

# Fast Relational Learning using Neural-Symbolic Systems (CILP++)

Artur S. d'Avila Garcez

City, University of London  
[a.garcez@city.ac.uk](mailto:a.garcez@city.ac.uk)

ACAI, Ferrara  
30 August 2018

# Relational Learning

- Learning a *first-order logic* theory from examples (in the presence of uncertainty)
  - By searching for candidate hypotheses at first-order level or through *propositionalization*
  - Relevant for the analysis of complex networks: drug design in bioinformatics, link analysis in social networks, etc.
  - Related work: (probabilistic) ILP, deep networks, MLNs, BLOG, StarAI (lifted inference), etc.
- Hypothesis search (sound and sometimes complete at inducing theories) can be costly (but see streaming ILP (ILP'2013), online relational learning (ECML'12) and more recent developments (including differentiable (d)ILP, JAIR 2018).



# Relational Learning in Neural-Symbolic Systems

- A number of attempts at first-order (and higher-order) logic learning in neural nets: semantic approach (Hitzler et al), fibring, topos (Osnabruck), association (SHRUTI), unification, etc.
- Propositionalization offers a trade-off between information loss and efficiency; it seems a natural choice for use with neural nets.
- CILP++ is a neural-symbolic system that can solve ILP problems efficiently (through propositionalization) using a neural network trained with backpropagation (França, Zaverucha and d'Avila Garcez, Mach. Learn., July 2013)


# Efficient Relational Learning using CILP++

- CILP++ is composed of:
  - Bottom Clause Propositionalization (BCP): creates a bottom clause for each positive and negative ILP example
  - Artificial Neural Networks (ANNs) trained with backpropagation: generalising from a set of bottom clauses given as binary vectors (c.f. Muggleton and Tamaddoni-Nezhad, QG-GA, Mach. Learn. 2008)
  - Presenting the trained knowledge in relational form: mapping trained features into first-order logic representations (ICCSW 2013, Dagstuhl OASlcs, 2013)





## Summary of Results on Relational Knowledge Extraction

- Initial evaluation of FOL rules learned by CILP++ (BCP-rules), extracted from the neural network trained using BCP propositionalization:
    - We compared BCP-rules with rules produced by Aleph (we used RIPPER (Cohen, 1995) as the rule learner, but other propositional rule learners can be used)
    - As expected, there is information loss in comparison with Aleph, but in exchange for considerable speed-ups and smaller (more compact) rule sets
    - Although our extraction algorithm shows good fidelity to network, the lifting of FOL rules from trained neural nets can improve accuracy!
  - Experiments on FOL knowledge extraction from neural networks have been limited and are ongoing. Next step: use macro operators (R. Mooney) with CILP++
- 
- CITY U



# A simple example: Family Relationship

- BCP: generates a bottom clause for each (positive or negative) example; converts each bottom clause into a binary vector (similarly to QG-GA)

## Background Knowledge:

*mother(mom1, daughter1)*  
*wife(daughter1, husband1)*  
*wife(daughter2, husband2)*

## Positive Examples:

*motherInLaw(mom1, husband1)*

## Negative Examples:

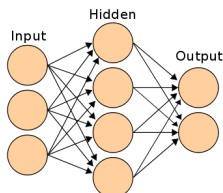
*motherInLaw(daughter1, husband2)*

- Using Prolog's bottom clause algorithm (Muggleton, 1995):
  - $\perp_+ = [\text{motherInLaw}(A, B) \leftarrow \text{mother}(A, C), \text{wife}(C, B)]$
  - $\perp_- = [\text{motherInLaw}(A, B) \leftarrow \text{wife}(A, C)]$
- BCP features: *mother(A, C), wife(C, B), wife(A, C)*
- Hypothesis: Is the use of a **set** of bottom clauses useful for learning and generalization?



# Neural networks and Backpropagation

- Popular connectionist models with application in many areas: pattern recognition, games, vision, speech, control
- CILP++ uses a multi-layer perceptron (MLP) with three layers of artificial neurons connected in a feed-forward manner; one could use a deep network instead, or even a recurrent network as used by CILP.
- Error back-propagation is a widely used training algorithm for MLPs; it seeks to minimize an error function through gradient descent
- When using error back-propagation, a common way of dealing with overfitting is to use **early stopping**

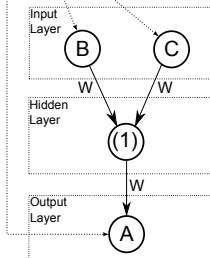




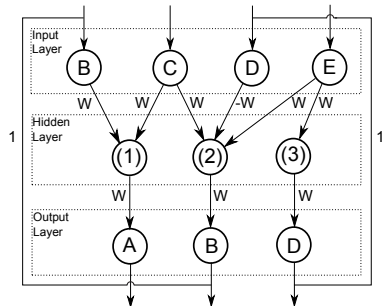
# The CILP Translation Algorithm (revisited)

$\text{BK} = \{A \leftarrow B, C; B \leftarrow C, \text{not } D, E; D \leftarrow E\}$

(1)                      (2)                      (3)



**N:**



# Bottom Clause Propositionalization (BCP)

- BCP uses Progol's bottom clause generation to create features for each ILP example
- Each positive example is saturated normally and labeled as positive (1); negative examples are processed just like positive ones, but are labeled as negative (-1)
- All unifications made during bottom clause generation are stored into *hash* sets
- A feature table  $F$  is created, consisting of all distinct body literals from all the bottom clauses created
- Finally, a binary input vector  $v$  of size  $|F|$  is created for each example (set of features) ( $\forall i, v(i) \in \{0, 1\}$ )



# BCP (cont.)

## Background Knowledge:

*mother(mom1, daughter1)*  
*wife(daughter1, husband1)*  
*wife(daughter2, husband2)*

## Positive Examples

*motherInLaw(mom1, husband1)*

## Negative Examples

*motherInLaw(daughter1, husband2)*

- Continuing the family relationship example:

$\perp_+ = [\text{motherInLaw}(A, B) \leftarrow \text{mother}(A, C), \text{wife}(C, B)]$

$\perp_- = [\text{motherInLaw}(A, B) \leftarrow \text{wife}(A, C)]$

- From  $\perp_+$ ,  $\text{hash}_+$ :

Key	Value
<i>mom1</i>	A
<i>husband1</i>	B
<i>daughter1</i>	C

- From  $\perp_-$ ,  $\text{hash}_-$ :

Key	Value
<i>daughter1</i>	A
<i>husband2</i>	B
<i>husband1</i>	C

- $F = \{\text{mother}(A, C), \text{wife}(C, B), \text{wife}(A, C)\}$
- $v_+ = (1, 1, 0)$
- $v_- = (0, 0, 1)$



- 4 Both standard stopping criteria and early stopping have been evaluated

# Stopping criteria

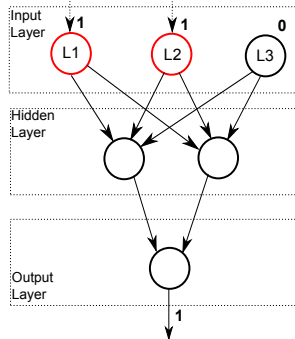
- Standard stopping criteria:
  - 300 training epochs have elapsed
  - 95% of the training data satisfies  $E < 0.1$
  - 90% of the training data is correctly classified by the network and no improvements can be seen for 5 training epochs
- Early stopping:
  - A validation set with  $\frac{1}{5}$  of the total set of examples was created
  - At time  $t$ , the best model from epochs 1 to  $t$  is stored
  - If after epoch  $e$ , the condition (first criterion of Prechet, 1999)
 
$$GL(t) > \alpha, GL(t) = 0.1 * \left( \frac{E_{va}(e)}{E_{opt}(t)} - 1 \right)$$
 is satisfied, training stops



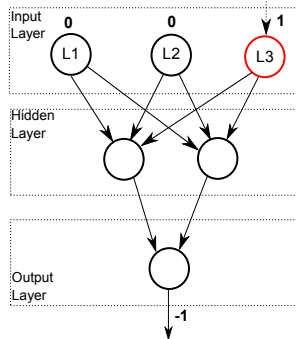
# Building the network and training

1)  $\text{motherInLaw}(A,B) \text{ :- } \underline{\text{mother}(A,C)}, \underline{\text{wife}(C,B)}$

L1:  $\text{mother}(A,C)$   
 L2:  $\text{wife}(C,B)$   
 L3:  $\text{wife}(A,C)$



2)  $\sim\text{motherInLaw}(A,B) \text{ :- } \underline{\text{wife}(A,C)}$



# Description of the Experiments

- Four CILP++ configurations have been tested:
  - *st*: uses standard backpropagation stopping criteria
  - *es*: uses early stopping
  - *n%bk*: the network is created using *n%* of the examples in ( $E_{\perp}^{train}$ ) as BK<sup>1</sup>
  - *2h*: when *n* = 0, the network uses 2 hidden neurons only!
- Why two hidden neurons in the *2h* configurations?
  - It can help avoid overfitting in large networks
  - Accuracy doesn't seem to increase substantially with more than 2 hidden neurons when the input is a binary vector (Haykin, 2009)

---

<sup>1</sup> *n* = 0, 2.5 and 5 were used.



# Experimental Results – Accuracy vs. Runtime

Dataset	<i>Aleph</i>	<i>CILP++<sub>st,2.5%bk</sub></i>	<i>CILP++<sub>st,5%bk</sub></i>	<i>CILP++<sub>st,2h</sub></i>
<i>mutagenesis</i>	80.85*( $\pm 10.5$ ) <b>0:08:15</b>	<b>91.70</b> ( $\pm 5.84$ ) 0:10:34	90.65( $\pm 8.97$ ) 0:11:15	89.20( $\pm 8.92$ ) 0:10:16
<i>krk</i>	<b>99.6</b> ( $\pm 0.51$ ) 0:11:03	98.31*( $\pm 1.23$ ) 0:04:38	98.32*( $\pm 1.25$ ) <b>0:04:34</b>	98.42( $\pm 1.26$ ) 0:04:40
<i>uw-cse</i>	<b>84.91</b> ( $\pm 7.32$ ) 0:45:47	66.24*( $\pm 7.01$ ) <b>0:08:47</b>	66.08*( $\pm 2.48$ ) 0:10:19	70.01*( $\pm 2.2$ ) 0:08:54
<i>alz-amine</i>	78.71( $\pm 5.25$ ) 1:31:05	<b>78.99</b> ( $\pm 4.46$ ) 1:23:42	76.02*( $\pm 3.79$ ) 2:07:04	77.08( $\pm 5.17$ ) <b>1:14:21</b>
<i>alz-acetyl</i>	<b>69.46</b> ( $\pm 3.6$ ) 8:06:06	63.64*( $\pm 4.01$ ) 4:20:28	63.49*( $\pm 4.16$ ) 5:49:51	63.30*( $\pm 5.09$ ) <b>2:47:52</b>
<i>alz-memory</i>	<b>68.57</b> ( $\pm 5.7$ ) 3:47:55	60.44*( $\pm 4.11$ ) 1:41:36	59.19*( $\pm 5.91$ ) 2:12:14	59.82*( $\pm 6.76$ ) <b>1:19:27</b>
<i>alz-toxic</i>	80.5( $\pm 3.98$ ) 6:02:05	79.92( $\pm 3.09$ ) 3:04:53	80.49( $\pm 3.65$ ) 3:33:17	<b>81.73</b> ( $\pm 4.68$ ) <b>2:12:17</b>

bold = best result; \* indicates statistically significant difference

- CILP++ achieves better accuracy and runtime in one *st* configuration
- CILP++{*st,2h*} is generally faster than Aleph



# Experimental Results with early stopping

Dataset	<i>Aleph</i>	<i>CILP++<sub>es,2.5%bk</sub></i>	<i>CILP++<sub>es,5%bk</sub></i>	<i>CILP++<sub>es,2h</sub></i>
<i>mutagenesis</i>	80.85(±10.51) 0:08:15	83.48(±7.68) <b>0:01:25</b>	83.01(±10.71) 0:01:43	<b>84.76</b> (±8.34) 0:01:50
<i>krk</i>	<b>99.6</b> (±0.51) 0:11:03	98.16*(±0.83) <b>0:04:08</b>	96.33*(±4.95) 0:04:28	98.31(±1.23) 0:04:18
<i>uw-cse</i>	<b>84.91</b> (±7.32) 0:45:47	68.16*(±4.77) <b>0:04:08</b>	65.69*(±1.81) 0:04:16	67.86*(±1.79) <b>0:04:08</b>
<i>alz-amine</i>	<b>78.71</b> (±3.51) 1:31:05	65.33*(±9.32) 0:35:27	65.44*(±5.58) <b>0:08:30</b>	70.26*(±7.1) 0:10:14
<i>alz-acetyl</i>	<b>69.46</b> (±3.6) 8:06:06	64.97*(±5.81) 3:04:47	64.88*(±4.64) 2:42:31	65.47*(±2.43) <b>0:25:43</b>
<i>alz-memory</i>	<b>68.57</b> (±5.7) 3:47:55	53.43*(±5.64) 1:40:51	54.84*(±6.01) 3:57:39	51.57*(±5.36) <b>1:33:35</b>
<i>alz-toxic</i>	<b>80.5</b> (±4.83) 6:02:05	67.55*(±6.36) <b>0:12:33</b>	67.26*(±7.5) 0:14:04	74.48*(±5.62) 0:28:39

- Considerable speed-ups are obtained, in exchange for accuracy at times



# Comparison with another Propositionalization

Dataset	<i>Aleph</i>	<i>BCP+ANN</i>	<i>RSD+ANN</i>	<i>BCP+C4.5</i>	<i>RSD+C4.5</i>
<i>muta</i>	80.85 * ( $\pm 10.51$ ) 0:08:15	<b>89.20</b> ( $\pm 8.92$ ) 0:10:16	67.63 * ( $\pm 16.44$ ) 0:11:11	85.43 * ( $\pm 11.85$ ) <b>0:02:01</b>	87.77 ( $\pm 1.02$ ) 0:02:29
<i>krk</i>	<b>99.6</b> ( $\pm 0.51$ ) 0:11:03	98.42 * ( $\pm 1.26$ ) 0:04:40	72.38 * ( $\pm 12.94$ ) 0:06:21	98.84 * ( $\pm 0.77$ ) <b>0:01:59</b>	96.1 * ( $\pm 0.11$ ) 0:05:54

- BCP achieves better accuracy than RSD overall, while being slightly faster
- Accuracy of BCP with ANNs is much higher than RSD with ANNs



# Relational Knowledge Extraction

- Since BCP features (and neurons) are first-order literals, it should be possible to extend the CILP knowledge extraction algorithm directly to obtain first-order rules from trained CILP++ networks
- *...but such rules will not obey any ILP restrictions as imposed by a language bias*
- França et. al., ICCSW'13, AAAI Spring Symposium 2015, contain a first proposal towards an efficient relational knowledge extraction algorithm using the propositional rule learner RIPPER (other rule learner could be used)



# From CILP++ networks to FOL rules

- Search for literals whose variables do not obey the ILP-style variable chaining
- Using BCP's **hash tables**, replace such variables by a disjunction of ground terms
- This transformation is provably correct: **BCP feature equivalence theorem**



# Example: Family relationship revisited

$$\begin{aligned}
 S_{\perp} = \{ & \text{motherInLaw}(A, B) : - \text{mother}(C, B), \text{wife}(C, D); \\
 & \text{motherInLaw}(A, B) : - \text{mother}(A, C), \text{wife}(C, B); \\
 & \sim \text{motherInLaw}(A, B) : - \text{wife}(C, B), \text{parents}(C, B, D), \text{dad}(E, F) \}.
 \end{aligned}$$

$$R_{\perp} = \{ \text{motherInLaw}(A, B) : - \text{mother}(A, C), \text{wife}(C, D), \text{not}(\text{dad}(E, F)) \}.$$

$$\begin{aligned}
 \text{hash} = \{ & D/\text{husband2}, D/\text{daughter1}, E/\text{husband1}, \\
 & F/\text{daughter1} \}
 \end{aligned}$$

- From the BCP feature equivalence:

$$\text{not}(\text{dad}(E, F)) \mapsto \text{not}(\text{dad}(\text{husband1}, \text{daughter1}))$$

$$\text{wife}(C, D) \mapsto \text{wife}(C, \text{daughter1}) \vee \text{wife}(C, \text{husband2})$$

- The resulting first-order rule is:

$$\begin{aligned}
 R_{\perp}^C = \{ & \text{motherInLaw}(A, B) : -\text{mother}(A, C), (\text{wife}(C, \text{daughter1}); \text{wife}(C, \text{husband2})), \\
 & \text{not}(\text{dad}(\text{husband1}, \text{daughter1})) \}.
 \end{aligned}$$



# First-order filtering

- We can improve results further by making the theory more compact
- A modified version of the theory filtering algorithm *T-reduce* (A. Srinivasan, The Aleph System, version 5) is applied
  - **Original T-reduce**: removes sets of rules without positive coverage or that contribute negatively to training set accuracy
  - **Modified T-reduce**: also removes literals that are always true, and literals containing no variables
- If applied on  $R_{\perp}^C$ :

$$R_{\perp}^{FOL} = \{motherInLaw(A, B) :- mother(A, C), (wife(C, daughter1); wife(C, husband2))\}$$



# Experimental settings

- We have evaluated the approach (named BCP+RIP<sub>FOL</sub>) in comparison with:
  - Aleph (as a baseline)
  - BCP + RIPPER (named BCP+RIP<sub>prop</sub>)
- We use the Alzheimers benchmark for comparison
- The used metrics for evaluation are:
  - Classification accuracy
  - Runtime
  - Theory size (i.e. total number of literals)



# Initial Results

	<i>Alz-ami</i>	<i>Alz-ace</i>	<i>Alz-mem</i>	<i>Alz-tox</i>
<i>Aleph</i> (baseline)	78.71(±5.25) 1:31:05, 36.1	69.46(±4.6) 8:06:06, 47.3	68.57(±5.7) 3:47:55, 45.7	80.5(±4.83) 6:02:05, 37.9
<i>BCP+RIP<sub>prop</sub></i>	73.35(±4.32) 0:19:49, 30	67.8(±3.77) 0:23:21, 20.1	65.27(±7.11) 0:25:11, 14.4	78.44(±5.44) 0:17:41, 35.2
<i>BCP+RIP<sub>FOL</sub></i>	77.73(±4.57) 0:21:59, 30.4	63.56(±5.06) 0:26:39, 18.7	57.64(±5.7) 0:28:45, 13.8	66.45(±6.93) 0:20:57, 18

Accuracy and theory size averaged over 10-fold cross-validation

- As expected, Aleph's accuracy is superior
- But we produce first-order theories that are more compact
- And the accuracy of BCP+RIP<sub>FOL</sub> is even higher than BCP+RIP in one case



# In Summary

- CILP++: Competitive accuracy results can be obtained efficiently using a neural-symbolic approach; in principle, any ILP dataset can be used with CILP++
- BCP can be used with any propositional learner (although some don't seem to like the idea of bottom clauses very much)
- Relational Knowledge Extraction: capable of generating compact first-order rule sets with negation from trained neural nets
- Considerably faster than Aleph, but sometimes with (sometimes considerable) accuracy loss



## Next Steps

- Evaluation of noise robustness and evaluation in theory revision tasks (with negation)
- Apply to first-order knowledge extraction from neural networks in general (compare with TREPAN and its variations)
- Compare accuracy results with other fast and hybrid ILP approaches, including differentiable ILP
- A lot of the magic seems to happen in the choice (and number of copies) of first-order literals to place in the input and output layers: use macro-operators!?