Lifted Relational Neural Networks

Gustav Sourek, Vojtech Aschenbrenner, Filip Zelezny

& Ondrej Kuzelka

Outline

- Motivation
 - From Neural Nets point of view (possibly)
 - From Markov Logic point of view
- What are Lifted Relational Neural Networks
 - Short version
 - Long version (possibly)
- Learning latent concepts with LRNNs

LRNN Motivation from Neural Networks' POV

Motivation (NN POV)

- How to learn with relational or graph-structured data?
 - Examples : molecules (networks, trees, etc.)
- How to represent data samples?
 - Sets of vertices & edges, relational logic clauses
 - Isomorphic samples should be treated the same!
- How to feed them into a classifier, a neural network?

Propositionalization

- Idea : turn arbitrary graph into a fixed-size vector
 - Through a predefined aggregation mapping



Powerful, yet need to predefine all useful patterns

Auxiliary concepts

- There may be useful sub-structures present
- For instance, halogen groups in a molecule (mutagenicity) classification problem
 - e.g., C-Br, C-Cl, C-F may be indicative
 - i.e., there is a useful pattern C-(halogen atom)
- We can predefine these in the feature-vector

Latent predicate invention

- What if we do not know any of the useful sub-structures of the problem in advance?
 - e.g., we do not know there is something like halogens or other indicative group of atoms
- \rightarrow We may design anonymous predicates for these patterns
 - And learn these in a way such that they are useful in different contexts (rules) (Muggleton, 1988)
- \rightarrow Neural learning of latent (non ground) patterns
 - This is beyond the scope of propositionalization

LRNNS

- We propose a framework avoiding the aforementioned limitation of propositionalization
- Lifted Relation Neural Networks (LRNNs)
- Inspiration:
 - Lifted (templated) graphical models: Markov Logic Networks (Richardson, Domingos, 2005), Bayesian Logic Programs (Kersting, De Raedt, 2000)
 - Neural-symbolic approaches: KBANN(Towel, Shavlik, 1994), CILP(Franca, Zaverucha, Garcez, 1999)

LRNN Motivation from Markov Logic POV

Markov Logic POV)

 How to learn with relational or graph-structured data in the presence of uncertainty?

> Lifted graphical models, e.g. Markov Logic

- How to *efficiently* learn latent concepts?
 - Neural Networks (propositional concepts)
- > How about latent relational concept learning?

Lifted Relational Neural Networks

What is LRNN? short version

What is LRNN? (short version)

- Syntactically: Set of weighted first-order Horn clauses
 - 0.5 : water :- bondOH(X,Y)
 - 1.0: bondOH(X,Y) :- H(X), O(Y), bond(X,Y)
 - LRNN encoding looks familiar like a weighted Prolog program...
- Semantically: Template for neural network construction
 - > We turn the template's Herbrand models into NNs as follows..

Network Construction

- 1. Every ground proposition (atom) which can be *derived* * from a given LRNN model corresponds to an **atom neuron**
- 2. Every **ground rule** $h \leftarrow (b_1, ..., b_k)$ such that $(b_1, ..., b_k)$ can be *derived* * from a given LRNN corresponds to a **rule neuron**
- 3. To aggregate different groundings derived with the same rule's ground head {h \leftarrow (b¹₁, ..., b¹_k), ..., h \leftarrow (bⁿ₁, ..., bⁿ_k)} there is an **aggregation neuron**

* meaning it is present in the least Herbrand model

Putting it all together...



Weight Learning

- LRNN model := grounding of {sample, template} clauses
 - Different samples result in different ground networks
 - This induces weight sharing across ground networks as their neurons are *tied* to the same template rules
- Different aggregation functions are used as neurons' activations so as to reflect the (fuzzy) logic of disjunction, conjunction, and different forms of aggregative reasoning over relational patterns
- Stochastic Gradient Descend can be used for training

What is LRNN? Long version

Data representation

- No propositionalization or feature vector transformation
- Similarly to LRNNs, we represent samples simply as raw sets of corresponding facts (typically ground unit clauses)



 A simple set union {} of a LRNN template with a relational sample can thus be though of simply as another LRNN

LRNN construction

LRNN := union of a sample and template clauses

→ Different samples result in different LRNNs

- Template remains the same
- We introduce building blocks of LRNN construction, these are 3 different types of neurons : atom neurons, rule neurons, aggregation neurons

Atom Neurons

- Every ground proposition (atom) which can be *derived** from a given LRNN corresponds to an atom neuron
- Example LRNN:
 - Template : 1.0 : bondOH(X,Y) :- H(X), O(Y), bond(X,Y).
 - Sample : H(h1), H(h2), O(o1), bond(h1,o1), bond(h2,o1)
- \rightarrow Set of all atom neurons:



• $\{N_{H(h1)}, N_{H(h2)}, N_{O(o1)}, N_{bond(h1,o1)}, N_{bond(h2,o1)}, N_{bond(h1,o1)}, N_{bondOH(h2,o1)}\}$

(* Meaning present in the least Herbrand model of it)

Atom Neurons

- Every ground proposition (atom) which can be *derived** from a given LRNN corresponds to an <u>atom neuron</u>
- Example LRNN:
 - Template : 1.0 : bondOH(X,Y) :- H(X), O(Y), bond(X,Y).
 - Sample : H(h1), H(h2), O(o1), bond(h1,o1), bond(h2,o1)
- \rightarrow Set of all atom neurons:



Rule neurons

- Every ground rule h ← (b₁, ..., b_k) such that (b₁, ..., b_k) can be *derived** from a given LRNN corresponds to a rule neuron
- Example LRNN:
 - Template : 1.0 : bondOH(X,Y) :- H(X), O(Y), bond(X,Y)
 - Sample : H(h1), H(h2), O(o1), bond(h1,o1), bond(h2,o1)
- \rightarrow Set of all rule neurons:

 $N_{bondOH(h1,o1)} \leftarrow H(h1), O(o1), bond(h1,o1), N_{bondOH(h2,o1)} \leftarrow H(h2), O(o1), bond(h2,o1)$ (*Meaning the atoms are true in the least Herbrand model)

Rule neurons

- Every ground rule h ← (b₁, ..., b_k) such that (b₁, ..., b_k) can be *derived** from a given LRNN corresponds to a rule neuron
- Example LRNN:
 - Template : 1.0 : bondOH(X,Y) :- H(X), O(Y), bond(X,Y)
 - Sample : H(h1), H(h2), O(o1), bond(h1,o1), bond(h2,o1)
- \rightarrow Set of all rule neurons:



Rule neuron activation

- Rule neuron basically represents conjunctive If-Then rule
 - This should be reflected in its activation function
- → Rule neuron has high output if and only if all the input atom neurons (rule's body) have high outputs
- Fuzzy logic inspiration :



Aggregation neurons

- We need to aggregate different groundings of the same non-ground rule having the same ground literal in the head. For each such aggregation there is an aggregation neuron.
- Example LRNN:
 - Template : 1.0 : hasOH :- bondOH(X,Y)
 1.0 : bondOH(X,Y) :- H(X), O(Y), bond(X,Y)
 - Sample : H(h1), H(h2), O(o1), bond(h1,o1), bond(h2,o1)
- Set of different ground rules for hasOH :- bondOH(X,Y) corresponds to neurons:
 - $N_{hasOH \leftarrow bondOH(h1,o1)}$, $N_{hasOH \leftarrow bondOH(h2,o1)}$
- Aggregation neuron $N_{hasOH \leftarrow bondOH(X,Y)}$ aggregates over these

Aggregation functions

- Different aggregation functions might be used for different logic of aggregation neurons
- MAX corresponds to "best pattern" matching



Possibilities in other contexts include, e.g., AVG

Atom neuron inputs

- There may be multiple weighted rules with the same ground head, yet with different weights
- Example template:
 - 1.0 : Group1 :- hasOH
 - 0.2 : Group1 :- hasHCl
- I.e. we end up with two different aggregation neurons with different weights:
 - 1.0 : $N_{Group1:-hasOH}$ and 0.2 : $N_{Group1:-hasHCI}$
- These finally form the inputs of atom neuron N_{Group1}

Atom neuron activation

- Combining different rules implying the same atom naturally corresponds to disjunction
- Atom neuron output should be high if and only if at least one of the rule neurons has high output
- Fuzzy logic inspiration :



Putting it all together...



Weight Learning

- The constructed ground LRNN can be thought of as a regular neural network with shared weights
- The shared weights come from grounding of same template's clause and exploit sample regularities
 - Similarly to convolutional neural networks this does not pose any problem to weight learning
- Stochastic Gradient Descend (SGD) with mild adaptions can be efficiently used for training

Experiments

Experiment template

```
0.0 atomGroup1(X) :- o(X).
0.0 atomGroup1(X) :- cl(X).
```

```
0.0 atomGroup3(X) :- cl(X).
```

```
0.0 bondGroup3(X) :- 2=(X).
```

```
graphlet0 :- atomGroup2(X), bond(X,Y,B1),
bondGroup1(B1), atomGroup3(Y)...
```

```
0.0 class1 :- graphlet0.
```

```
0.0 class1 :- graphlet242.
```

Samples



Results



Where was latent predicate Invention?

- Different modeling concepts exploiting predicate invention
- Particularly, implicit soft clustering:

atomKappa_1	4.58	2.62 1	11.23	0.39	3.39	0.43	0.68	0.76	-0.70	-0.30	-2.36	-15.35	1.59 -	-0.52	-1.11	2.43-	10.26	1.12 1	14.44	-0.17	-9.22	4.95 ·	-1.66	1.64	4.09 -	-0.22	2.08	-3.82	-0.37	0.69	-21.83	10.22	12.23 1	10.91	6.98-6	.19
atomKappa_2	1.16	0.13	0.55	0.85	0.10	0.08	0.12	0.30	0.10	-0.28	0.02 -	-0.50	0.35 -	-0.03	0.22	0.31 ·	-0.69	0.22	1.96	0.03 -	-0.38	1.59	0.01	0.17	0.28	0.10	0.08 -	-1.44	0.65	0.14	-1.46	-0.17	0.68	0.98 -	0.40-0	A 3
atomKappa_3	4.89	2.56 1	10.04	0.46	2.12	0.71	0.87	0.24	-0.91	-3.84	-2.15	14.09	1.45 -	-0.30-	-0.68	2.04	-9.40	1.30 1	15.14	-0.70	-8.74	6.50	-0.75	0.58	2.43	0.17	2.24	-3.36	-0.03	0.59	-19.47	4.95	15.321	0.38	6.33-5	.65
	c_28	c_14	c_16	n_31	0-40	195	0_51	c_25	n_36	c_21	d_93	n_34	0_52	0_41	br_94	c_194	n_35	0_49	c_29	0_45	n_32	c_27	f_92	c_19	n_38	h_8	c_232	c_22	0_42	c_230	h_3	c_26	c_10	c_195	4	2000

 Other concepts include soft-matching, hypergraph approximation, relational autoencoders,...

Learning Predictive Categories Using Lifted Relational Neural Networks

Gustav Sourek¹, Suresh Manandhar², Filip Zelezny¹, Steven Schockaert³, and Ondrej Kuzelka³

> Czech Technical University in Prague, Czech Republic {souregus, zelezny}@fel.cvut.cz
> Department of Computer Science, University of York, UK suresh.manandhar@york.ac.uk
> School of CS & Informatics, Cardiff University, UK {SchockaertS1, KuzelkaO}@cardiff.ac.uk

Learning Predictive Categories with LRNNs

Learning Predictive Categories

We consider a following (learning) scenario with latent categories:

- 1. Entities
 - a) Have properties, b) Belong to **categories**
 - > Categories largely determine belonging entities' properties
- 2. Properties
 - a) Belong to entities, b) Belong to categories
 - > **Categories** largely determine entities satisfying a property

Learning Predictive Categories

- 1. Given : a set of entities and corresponding lists of their properties
- 2. Assumption : there exists some latent hierarchy of categories that are predictive of their corresponding object's properties
 - The hierarchy should allow for property inheritance
 - Similarly we induce latent hierarchy on properties
- 3. Goal: Learn suitable category structures from data

Encoding in LRNN

- Given input samples : {1/0 HasProperty(e, p)
- Membership to categories : w_{ec} : IsA(e, c)
- Category hierarchy : W_{c1c2} : IsA(c₁, c₂)
- Category properties : *w_{cecp}* : HasProperty(c_e, c_p)
- Transitivity : w_{isa} : IsA(A, C) \leftarrow IsA(A, B),IsA(B, C)
- Categories determine their entities' properties : *w*'_{cecp} : HasProperty(A, B) ← IsA(A, c_e), IsA(B, c_p), HasProperty(c_e, c_p)

Learning Setting

- We minimize MSE of the query atom neuron outputs and their targets {1/0 HasProperty(e, p)} via SGD
- The activation functions used were
 - Conjunction $\Lambda_{(b1,\ldots,bk)} = \text{sigm} \left(\sum_{i=1}^{k} b_i k + b_0 \right)$
 - Disjunction $V_{(b1,\ldots,bk)} = \text{sigm} \left(\sum_{i=1}^{k} b_i + b_0 \right)$
 - Aggregation $*_{(b1, \ldots, bm)} = \max_{i} b_{i}$
- We set up 2 level hierarchy with [3, 2] hidden categories for both objects and properties

Evaluation

- Animals dataset (https://alchemy.cs.washington.edu/data/animals)
 - 50 animals + 65 properties (e.g., large, smelly, strong,...)

- Predictive ability : AUC PR 0.8, AUC ROC 0.86
- Same as with second order Markov Logic Networks, reported in (Statistical Predicate Invention, Kok and Domingos, 2007)
 - Which is related to the introduced model, while jointly clustering objects and relations

Embeddings of entities



Embeddings of properties



Outlook

- We have learned implicit similarity measure via latent category membership degrees
- - Where I denotes some level of similarity, e.g. based on externally obtained embeddings
 - With that we might emulate 1-NN or kernel regression
- Also whole triples of (subject, *predicate*, object) might be considered to learn soft categories of predicates, too

Conclusions

- LRNNs are a flexible framework to easily encode non-trivial SRL scenarios
 - e.g., a joint learning of predictive categories of entities and their properties
- More complicated settings might be easily reached with just mild extensions of the template
 - e.g., semi-supervised learning, embeddings, etc.
- We plan for thorough comparison with MLNs and incorporation of LRNNs into NLP tasks pipelines

General conclusions

- LRNNs may be thought of as a neural analogy to lifted (templated) graphical models (e.g., Markov Logic Networks)
 - Both are template languages for defining and weight tying in the corresponding ground models
- Benefits
 - Latent (deep) relational concepts, Flexible templates (e.g., convolutional NN), Explicit variable binding
 - Future work

 \bullet

 Different modeling concepts, recurrent NNs, ASP, optimization, Structure learning inspired by meta-interpretive learning

Thank you!

See "Lifted Relational Neural Networks" at arXiv.org for more details

Convolutional NN







image I as a set of weighted facts