clause C_i corresponds to a random variable X_{ij} with values $\{1, \dots, v_i\}$
- The random variables X_{ij} are independent of each other.

Example

- UW-CSE domain:

advisedby(A, B) : 0.3 : –

student(A), professor(B), project(C, A), project(C, B).

advisedby(A, B) : 0.6 : –

student(A), professor(B), ta(C, A), taughtby(C, B).

Hierarchical PLP

- We want to compute the probability of atoms for a predicate r : $r(\mathbf{t})$, where \mathbf{t} is a vector of constants.
- $r(\mathbf{t})$ can be an example in a learning problem and r a **target predicate**.
- A specific form of an LPADs defining r in terms of the input predicates.
- The program defined r using a number of input and **hidden predicates** disjoint from input and target predicates.
- Each rule in the program has a single head atom annotated with a probability.
- The program is hierarchically defined so that it can be divided into layers.
- Each layer contains a set of hidden predicates that are defined in terms of predicates of the layer immediately below or in terms of input predicates.

Hierarchical PLP

- Generic clause C :

$$C = p(\mathbf{X}) : \pi : - \phi(\mathbf{X}, \mathbf{Y}), b_1(\mathbf{X}, \mathbf{Y}), \dots, b_m(\mathbf{X}, \mathbf{Y})$$

where $\phi(\mathbf{X}, \mathbf{Y})$ is a conjunction of literals for the input predicates using variables \mathbf{X}, \mathbf{Y} .

- $b_i(\mathbf{X}, \mathbf{Y})$ for $i = 1, \dots, m$ is a literal built on a hidden predicate.
- \mathbf{Y} is a possibly empty vector of variables existentially quantified with scope the body.
- Literals for hidden predicates must use the whole set of variables \mathbf{X}, \mathbf{Y} .
- The predicate of each $b_i(\mathbf{X}, \mathbf{Y})$ does not appear elsewhere in the body of C or in the body of any other clause.

Hierarchical PLP

- A generic program defining r is thus:

$$C_1 = r(\mathbf{X}) : \pi_1 \quad :- \quad \phi_1, b_{11}, \dots, b_{1m_1}$$

...

$$C_n = r(\mathbf{X}) : \pi_n \quad :- \quad \phi_n, b_{n1}, \dots, b_{nm_n}$$

$$C_{111} = r_{11}(\mathbf{X}) : \pi_{111} \quad :- \quad \phi_{111}, b_{1111}, \dots, b_{111m_{111}}$$

...

$$C_{11n_{11}} = r_{11}(\mathbf{X}) : \pi_{11n_{11}} \quad :- \quad \phi_{11n_{11}}, b_{11n_{11}1}, \dots, b_{11n_{11}m_{11n_{11}}}$$

...

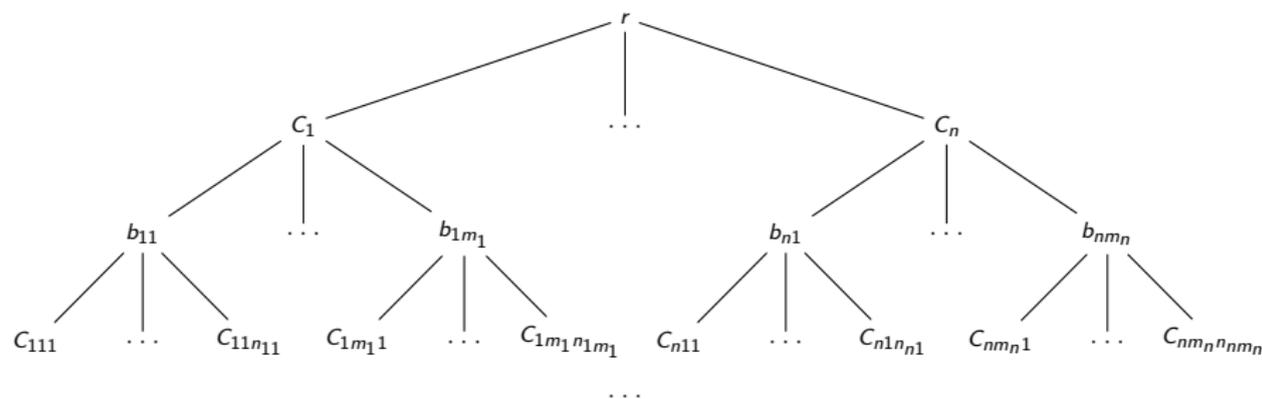
$$C_{n11} = r_{n1}(\mathbf{X}) : \pi_{n11} \quad :- \quad \phi_{n11}, b_{n111}, \dots, b_{n11m_{n11}}$$

...

$$C_{n1n_{n1}} = r_{n1}(\mathbf{X}) : \pi_{n1n_{n1}} \quad :- \quad \phi_{n1n_{n1}}, b_{n1n_{n1}1}, \dots, b_{n1n_{n1}m_{n1n_{n1}}}$$

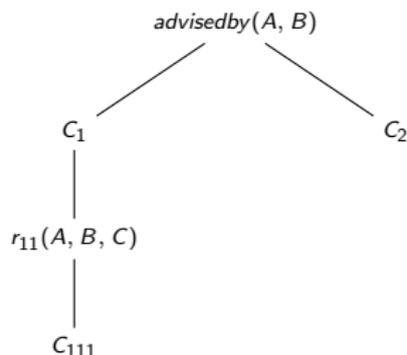
...

Program Tree



Example

- $C_1 = \text{advisedby}(A, B) : 0.3 :-$
 $\text{student}(A), \text{professor}(B), \text{project}(C, A), \text{project}(C, B),$
 $r_{11}(A, B, C).$
- $C_2 = \text{advisedby}(A, B) : 0.6 :-$
 $\text{student}(A), \text{professor}(B), \text{ta}(C, A), \text{taughtby}(C, B).$
- $C_{111} = r_{11}(A, B, C) : 0.2 :-$
 $\text{publication}(D, A, C), \text{publication}(D, B, C).$



Inference

- Generate the grounding.
- Each ground probabilistic clause is associated with a random variable whose probability of being true is given by the parameter of the clause and that is independent of all the other clause random variables.
- Ground clause $C_{\mathbf{p}i} = a_{\mathbf{p}} : \pi_{\mathbf{p}i} :- b_{\mathbf{p}i1}, \dots, b_{\mathbf{p}im_{\mathbf{p}}}$. where \mathbf{p} is a path in the program tree
- $P(b_{\mathbf{p}i1}, \dots, b_{\mathbf{p}im_{\mathbf{p}}}) = \prod_{i=k}^{m_{\mathbf{p}}} P(b_{\mathbf{p}ik})$ and $P(b_{\mathbf{p}ik}) = 1 - P(a_{\mathbf{p}ik})$ if $b_{\mathbf{p}ik} = \neg a_{\mathbf{p}ik}$.
- If a is a literal for an input predicate, then $P(a) = 1$ if a belongs to the example interpretation and $P(a) = 0$ otherwise.

Inference

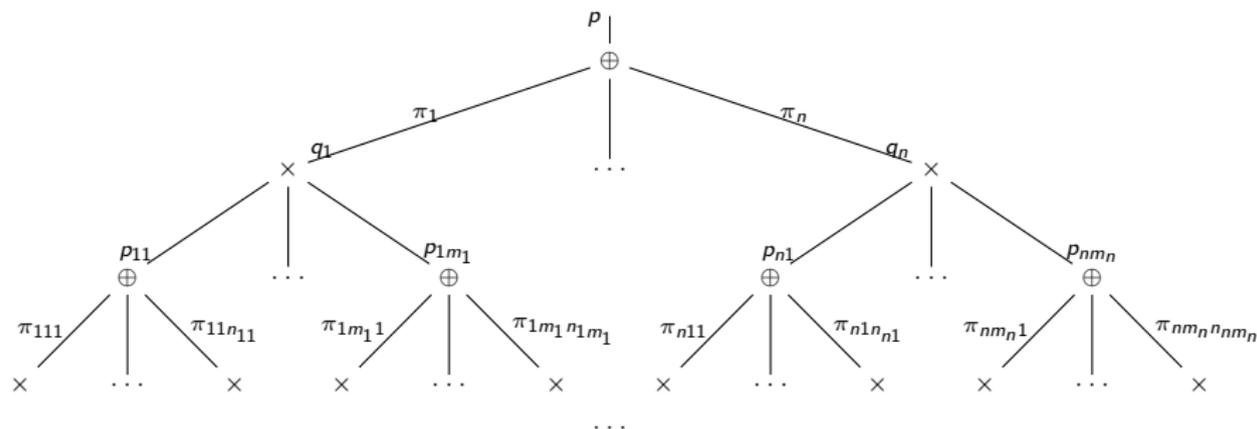
- Hidden predicates: to compute $P(a_p)$ we need to take into account the contribution of every ground clause for the predicate of a_p .
- Suppose these clauses are $\{C_{p1}, \dots, C_{pn}\}$.
- For one clause, $P(a_p) = \pi_{p1} \cdot P(\text{body}(C_{p1}))$
- If we have two clauses,
$$P(a_p) = 1 - (1 - \pi_{p1} \cdot P(\text{body}(C_{p1})) \cdot (1 - \pi_{p2} \cdot P(\text{body}(C_{p2}))))$$
- $p \oplus q \triangleq 1 - (1 - p) \cdot (1 - q)$.
- This operator is commutative and associative:

$$\bigoplus_i p_i = 1 - \prod_i (1 - p_i)$$

- The operators \times and \oplus are respectively the t-norm and t-conorm of the product fuzzy logic [Hajek 98]: **product t-norm** and **probabilistic sum**.

Inference

- If the probabilistic program is ground, the probability of the example atom can be computed with the arithmetic circuit:

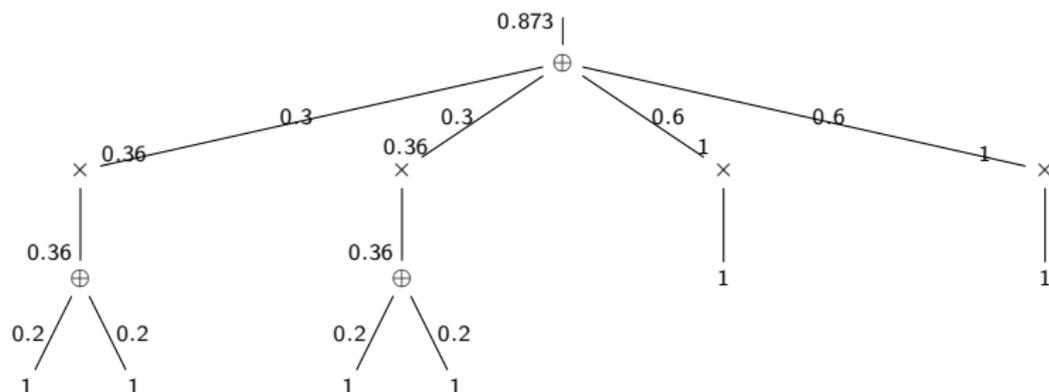
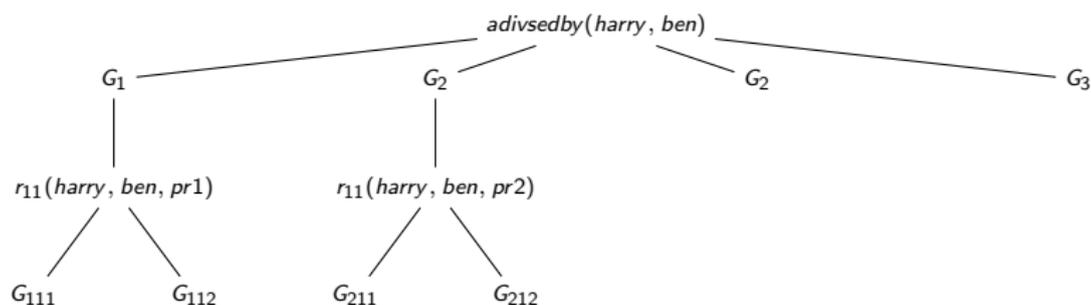


- The arithmetic circuit can be interpreted as a deep neural network where nodes have the activation functions \times and \oplus

Example

- $G_1 = \text{advisedby}(\text{harry}, \text{ben}) : 0.3 : -$
 $\text{student}(\text{harry}), \text{professor}(\text{ben}), \text{project}(\text{pr1}, \text{harry}),$
 $\text{project}(\text{pr1}, \text{ben}), r_{11}(\text{harry}, \text{ben}, \text{pr1}).$
- $G_2 = \text{advisedby}(\text{harry}, \text{ben}) : 0.3 : -$
 $\text{student}(\text{harry}), \text{professor}(\text{ben}), \text{project}(\text{pr2}, \text{harry}),$
 $\text{project}(\text{pr2}, \text{ben}), r_{11}(\text{harry}, \text{ben}, \text{pr2}).$
- $G_3 = \text{advisedby}(\text{harry}, \text{ben}) : 0.6 : -$
 $\text{student}(\text{harry}), \text{professor}(\text{ben}), \text{ta}(\text{c1}, \text{harry}), \text{taughtby}(\text{c1}, \text{ben}).$
- $G_4 = \text{advisedby}(\text{harry}, \text{ben}) : 0.6 : -$
 $\text{student}(\text{harry}), \text{professor}(\text{ben}), \text{ta}(\text{c2}, \text{harry}), \text{taughtby}(\text{c2}, \text{ben}).$
- $G_{111} = r_{11}(\text{harry}, \text{ben}, \text{pr1}) : 0.2 : -$
 $\text{publication}(\text{p1}, \text{harry}, \text{pr1}), \text{publication}(\text{p1}, \text{ben}, \text{pr1}).$
- $G_{112} = r_{11}(\text{harry}, \text{ben}, \text{pr1}) : 0.2 : -$
 $\text{publication}(\text{p2}, \text{harry}, \text{pr1}), \text{publication}(\text{p2}, \text{ben}, \text{pr1}).$
- $G_{211} = r_{11}(\text{harry}, \text{ben}, \text{pr2}) : 0.2 : -$
 $\text{publication}(\text{p3}, \text{harry}, \text{pr2}), \text{publication}(\text{p3}, \text{ben}, \text{pr2}).$
- $G_{212} = r_{11}(\text{harry}, \text{ben}, \text{pr2}) : 0.2 : -$
 $\text{publication}(\text{p4}, \text{harry}, \text{pr2}), \text{publication}(\text{p4}, \text{ben}, \text{pr2}).$

Example



Example

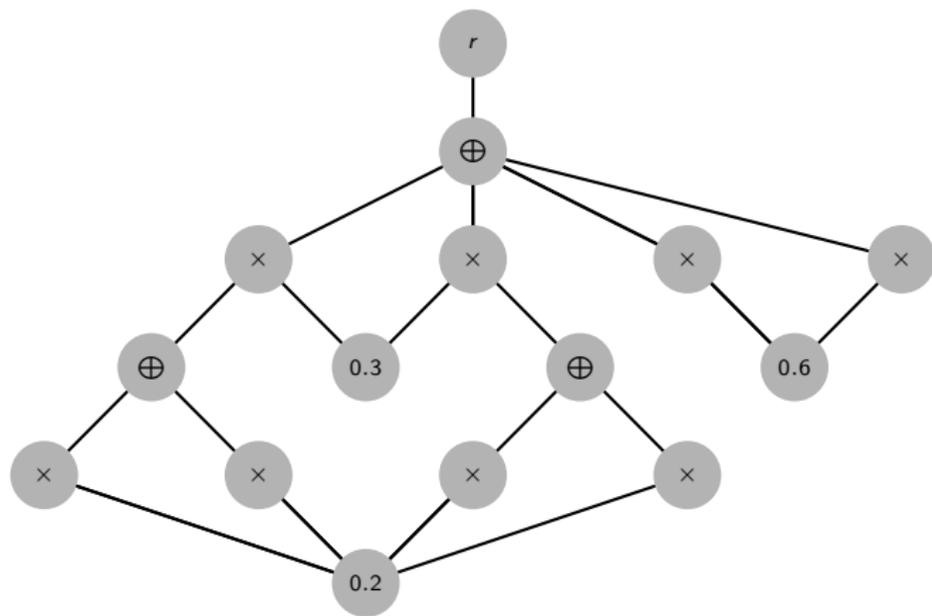


Figure: Arithmetic circuit.

Parameter Learning

- Given a Hierarchical PLP T with parameters Π , an interpretation I defining input predicates and a training set $E = \{e_1, \dots, e_M, \text{not } e_{M+1}, \dots, \text{not } e_N\}$ find the values of Π that maximize the log likelihood:

$$\arg \max_{\Pi} \sum_{i=1}^M \log P(e_i) + \sum_{i=M+1}^N \log(1 - P(e_i)) \quad (1)$$

where $P(e_i)$ is the probability assigned to e_i by $T \cup I$.

- Maximizing the log likelihood can be equivalently seen as minimizing the sum of *cross entropy errors* err_i for all the examples

$$err_i = -y_i \log(p_i) - (1 - y_i) \log(1 - p_i) \quad (2)$$

where $y_i = 1$ for positive example, $y_i = 0$ otherwise and p_i the probability that the atom is true.

Parameter Learning

- Partial derivative of the error with respect to each node $v(n)$:

$$\frac{\partial \text{err}}{\partial v(n)} = \begin{cases} -\frac{1}{v(r)} d(n) & \text{if } e \text{ is positive,} \\ \frac{1}{1-v(r)} d(n) & \text{if } e \text{ negative.} \end{cases}$$

where

$$d(n) = \begin{cases} d(p_n) \frac{v(p_n)}{v(n)} & \text{if } n \text{ is a } \oplus \text{ node,} \\ d(p_n) \frac{1-v(p_n)}{1-v(n)} & \text{if } n \text{ is a } \times \text{ node} \\ \sum_{p_n} d(p_n) \cdot v(p_n) \cdot (1 - \Pi_i) & \text{if } n \text{ is a leaf node } \Pi_i \\ -d(p_n) & p_n = \text{not}(n) \end{cases} \quad (3)$$

and $v(n)$, p_n are respectively the value and the parent of the node n .

Parameter Learning

- Build the ACs and initialize the parameters and the gradients.
- Perform the forward pass by computing the output of each node ($v(n)$) in the AC.
- Compute the gradient of the error w.r.t the output and back-propagate.
- Update the parameters using Adam optimizer.
- Until convergence or a certain condition is satisfied.

Conclusion and Future Work

- Conclusion
 - hierarchical PLP: a restriction of the language of LPADs that allows to perform inference quickly using a simple and cheap dynamic programming algorithm such as PITA(IND,IND).
 - Programs can be seen as arithmetic circuits/neural networks.
 - Parameters can be trained by gradient descent and back-propagation.
- Future work
 - Perform experiments.
 - Parameter learning by Expectation Maximization.
 - Perform also structure learning.



**THANKS FOR
LISTENING
AND
ANY
QUESTIONS ?**