

Neuro-Symbolic Integration (NeSy)

Andrea Passerini
andrea.passerini@unitn.it

Advanced Topics in Machine Learning and Optimization

PROs

- Efficient processing of high-dimensional data
- Robust to noise and ambiguity
- Does not require extensive background knowledge and feature engineering

CONs

- Data hungry (large training sets needed)
- Non-interpretable models and predictions
- Hard to incorporate complex domain knowledge

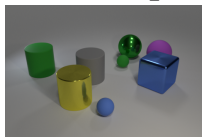
PROs

- Expressive, can formalize complex domain knowledge
- Interpretable, inference can be explained in terms reasoning steps (proofs)
- Can generalize from few examples

CONS

- Inefficient, inference is typically expensive
- No support for noise or ambiguity
- Difficult to deal with high-dimensional data

Neuro-Symbolic Integration (NeSy)



Q: How many objects are both right of the green cylinder and have the same material as the small blue ball?

A: 3

Best of both worlds

- Deep networks for low-level data processing and “atomic” predictions
- Symbolic approaches for reasoning on top of atomic predictions
- Probabilities (or scores) for dealing with uncertainty

Image from Mao et al. 2019

Dimensions: directed vs undirected models

Directed models

- Generalize Bayesian Networks to deal with (first-order) logic
- Generalize Logic Programs to deal with probabilities
- Incorporate Neural “primitives” (e.g., predicates)

Undirected models

- Generalize Markov Networks to deal with (first-order) logic
- Enforce logical constraints over neural predictions
- Relax logical constraints to deal with uncertainty

Dimensions: integration vs regularization

Integration

- Neural primitives inside reasoning framework (typically logic program)
- Differentiability via probability of worlds or proof score.

Regularization

- Logical Constraints are used as regularizers for neural network training
- Differentiability by relaxed constraints or consistency in expectation

Dimensions: semantics

Probabilistic semantics

- Extends Boolean logic with probabilities
- Defines a probability distribution over possible worlds
- Allows to perform inference under uncertainty (expensive)

Fuzzy semantics

- Relax Boolean variables in $[0,1]$ interval
- Relies on t-norms for relaxing Boolean connectives
- Efficient inference, Boolean semantics not preserved

Semantic-based Regularization

Setting

- Model problems with multiple related predictions
- Incorporate knowledge as constraints over related predictions

Solution

- Model each prediction task with a statistical learner (kernel machine, neural network)
- Represent constraints over predictions in fuzzy logic
- Combine regularization with loss on fuzzy constraint satisfaction (including label supervision)

Semantic-based Regularization: Fuzzy logic

Boolean	Gödel	Product	Łukasiewicz
$X \wedge Y$	$\min(X, Y)$	$X Y$	$\max(0, X + Y - 1)$
$X \vee Y$	$\max(X, Y)$	$1 - (1 - X)(1 - Y)$	$\min(1, X + Y)$
$\neg X$	$1 - X$	$1 - X$	$1 - X$

Fuzzy logic

- Boolean variables relaxed into real variables in $[0, 1]$.
- Conjunction relaxed using **t-norm**
- Disjunction relaxed using **t-conorm**
- Existential quantifier relaxed as maximum (over dataset)
- Universal quantifier relaxed as minimum (over dataset, usually replaced by average)

$$\mathcal{L}(\mathbf{f}, \Phi) = \sum_{k=1}^{|\mathbf{f}|} \|f_k\|^2 + \sum_{h=1}^{|\Phi|} \lambda_h (1 - \hat{\Phi}_h(\mathbf{f}))$$

Objective function

- \mathbf{f} is a vector of parameterized predictors (one per task)
- Φ is a set of logic formulas (the constraints)
- $\|f_k\|$ is the norm of f_k (e.g. norm of the weights for kernel machines)
- λ_h is a weight associated to constraint h
- $\hat{\Phi}_h$ is the fuzzy version of formula Φ_h

Semantic-based Regularization: example

positive supervision

manifold regularisation

$$F := \forall d P_A(d) \Rightarrow A(d)$$

$$F_R := \forall d \forall d' R(d, d') \Rightarrow ((A(d) \wedge A(d')) \vee (\neg A(d) \wedge \neg A(d')))$$

$$C = \{d_1, d_2\}$$

Evidence Predicate
Groundings

$$P_A(d_1) = 1$$

$$R(d_1, d_2) = 1$$

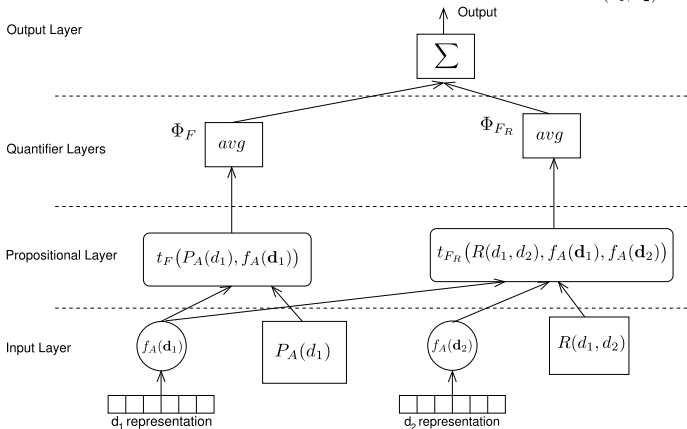


Image adapted from Diligenti et al., 2017

Semantic-based Regularization: learning

$$\frac{\partial \mathcal{L}(\mathbf{f}, \Phi)}{\partial \mathbf{w}_{k,j}} = \frac{\partial \|\mathbf{f}_k\|^2}{\partial \mathbf{w}_{k,j}} + \sum_{h=1}^{|\Phi|} \lambda_h \frac{\partial (1 - \hat{\Phi}_h)}{\partial \hat{\Phi}_h} \cdot \left(\sum_{t_{\Phi_h}} \frac{\partial t_{\Phi_h}}{\partial \mathbf{f}_k} \cdot \frac{\partial \mathbf{f}_k}{\partial \mathbf{w}_{k,j}} \right)$$

Gradient-based learning

- $w_{k,j}$ is a parameter of a predictor f_k
- t_{Φ_h} is a grounding of formula Φ_h

Note

Learning problem is convex if:

- f_k are kernel machines (or similar)
- A convex fragment of the Łukasiewicz logic is used

Semantic-based Regularization: MAP inference

$$\mathcal{L}(\bar{\mathbf{f}}(\mathcal{X}), \mathbf{f}(\mathcal{X})) = \frac{1}{2} \|\bar{\mathbf{f}}(\mathcal{X}) - \mathbf{f}(\mathcal{X})\|^2 + \sum_h \lambda_h (1 - \hat{\Phi}_h(\bar{\mathbf{f}}(\mathcal{X})))$$

Gradient-based MAP inference

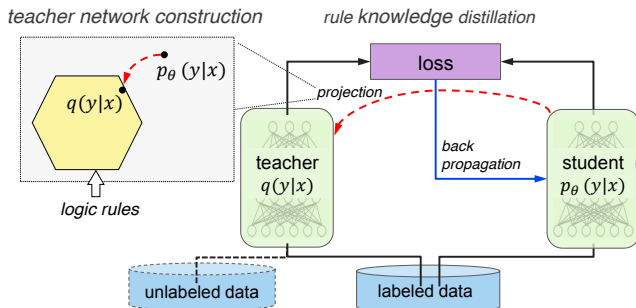
- \mathcal{X} set of (related) test examples
- $\mathbf{f}(\mathcal{X})$ set of independent predictions over test examples
- $\bar{\mathbf{f}}(\mathcal{X})$ set of collective predictions over test examples (accounting for constraints)
- Inference of $\bar{\mathbf{f}}(\mathcal{X})$ is performed by gradient descent:

$$\frac{\mathcal{L}(\bar{\mathbf{f}}(\mathcal{X}), \mathbf{f}(\mathcal{X}))}{\partial \bar{\mathbf{f}}_k(\mathcal{X}_i)} = \bar{\mathbf{f}}_k(\mathcal{X}_i) - \mathbf{f}_k(\mathcal{X}_i) + \sum_h \lambda_h \left(\frac{\partial 1 - \hat{\Phi}_h(\bar{\mathbf{f}}(\mathcal{X}))}{\partial \bar{\mathbf{f}}_k(\mathcal{X}_i)} \right)$$

dimensions

- **Undirected model:** constraints as set of FOL formulas (probabilistic variant as deep Markov Logic Network exists)
- **Regularization approach:** soft consistency is a regularization term in training loss
- **Fuzzy semantics:** fuzzy logic is employed as relaxation

Knowledge distillation



Teacher-student distillation

- Student learns to fit data and satisfy rules
- Teacher “shows” student how to change predictions to satisfy rules (projection in feasible space)
- Student should learn to implicitly satisfy rules (no rule enforcement at prediction time)

$$\mathcal{L}(\mathcal{D}; \Phi) = \sum_{(\mathbf{x}_n, \mathbf{y}_n) \in \mathcal{D}} (1 - \pi) \ell(\mathbf{y}_n, f_p(\mathbf{x}_n)) + \pi \ell(f_q(\mathbf{x}_n), f_p(\mathbf{x}_n))$$

Iterative procedure

- $f_p(\mathbf{x}_n)$ are the student predictions for \mathbf{x}_n (i.e., according to $p_\theta(\mathbf{y}|\mathbf{x}_n)$)
- $f_q(\mathbf{x}_n)$ is the teacher projection of those predictions in the feasible space Φ (i.e., according to $q(\mathbf{y}|\mathbf{x}_n)$)
- π is a parameter trading-off data fitting and constraint satisfaction (possibly on unlabelled data too)
- At each iteration θ is updated minimizing the loss

Knowledge distillation: teacher projection

$$\begin{aligned} \min_{q, \xi} \quad & KL(q(Y|X) || p_{\theta}(Y|X)) + C \sum_h \sum_g \xi_{h,g} \\ \text{s.t.} \quad & \lambda_h (1 - E_q[\hat{\Phi}_{h,g}(X, Y)]) \leq \xi_{h,g} \end{aligned}$$

Projection as constrained optimization

- KL divergence between student and teacher predictions
- $\hat{\Phi}_{h,g}(X, Y)$ is the g -th grounding of a fuzzy version of formula Φ_h on (X, Y) .
- $E_q[\hat{\Phi}_{h,g}(X, Y)]$ is satisfaction of $\hat{\Phi}_{h,g}(X, Y)$ in expectation over $q(Y|X)$.
- λ_h is the weight of formula Φ_h
- $\xi_{h,g}$ is a slack variable to penalize unsatisfied constraints
- C is a parameter trading-off divergence with student prediction and satisfaction of formulas

$$q^*(Y|X) \propto p_\theta(Y|X) \cdot \exp \left(\sum_h \sum_g C \lambda_h (1 - \hat{\Phi}_{h,g}(X, Y)) \right)$$

Closed form solution

- The constrained optimization problem has a closed form solution.
- The normalization term is computed by dynamic programming if relationship between constraints allows for it, or approximated with sampling approaches otherwise.

dimensions

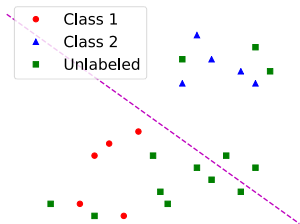
- **Undirected model:** constraints as set of FOL formulas
- **Regularization approach:** projection on consistent predictions is a regularization term in training loss
- **Fuzzy semantics:** fuzzy logic is employed as relaxation

Semantic Loss

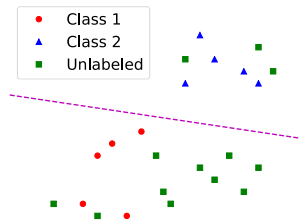
$$\mathcal{L}_s(\phi, \mathbf{p}) \propto -\log \sum_{\mathbf{y} \models \phi} \prod_{\mathbf{y} \models Y_i} p_i \prod_{\mathbf{y} \models \neg Y_i} (1 - p_i)$$

- ϕ is a propositional formula (a constraint that should hold)
- \mathbf{p} is a vector of probabilities associated to \mathbf{Y} variables (e.g. outputs of a neural network)
- The semantic loss is proportional to the negative logarithm of the probability that sampling \mathbf{Y} according to \mathbf{p} produces a value \mathbf{y} satisfying the constraint ϕ .

Semantic Loss Regularization



$$\phi = T$$



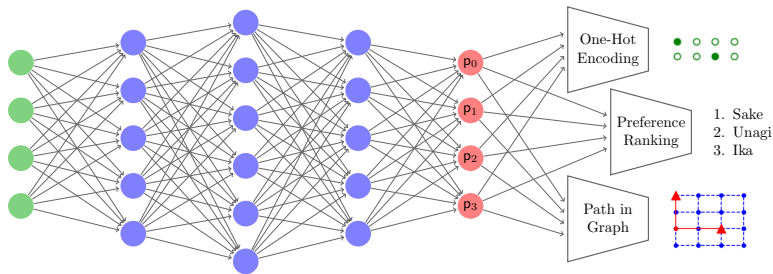
$$\phi = \neg Y_1 \vee \neg Y_2$$

Regularizing with semantic Loss

$$\mathcal{L}_{reg} = \text{training_loss} + \lambda \text{ semantic_loss}$$

- Semantic loss as regularizer of training loss (encourages predictions satisfying constraints)

Semantic Loss Regularization



End-to-end training with semantic Loss

- Semantic loss can be compiled into an algebraic circuit
- Partial derivatives can be computed on the circuit (see e.g. Deep ProbLog)

dimensions

- **Undirected model:** constraints as set of propositional formulas
- **Regularization approach:** semantic loss is additional term to training loss
- **Probabilistic semantics:** constraints are enforced in expectation over probabilities of possible worlds

```
nn(m_digit, [X], Y, [0, . . . , 9]) :: digit(X, Y).
```

From ProbLog to Deep ProbLog

- Introduce neural networks to process low-level data (softmax output layer)
- **neural annotated disjunction** (nAD) maps inputs to distributions over candidate outputs
- `nn` is a reserved word (stands for neural network)
- `m_digit` is the identifier of a neural network (CNN classifying digit images)
- `digit` is a **neural predicate** evaluated via `m_digit`.

Deep ProbLog: nAD example

`nn(m_digit, [X], Y, [0, ..., 9]) :: digit(X, Y).`

ground on **3**

`nn(m_digit, 3, 0) :: digit(3, 0) ; ... ; nn(m_digit, 3, 9) :: digit(3, 9).`

evaluate m_digit on **3**

`p0 :: digit(3, 0) ; ... ; p9 :: digit(3, 9).`

Inference by knowledge compilation

- 1 Ground relevant part of the program to answer query (including nADs).
- 2 Run forward step in neural nets to turn ground nAD into ground AD.
- 3 Compile resulting formula (same as ProbLog)
- 4 convert into AC (same as ProbLog)
- 5 evaluate AC (same as ProbLog)

Deep ProbLog: grounding example

```
nn(m_digit, [X], Y, [0...9]) :: digit(X,Y).  
addition(X,Y,Z) :- digit(X,N1), digit(Y,N2), Z is N1+N2.
```

DeepProbLog
program

ground on **0 1**

query
addition(**0**, **1**, 1)

```
nn(m_digit, [0], 0) :: digit(0, 0); nn(m_digit, [0], 1) :: digit(0, 1).  
nn(m_digit, [1], 0) :: digit(1, 0); nn(m_digit, [1], 1) :: digit(1, 1).  
addition(0, 1, 1) :- digit(0, 0), digit(1, 1).  
addition(0, 1, 1) :- digit(0, 1), digit(1, 0).
```

ground
DeepProbLog
program

forward step of nn

```
0.8 :: digit(0, 0); 0.1 :: digit(0, 1).  
0.2 :: digit(1, 0); 0.6 :: digit(1, 1).  
addition(0, 1, 1) :- digit(0, 0), digit(1, 1).  
addition(0, 1, 1) :- digit(0, 1), digit(1, 0).
```

ground
ProbLog
program

Image adapted from Manhaeve et al., 2019

Learning by gradient descent in ProbLog

- Gradient computation can be done over algebraic circuit used for inference.
- Need to replace probability semiring used for inference with **gradient semiring** (algebraic Problog)
- Gradient update followed by normalization to get valid probabilities

Deep ProbLog: probability vs gradient semiring

probability

gradient

$$a \oplus b = a + b$$

$$(a, \mathbf{a}_{\nabla}) \oplus (b, \mathbf{b}_{\nabla}) = (a + b, \mathbf{a}_{\nabla} + \mathbf{b}_{\nabla})$$

$$a \otimes b = ab$$

$$(a, \mathbf{a}_{\nabla}) \otimes (b, \mathbf{b}_{\nabla}) = (ab, a\mathbf{b}_{\nabla} + b\mathbf{a}_{\nabla})$$

$$e^{\oplus} = 0$$

$$e^{\oplus} = (0, \mathbf{0}_{\nabla})$$

$$e^{\otimes} = 1$$

$$e^{\otimes} = (1, \mathbf{0}_{\nabla})$$

$$L(f) = p$$

$$L(f) = (p, \mathbf{0}_{\nabla}) \quad (\text{fixed } p)$$

$$L(f_i) = p_i$$

$$L(f_i) = (p_i, \mathbf{e}_i) \quad (\text{learnable } p_i)$$

$$L(\neg f) = 1 - p$$

$$L(\neg f) = (1 - p, -\nabla p) \quad (\text{with } L(f) = (p, \nabla p))$$

ProbLog: gradient semiring example

```
0.2::earthquake.  
0.1::burglary.  
0.5::hears_alarm(mary).  
0.4::hears_alarm(john).  
alarm :- earthquake.  
alarm :- burglary.  
calls(X):-alarm,hears_alarm(X).
```

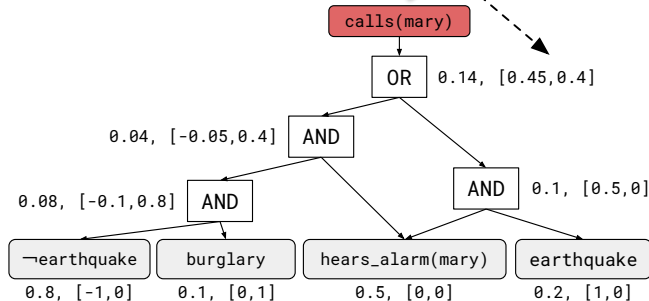
```
0.2::earthquake.  
0.1::burglary.  
0.5::hears_alarm(mary).  
alarm :- earthquake.  
alarm :- burglary.  
calls(mary):-alarm,hears_alarm(mary).
```

grounding for

`calls(mary)`

learnable parameters

compilation using gradient semiring



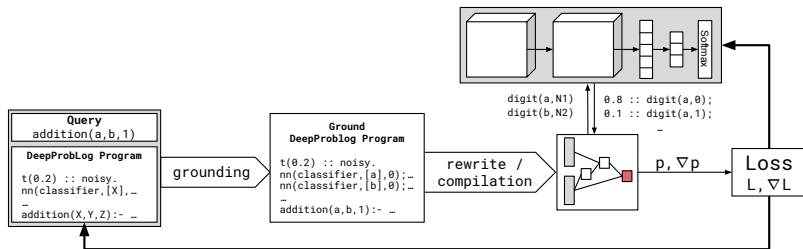
Learning by gradient descent in DeepProbLog

- Use gradient semiring as for ProbLog (considering outputs of neural predicates as abstract parameters).
- Backpropagate gradient from abstract parameters into the corresponding neural network

$$\frac{d\mathcal{L}}{d\theta_k} = \frac{d\mathcal{L}}{dP(q)} \sum_{i=1}^m \frac{dP(q)}{d\hat{p}_i} \frac{d\hat{p}_i}{d\theta_k}$$

- \mathcal{L} is a loss function
- $P(q)$ is the probability of a training example q (query)
- m is the number of outputs of a neural network (alternatives)
- \hat{p}_i is the i -th output of the network for example q .
- θ_k is the k -th parameter of a neural network

Deep ProbLog: learning pipeline



dimensions

- **Directed model:** probabilistic logic program (definite clauses)
- **Integration approach:** probabilistic logic program enriched with neural predicates
- **Probabilistic semantics:** constraints are enforced in expectation over probabilities of possible worlds

Motivation

- Theorem proving allows to infer novel facts entailed by a KB, but fails with noisy or ambiguous knowledge (e.g. slightly different names for the same relation)
- Neural models are robust to noise and ambiguity but have limited reasoning capabilities
- Neural theorem proving aims at combining the best of both worlds

In a nutshell

- End-to-end differentiable deductive reasoner
- Use Prolog backward-chaining algorithm for proving goals
- Replace symbolic unification between atoms with a differentiable similarity between their embeddings
- Collect the highest scoring proof as the goal proof
- Embeddings are learned by gradient descent over goal proofs for true (positive) and false (negative) facts.

Neural Theorem Proving: Prolog backward chaining

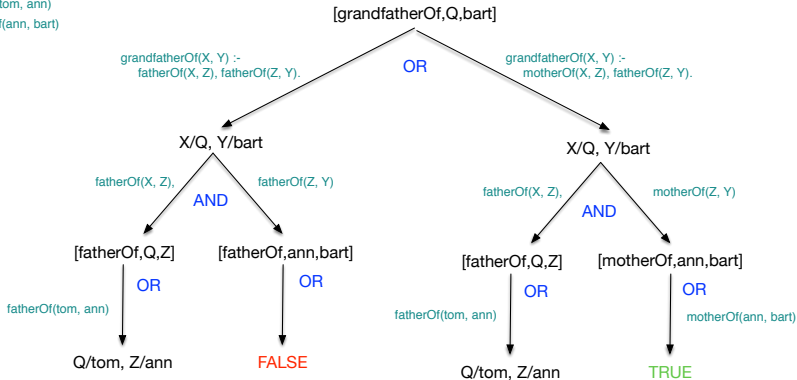
```
grandfatherOf(X, Y) :- fatherOf(X, Z), fatherOf(Z, Y).  
grandfatherOf(X, Y) :- fatherOf(X, Z), motherOf(Z, Y).  
fatherOf(tom, ann).  
motherOf(ann, bart).
```

OR / AND search

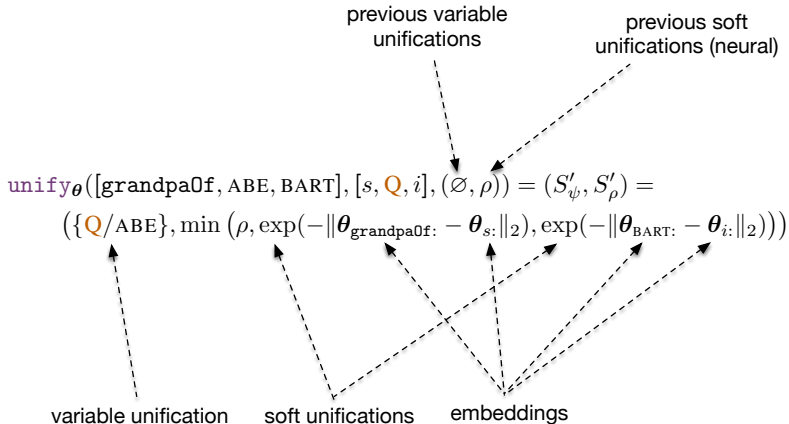
- OR iterates over all rules and unifies the rule head with the goal (one rule suffice)
- AND iterates over all atoms in the body of the rule (all atoms should be proved)
- OR is recursively applied to each atom in the body

Prolog backward chaining: example

grandfatherOf(X, Y) :- fatherOf(X, Z), fatherOf(Z, Y).
grandfatherOf(X, Y) :- fatherOf(X, Z), motherOf(Z, Y).
fatherOf(tom, ann)
motherOf(ann, bart)



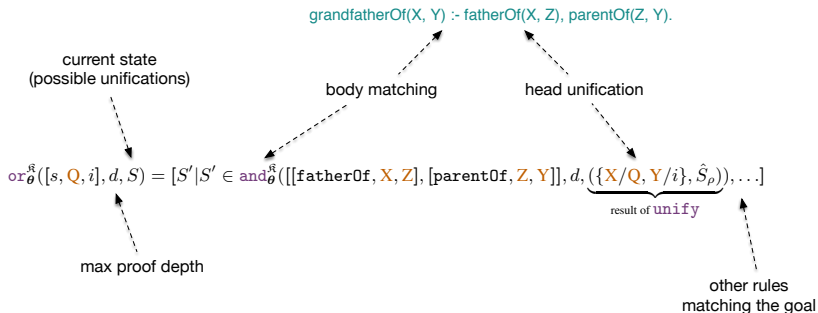
Neural Theorem Proving: unification



Soft unification

- Variables unify with variables or symbols as in Prolog
- Constants and predicates unify softly via similarity of their embeddings

Neural Theorem Proving: OR



OR module

- The goal is (soft) unified with the head of a rule (for all possible rules that soft unify)
- The AND module is called for all atoms in the body

Neural Theorem Proving: AND

$$\text{and}_{\hat{\theta}}^{\hat{\rho}}([\text{fatherOf}, X, Z], [\text{parentOf}, Z, Y], d, \underbrace{(\{X/Q, Y/i\}, \hat{S}_{\rho})}_{\text{result of unify in or}}) =$$
$$[S'' | S'' \in \text{and}_{\hat{\theta}}^{\hat{\rho}}([\text{parentOf}, Z, Y], d, S') \text{ for } S' \in \text{or}_{\hat{\theta}}^{\hat{\rho}}(\underbrace{[\text{fatherOf}, Q, Z]}_{\text{result of substitute}}, d-1, \underbrace{(\{X/Q, Y/i\}, \hat{S}_{\rho})}_{\text{result of unify in or}})]$$

AND called on remaining atoms

OR called on first atom

max depth is reduced

AND module

- The AND module fails if the maximum depth is reached (or the upstream OR failed)
- The AND module succeeds if it reaches the end of the list of atoms
- Otherwise it recurs over the atoms substituting variables wherever possible and calling OR

$$\text{ntp}_{\theta}^{\mathbb{R}}(\mathbf{G}, d) = \underset{\substack{S \in \text{or}_{\theta}^{\mathbb{R}}(\mathbf{G}, d, (\emptyset, 1)) \\ S \neq \text{FAIL}}}{S} \arg \max S_{\rho}$$

Proof with maximal score

- The search is initialized with an empty substitution set and a score of 1
- The maximization is over all possible goal proofs
- The score of a proof is the minimal score of all soft unifications in the proof

Neural Theorem Proving: proof example

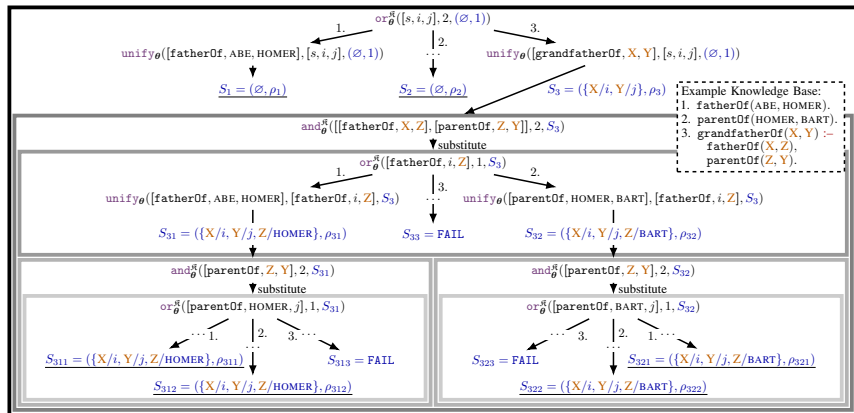


Image from Rocktäschel and Riedel, 2017

Neural Theorem Proving: prediction examples

QUERY: `part_of(CONGO.N.03,AFRICA.N.01)`

Score	Proofs
-------	--------

0.995	<code>part_of(X, Y) :- has_part(Y, X)</code> <code>has_part(AFRICA.N.01, CONGO.N.03)</code>
-------	--

0.787	<code>part_of(X, Y) :- instance_hyponym(Y, X)</code> <code>instance_hyponym(AFRICAN_COUNTRY.N.01, CONGO.N.03)</code>
-------	---

QUERY: `hyponym(EXTINGUISH.V.04, DECOUPLE.V.03)`

Score	Proofs
-------	--------

0.987	<code>hyponym(X, Y) :- hypernym(Y, X)</code> <code>hypernym(DECOUPLE.V.03, EXTINGUISH.V.04)</code>
-------	---

dimensions

- **Directed model:** logic program (definite clauses)
- **Integration approach:** logic program enriched with neural similarity in place of symbolic unification
- **“Fuzzy” semantics:** a score is associated to a proof, no explicit probabilistic interpretation

Bibliography

- Luc de Raedt, Sebastijan Dumancic, Robin Manhaeve, Giuseppe Marra, *From Statistical Relational to Neuro-Symbolic Artificial Intelligence*, In IJCAI 2020.
- Michelangelo Diligenti, Marco Gori, Claudio Saccà, *Semantic-based regularization for learning and inference*, Artificial Intelligence, Volume 244, 2017.
- Zhiting Hu, Xuezhe Ma, Zhengzhong Liu, Eduard Hovy, Eric Xing, *Harnessing Deep Neural Networks with Logic Rules*, Proc. of ACL, 2016.
- Xu, J., Zhang, Z., Friedman, T., Liang, Y. and Van den Broeck, G., *A Semantic Loss Function for Deep Learning with Symbolic Knowledge*. In ICML 2018.
- Manhaeve R., Dumancic S., Kimmig A., Demeester T. and De Raedt, Luc., *DeepProbLog: Neural Probabilistic Logic Programming*. In NeurIPS 2018.
- T. Rocktäschel and S. Riedel, *End-to-End Differentiable Proving*, Proc. of NIPS 2017.

Software Libraries

- **Semantic-based regularization (SBR)**
[<https://sites.google.com/site/semanticbasedregularization/home/software>]
- **Deep ProbLog** [<https://bitbucket.org/problog/deepproblog/src/master/>]
- **Greedy Neural Theorem Provers (GNTP)**
[<https://github.com/uclnlp/gntp>]