

---

# Learning Concepts through Differentiable Logical Induction

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We develop a framework where an algorithm learns concepts that have symbolic  
2 and subsymbolic content. The generative algorithm can learn the logical structure  
3 underlying a set of observed data while simultaneously operating on dense vector  
4 representations of concepts. We revisit the framework studied by previous work that  
5 models semantic cognition as a form of logical dimensionality reduction inducing  
6 a set of abstract rules, while simultaneously benefiting from the recent  
7 substantial progress of artificial neural networks. Our algorithm inherits some  
8 of the advantages of both perspectives by combining the simplicity of forward  
9 chaining with the power of parametrized unification. We show that these Neural  
10 Logical Reasoners can perform inductive and deductive inferences while handling  
11 ambiguity and noise, at much better speed than previous symbolic approaches.

## 12 Preamble. About Concepts

13 There is much to say about concepts. Crucially, the meaning of a concept depends in part of their  
14 relationship with other concepts and in part of their relation with the world. Recent theories from  
15 cognitive science research posit two aspects for concepts and their determination. 'Narrow' content  
16 is related to conceptual role, this part of the meaning of a concept comes from how it relates to other  
17 concepts. The Theory Theory of concepts (TODOREFERENCE) posits that you cannot learn and  
18 represent concepts in isolation, instead concepts acquire their meaning arises in the context of a  
19 whole conceptual system. Thus it has been argued that children don't learn concepts like ENERGY,  
20 MOMENTUM and ACCELERATION in terms of concepts they already know. Rather, they gradually  
21 learn how to use the whole new terminology. It is worth noting this is very aligned with predicate  
22 invention in the context of Inductive Logic Programming (ILP), in those cases the meaning of a  
23 predicate like EVEN comes from how it relates to other predicates.

24 On the other hand, 'wide' content is related to the referents picked up out there in the world by the  
25 representations. This Dual-Factor Theory of concepts is influenced by Kripke and Putnam's subtle  
26 analysis suggesting that references cannot be entirely determined by internal relations of concepts  
27 with other concepts. That is, what we know about entities picked out by a concept cannot be entirely  
28 what determines which entities those are, that is because what we know about entities is always  
29 subject to revision and to be mistaken. Instead of conceptual role semantics, wide content is mediated  
30 by informational semantics which is more related to statistical covariation between representations  
31 and reality. To not get into a philosophical digression, we refer the reader to relevant extensive  
32 literature(TODOREFERENCE). Instead we just mention that this duality of factors can perhaps  
33 help accommodate a wide range of psychological data from the 1970's showing that people can rank  
34 concepts in terms of their typicality and let them to propose that concepts are better characterized  
35 as lists of properties of features. These rankings are robust and have a direct effect on the speed of  
36 categorization(TODOREFERENCE) and also provides one candidate explanation of why concepts  
37 are rather fuzzy and inexact.

Most representations in artificial intelligence don't have rich dual-factor components. Symbolic approaches seem intuitively more related to conceptual role, often explicitly specifying the relationships between variables. In contrast, subsymbolic approaches compute representations with little conceptual role, at the very least, the "meaning" of the representations is obtained through statistical covariation rather than logical composition. The two approaches seem strikingly complementary in their strengths and weaknesses. Just like in the philosophical debates around concepts, both approaches have managed to appeal some of the smartest people in the field, often opposing them. Just like with the history of the theory of concepts, perhaps this suggests a combination.

## 1 Introduction

Until recently, most research on the problem of learning logical rules such as  $grandfatherof(x, y) \leftarrow fatherof(x, z), parentof(z, y)$  from data was in the context of symbolic systems like Prolog(REFERENCE). These algorithms have a number of desirable properties. First, they are interpretable, as they learn explicit relations between the concepts. Second, they tend to be very data efficient, able to generalize well from a handful of examples. While Neural Networks tend to require a lot of training examples and struggle to generalise, ILP approaches are doomed to generalise by learning general rules involving free variables that apply for all concepts. Third, as (EVANS) argue, these models easily support transfer learning, being explicit, learned rules can be transferred to other tasks to continue further training. Relatedly, these models present a straightforward way of incorporating domain-specific knowledge in the form of explicit logical rules. In contrast, traditional ILP models are unable to handle noisy, erroneous or ambiguous data. A recent paper (REFERENCE) proposed a differentiable version of ILP. This work combined some of the advantages of ILP and of neural network-based systems which allowed them to use gradient descent and backpropagation to learn rules that were interpretable and data-efficient while robust to noisy and ambiguous data. While our work resembles in several senses this work, it also differs in several dimensions. Fundamentally, while their approach, like previous ILP solutions, requires to first generate all the possible rules and then learns to select the relevant rules with a mechanism akin to attention; our approach operates at the more modular and compositional level of individual concepts, and directly learns the atoms that conform the logical rules.

In a less explored domain, previous symbolic approaches have also been embedded in hierarchical Bayesian models that are capable of performing a kind of dimensionality reduction for structured logical theories(REFERENCE). These probabilistic generative models learn theories by inducing a set of logical rules along with a set of core relations that form a compression of the data which can be recovered using the logical rules. As an example of how theories support compression consider an animal taxonomy like that of Figure 1. The algorithm can learn to compress all the information about salmons into the core relation  $IS(salmon, fish)$  which is sufficient to recover and infer all sorts of things about salmons ( $IS(salmon, animal)$ ,  $HAS(salmon, fins)$ , etc.) using other core relations and a set of learned logical rules ( $IS(x, y) \leftarrow IS(x, z), IS(z, y)$ ;  $HAS(x, y) \leftarrow IS(x, z), HAS(z, y)$ ). Thus through the induction of logical rules, the algorithm can learn to make deductive inferences. The additional induction of core relations allows the very interesting capability to the argument of making inductive inference. When observing that salmons have fins, and gills and are animals, these probabilistic algorithms can be incentivised to compress the information into the single core relation that salmons are fishes. This capability of making very rich inductive and deductive inferences from very sparse data and some general abstract knowledge is a landmark of human learning. While the Bayesian symbolic presented a very promising direction with interesting results it showed limitations in terms of speed and scalability to real datasets. It also suggests challenges in terms of connecting to lower level perception modules.

In this work we explore to what extent the substantial progress made in the last year(s) in terms of modelling logical reasoning with neural approaches can capture the inductive and deductive inferences shown in previous work. We build on (SEBASTIAN'S REFERENCE) and work with a new version of neural logical reasoning that combines the simplicity of forward chaining with the flexibility of parametrized unification by operating in dense vector representations. This allows to apply rules when the symbols of the atoms are not equal but similar in meaning and thus replace symbolic comparison with a graded notion of similarity. This approach can also connect seamlessly to upstream perception units and scale better due to the efficacy of SGD and backpropagation. We show that

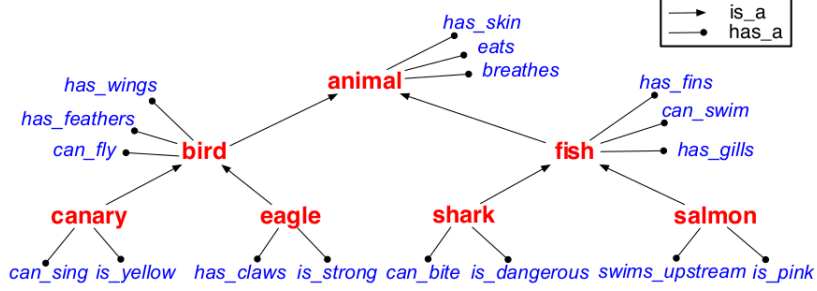


Figure 1: Animal Taxonomy. Constants in red and blue, relations indicated with lines and arrows.

neural logical reasoners can perform induction and deduction at much better speed than previous symbolic approaches.

The structure of the paper is as follows. We begin with an overview of some relevant background. Covering ILP and other recent neural models. Then in section 3 we propose our algorithm that performs a differentiable form of ILP where dense vector representations are learned through back-propagation while forming the building blocks of logical rules. In section 4 we discuss three different case studies that highlight different capabilities of our approach. We conclude in section 5 with a discussion about the limitations and possible future directions of the proposed framework.

## 2 Background

### 2.1 Inductive Logic Programming

A Logic Program is a set of logical if-then rules and background facts which can be used to answer queries. A rule in these systems is typically of the form:

$$h \leftarrow b_1, b_2, \dots, b_k \quad (1)$$

Here  $h$  is called the *head* of the rule and  $b_1, b_2, \dots, b_k$  constitute the *body*. Intuitively the head of the rule is true if each of the  $b_i$  in the body are. For example a rule might be:

$$\text{grandfather}(X, Y) \leftarrow \text{father}(X, Z), \text{parent}(Z, Y) \quad (2)$$

which is read:  $X$  is the grandfather of  $Y$  if  $X$  is the father of  $Z$  AND  $Z$  is the parent of  $Y$ . Here  $X$  and  $Y$  are universally quantified and  $Z$  is existentially quantified. This means that the rule holds for any  $X$  and  $Y$  as long there is some individual  $Z$  for which  $\text{father}(X, Z)$  and  $\text{parent}(Z, Y)$  are both true facts. Given background facts such as  $\text{father}(\text{Bill}, \text{Mary})$  and  $\text{parent}(\text{Mary}, \text{Liz})$ , a logic programming system can use the rule to prove a goal fact  $\text{grandfather}(\text{Bill}, \text{Liz})$ . This is an example of a *forward chaining* deduction because it starts from a set of facts and matches (unifies) the body of a rule to derive the goal. It is also possible to do *backward chaining* in which we start with a goal and work backwards by unifying it with the head of a rule, recursively trying to prove the body.

But specifying all the rules is cumbersome and brittle and we would like instead to be able to learn the rules from examples. Learning the rules necessary for a logic program from a set of background facts and positive and negative examples of the goal is called *inductive logic programming* (ILP) system. In practice, the necessary rules are chosen from a human-supplied template or meta-rule which narrows the space of possibilities. ILP has been extensively developed over the last three decades for symbolic systems but only recently recast as a continuous optimization problem amenable to solution using neural networks and stochastic gradient descent [FIX: citations].

Symbolic ILP systems do very well at generalizing from just a few examples. This is because they are learning universal rules. They are however susceptible to noisy inputs and even a single bad fact can cause them to fail. On the other hand, neural systems generally are very robust to noisy input but sample inefficient and prone to overfitting on small amounts of data. Differentiable ILP systems aim for the best of both worlds. They can be made robust to noisy inputs while still retaining some of the strong generalization properties typically associated with symbolic systems. However, current systems suffer from poor memory scalability since they must create ground versions of their rules.

## 128 2.2 Differentiable ILP systems

129 There are two differentiable ILP approaches of which we are aware [FIX: cite, evans and rock]. Both  
130 of these approaches assume a rule template which describes the structure of any candidate rule to be  
131 induced. Both construct a differentiable function which implements a proof of the desired goal. Both  
132 require grounding the constructed rules on all the constants – effectively creating a family of ground  
133 rules to evaluate. And both offer interpretable rules after training. However they differ in several key  
134 respects.

135 First, [cite: rock] constructs a function representing a backward-chained proof of the goal, while  
136 [cite: evans] do a forward chained proof of the goal from the initial facts. [cite: evans] requires a  
137 representation of the truth values of all possible facts and non-facts, which may require considerable  
138 space. By contrast, [cite: rock] only requires a representation of the true facts. A more conceptual  
139 distinction arises in their parameterizations. In [cite: evans] the parameters are weights on the set  
140 of possible choices for each atom in the body of the rule – the rule structure. On the other hand in  
141 [cite:rock], the rule structure is, in effect, fixed and what is learned instead is the embeddings of the  
142 goal predicate and arguments.

143 In our approach we follow [cite:rock] in parameterizing with embeddings but use the forward rather  
144 than backward chaining approach so that we don't have to represent a proof tree explicitly. This  
145 greatly improves memory scalability since we do not need to represent all possible groundings of a  
146 set of symbolic rules.

## 147 2.3 Semantic Cognition

148 Joshs models for importance of:

- 149 - Deductive and inductive inferences.
- 150 - Core relations in theories

## 151 2.4 Neural Networks for Knowledge

152 -Sebastians paper

153 Ours better... not rule enforcement

## 154 2.5 Other Related Work

155 Go for a walk? Already in limitations.

156 This could be a substitute for the usual section of related work

## 157 3 Neural Logical Reasoner

158 In this section we describe the Neural Logical Reasoner. A model that induces logical rules through  
159 the learning of the vector embeddings that constitute the atoms. The algorithm is conceptually very  
160 simple: It starts with a set of known facts that constitutes its current Background Knowledge. For  
161 each forward step, the algorithm generates all the consequences implied by its known logical rules.  
162 Implication is done through the unifications of the predicates of the rules and those of the known  
163 facts. The generated consequences after a fixed number of forward steps are compared with a set of  
164 positive and negative examples (Figure 2).

### 165 3.1 Inference

166 The setup is similar to that in (SEBASTIANS). 1) Facts consists of triplets of the form Rela-  
167 tion(Constant1, Constant2) (i.e Mother(Rosa, Andres)). 2) Concepts are associated with vector  
168 embeddings through a dictionary. For now only relations will have associated embeddings. Later  
169 we will consider the case where constants can also have embeddings. 3) Logical rules with a pre-  
170 specified structural form are initialized as randomly parameterized sets of vectors , for example:  
171  $[Head](X, Y) \leftarrow [Body_1](Y, Z), [Body_2](Z, Y)$  or  $[Head](X, Y) \leftarrow [Body_1](Y, X)$ . The content

of these embedding representations is learned from the data. 4) Finally, fact are associated with a valuation  $v \in [0, 1]$  which represents the belief of that fact being true.

**Forward Chaining** A forward step consists on taking a set of known facts and generating all the consequences implied by the rules. We implement two alternative methods.

1) When the problem is sufficiently simple and there are only a few concepts we use Method1. Method1 keeps a valuation for every potential fact (i.e  $r_1(s, o)$ ) and updates its valuation according to:

$$v_{r_1(s,o)}^{i+1} = \max_{r_1 \leftarrow r_2, r_3} \left( \max_Z \min(v_{r_2(s,Z)}^i, v_{r_3(Z,o)}^i) \right)$$

That is, maximising over all rules with the same head, maximizing over the free variable  $z$ , and minimising over the atoms in the rule (this implements logical conjunction, minimization can also be used).

2) Method2 considers only the set of known facts. It implements the same update, but instead of iterating through the space of all facts to find the pair that maximally implies it, it iterates through all the pairs of known facts and keeps the top  $K$  maximally implied facts.

**Parametrized Unification** To know if the relation of a particular fact is related to the atom of a particular rule. The unification score for the relevant embeddings  $U(r, f)$  is computed using a similarity metric (cosine\_similarity was selected after some exploration). With the unification scores, the equation above for a particular rule  $head \leftarrow body_1, body_2$  now becomes:

$$v_{head(s,o)}^{i+1} = U(body_1, fact_1) * U(body_2, fact_2) * \left( \max_Z \min(v_{fact_1(s,Z)}^i, v_{fact_2(Z,o)}^i) \right)$$

Implicitly adding the requirement that for an implication to happen the unification of the rule and the considered facts has to be high.

## 3.2 Learning

As mentioned, logical rules are initialized as randomly parameterized sets of vectors. And the content of these embedding representations is learned from the data with backpropagation through all the forward chaining steps.

Vector and Conceptual role influence each other

- **Loss**
  - **One-Hot Vectors**
  - **General Embeddings**
- Note that intelligent sampling helps with that problem. Like humans

## 3.3 Learning Theories

- **Learning Background Knowledge** Same mechanism of induction. But for core facts
- **Constant Embeddings**

## 4 Case Studies

- **Space of tasks**
- A wide range of previous work has focused on different aspects of logical induction and knowledge base completion. Here we consider three different case studies to highlight the range of capabilities of our algorithm. We consider three case studies where our algorithm is better in some sense than previous considered approaches: Forward chaining simple Josh., fast Grefenstette compositional, no memory, cleaner When compared against such a large range it has limitations, discussed further down

### 4.1 Standard ILP Tasks

- **Evans and Grefenstette Ambiguity. Noise?**

## 206 4.2 Inferring Core Relations

- 207 • **Josh's**
- 208 • **Constant Embeddings?** The space of constants
- 209 • **Sparse Inference**

## 210 4.3 Completing a Knowledge Base

211 better: no enforcement + no ComplEX+ Simplicity

- 212 • **Countries and other Sebastian's tasks** Contrast with Josh's models, ours is more scalable
- 213 than those

## 214 5 Conclusion and Future work

215 We see our work as a proposition of a research direction that combines the simplicity of logical  
216 forward chaining with parametrized unification using dense embedding representations. Such a  
217 combination allows for concepts that acquire their meaning simultaneously from both their relations  
218 to other concepts through logical rules like in ILP, and from their subsymbolic embeddings like in  
219 Neural Networks. We suggest how such representations consists of a move towards representations  
220 with parallelisms to a contemporary notion of a "Concept" which has been the subject of much study  
221 in cognitive science and philosophy. The algorithm we proposed inherits many of the advantages of  
222 both approaches. From the symbolic side, it acquire interpretable representations that are learned  
223 with little data and show great generalization. At the same side, its subsymbolic nature allows it to  
224 handle ambiguity, noise and graded similarity, while being able to be learn through gradient descent,  
225 which makes it more scalable to bigger real datasets. We highlight this different advantages in three case  
226 studies. First we evaluate it in a set of traditional ILP tasks, replicating results from a recent seminal  
227 paper. Second, we show that the structure of the algorithm allows it to be deployable in realistic  
228 datasets where previous ILP approaches fail. Third and perhaps most importantly we show that neural  
229 logical reasoners can perform induction and deduction from sparse data, through the additional  
230 induction of core relations, with much better speed than previous bayesian symbolic approaches.

231 **Limitations** Our algorithm tries to address a broad range of aspects in several dimensions and  
232 presents limitations in all of them. From an algorithmic perspective, like all the considered previous  
233 work, our algorithm is provided with templates that contain information of the structural form of  
234 the rules. This would ideally be part of the learning algorithm. From a scalability perspective.  
235 While greatly improving in terms of memory, speed and size of the data relative to previous models,  
236 the model remains far from being able to reach bigger tasks of knowledge completion, that have  
237 been attacked with other purely neural approaches (REFERENCE GO FOR A WALK ???Do we  
238 want this). From a cognitive science perspective, the model is still more limited than its bayesian  
239 symbolic counterparts. Specifically, while those models provide graded measures of confidence  
240 in their inferences, the neural logical reasoners do not currently provide meaningful estimates of  
241 uncertainty, but see below.

242 **Future Directions** We signal some concrete (and straightforward in some cases) ways of addressing  
243 the above limitations. We also point to some clear directions to enrich our current framework that  
244 we would like to explore in future work. First, a straightforward way of having the algorithm learn  
245 the templates would be to encode the structural information of the atoms in the rules (arity and  
246 variable order) by adding dimensions to the embeddings and have the algorithm use independent  
247 unifications that it would then interpret in the desired ways. This would constitute only a slightly  
248 more complicated learning task but would maintain the same structure and mechanism of the problem  
249 that could be trained through gradient descent. Second, more interesting sampling procedures and the  
250 integration of forward with backward chaining could perhaps yield regimes more similar to those that  
251 humans yield with that could help cope with scalability to larger datasets. Third, we would like to  
252 investigate different ways of providing better estimates of uncertainty: from a full neural probabilistic  
253 formulation, to a heuristic metric based on the number of initializations and on the unification scores.  
254 Finally, an interesting direction to be explored relies on the insight that when a generative model  
255 has a very particular form, in this case the forward application of logical rules, the latent variables

256 are forced to acquire a very particular form, in this case a set of core relations. We are starting to  
257 investigate the idea of Logical Autoencoders that because of their particular logical decoders, are  
258 forced to create powerful encoders that can perform things like inductive inference.

259 The current framework constitutes an attempt of a step in the direction of building the representations  
260 that let humans learn so much from so little.

## 261 **Acknowledgments**

262 Gracias.

## 263 **References**

264 Test of ? [1] Alexander, J.A. & Mozer, M.C. (1995) Template-based algorithms for connectionist  
265 rule extraction. In G. Tesauro, D.S. Touretzky and T.K. Leen (eds.), *Advances in Neural Information*  
266 *Processing Systems 7*, pp. 609–616. Cambridge, MA: MIT Press.

267 [2] Bower, J.M. & Beeman, D. (1995) *The Book of GENESIS: Exploring Realistic Neural Models*  
268 *with the GEneral NEural Simulation System*. New York: TELOS/Springer-Verlag.

269 [3] Hasselmo, M.E., Schnell, E. & Barkai, E. (1995) Dynamics of learning and recall at excitatory re-  
270 current synapses and cholinergic modulation in rat hippocampal region CA3. *Journal of Neuroscience*  
271 **15**(7):5249-5262.