



City Research Online

City, University of London Institutional Repository

Citation: Harder, F. & Besold, T. R. (2018). Learning Lukasiewicz logic. *Cognitive Systems Research*, 47, pp. 42-67. doi: 10.1016/j.cogsys.2017.07.004

This is the accepted version of the paper.

This version of the publication may differ from the final published version.

Permanent repository link: <https://openaccess.city.ac.uk/id/eprint/19865/>

Link to published version: <https://doi.org/10.1016/j.cogsys.2017.07.004>

Copyright: City Research Online aims to make research outputs of City, University of London available to a wider audience. Copyright and Moral Rights remain with the author(s) and/or copyright holders. URLs from City Research Online may be freely distributed and linked to.

Reuse: Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge. Provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.

Learning Łukasiewicz Logic

Frederik Harder^a, Tarek R. Besold^b

^a*Faculty of Science, University of Amsterdam, Amsterdam, The Netherlands*

^b*Digital Media Lab, Center for Computing and Communication Technologies (TZI),
University of Bremen, Bremen, Germany*

Abstract

The integration between connectionist learning and logic-based reasoning is a long-standing foundational question in artificial intelligence, cognitive systems, and computer science in general. Research into neural-symbolic integration aims to tackle this challenge, developing approaches bridging the gap between sub-symbolic and symbolic representation and computation. In this line of work the *core method* has been suggested as a way of translating logic programs into a multilayer perceptron computing least models of the programs. In particular, a variant of the core method for three valued Łukasiewicz logic has proven to be applicable to cognitive modelling among others in the context of Byrne’s suppression task. Building on the underlying formal results and the corresponding computational framework, the present article provides a modified core method suitable for the supervised learning of Łukasiewicz logic (and of a closely-related variant thereof), implements and executes the corresponding supervised learning with the backpropagation algorithm and, finally, constructs a rule extraction method in order to close the neural-symbolic cycle. The resulting system is then evaluated in several empirical test cases, and recommendations for future developments are derived.

Keywords: Neural networks, Logic programs, Neural-symbolic integration, Cognitive modelling, Reasoning

1. Introduction

Neural-symbolic integration attempts to bridge the gap between two prominent paradigms in artificial intelligence. Symbolic AI, the first of the two, encompasses ex-

5 plicit knowledge representation, logic programming and search-based problem solving
techniques which have been responsible for many of the early successes in artificial
intelligence such as game playing, automated theorem proving and natural language
processing ([1, 2, 3]). While the paradigm is still very much alive in expert systems
managing and reasoning over vast quantities of symbolic data, it is also at times referred
to as “good old-fashioned AI” or GOF AI ([4]), having lost some of its appeal as its lim-
10 itations have become apparent. Learning from, and finding structure in sets of noisy
data is something symbolic AI largely fails at. Unfortunately this means that whole
classes of problems which are integral to a common conception of intelligence, such
as image and voice recognition, on a general scale currently can hardly be addressed
using symbolic AI.¹ Also, while (mostly non-monotonic) logic-based cognitive mod-
15 elling is still being pursued with valuable results, the brittleness of the corresponding
models together with their necessary restriction to high-level cognition (leaving out the
bigger part of the actual representation and processing apparatus of human cognizers),
are clear drawbacks when compared to connectionist or statistical approaches.

The second paradigm is that of machine learning. As the name suggests, it refers
20 to a variety of methods for learning from data, artificial neural networks (ANN) be-
ing one prominent group of these methods. Aided by a leap in processing power and
available data, machine learning has been credited with most of the more recent ac-
complishments in AI, from the now commonplace feat of handwriting recognition to
self-driving cars and the fully autonomous learning of computer games ([6, 7, 8]).
25 Promising as the paradigm may be, there are areas in which, on its own, it performs
very poorly. While the learning of simple logical dependencies from data is achieved
with relative ease, the process becomes increasingly difficult when higher order con-
cepts are involved ([9]). Examples for the latter impasse are numerous, including con-
nectionist systems’ problems with high-level visual analysis taking into account partial

¹Recent logic-based approaches such as, for instance, Meta-Interpretive Learning for Logical Vision [5] might in the future help to mitigate this problem, but currently have only reached proof of concept state and still have to confirm their generalizability across tasks and domains, and their scalability to real-world problem sizes.

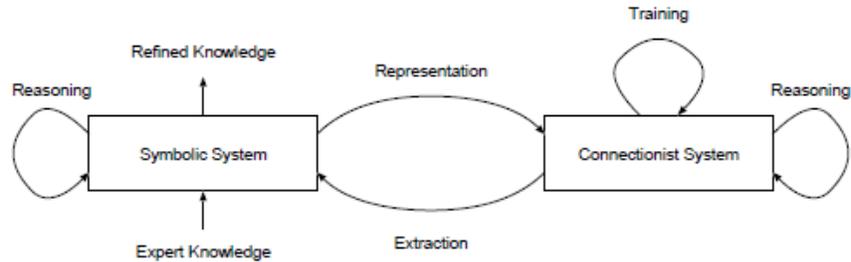


Figure 1: A conceptual overview of the neural-symbolic cycle (as introduced in [11]).

30 occlusion, light source identification, or shadow prediction, or with higher-level inference such as the recognition of intentions of depicted agents. Also, as knowledge is represented in connectionist systems in a distributed fashion that is hard to interpret from an outside perspective, it is usually difficult to provide background knowledge in a format which the machine learning algorithm can use, or to extract learned features
 35 from a network for instance for verification purposes. All of these are problems that often become trivial when tackled with a symbolic system.

Much stands to be gained from a unification of the two paradigms that could cancel out their respective weak spots and highlight their strengths. Neural-symbolic integration ([10]) offers some ideas in how this may be achieved, centering around the
 40 concept of the neural-symbolic cycle (see Figure 1). The cycle contains two reasoning systems. One is symbol-based, utilizing available expert knowledge, and the other is a connectionist system or ANN, which learns from data. The challenge of interfacing these systems is twofold. Coming from the symbolic side, the first task is to find a way of translating the existing symbolic knowledge into the connectionist system, finding a
 45 representation that is appropriate for the network. Secondly, one needs to devise methods for extracting the information gained by the connectionist system through learning and convert it back into a clean symbolic format. Equipped with these processes of representation and extraction the system as a whole is capable of incorporating both background knowledge and training data as either become available.

50 When asked about the feasibility of integrating both paradigms, the human brain and mind serve as prime examples and proof of concept. The brain has a neural struc-

ture which operates on the basis of low-level processing of perceptual signals, but cognition also exhibits the capability to efficiently perform abstract reasoning and symbol processing; in fact, processes of the latter type are taken to provide the foundations
55 for thinking, decision-making, and other mental activities ([12]). It is precisely this seamless coupling between learning and reasoning which is commonly considered the basis for intelligence in humans—see also, e.g., [13], p. 163: “While I do not regard intelligence as a unitary phenomenon, I do believe that the problem of reasoning from learned data is a central aspect of it.”—and, in close analogy, quite plausibly also for the
60 (re-)creation of cognitive capacities up to human-level intelligence in artificial systems.

Returning to the neural-symbolic cycle discussed above, it should be made clear, that the task of constructing such a cycle rapidly increases in difficulty when raising the expressive capacities of the involved systems. There are approaches for fragments of first order logic ([14, 15]), but most results focus on various propositional logics.
65 Furthermore, extraction algorithms for connectionist systems tend to be intractable. So while the general method of the field can be described in a few pages, the underlying problems are hard and there is still a long way to go before neural-symbolic integration may be applied to state-of-the-art methods of either paradigm.

As one of the currently most prominent and best understood methods, Hölldobler’s
70 and Kalinke’s *core method* ([16]) has since been developed as a neural-symbolic system for, among others, propositional modal ([17]) and covered first order logic programs ([15]). It provides a way of translating logic programs into a type of multilayer perceptron (MLP) which, embedded in the core architecture, computes least models of these programs. In [18], a variant of the core method for three valued Łukasiewicz logic is
75 presented, and it is suggested to apply the resulting approach to cognitive modelling tasks (see, e.g., [19] for a subsequent application to Byrne’s suppression task [20]). In the discussion of their work, the authors make the claim that the architecture they have used can be modified in such a way, as to allow training via the backpropagation algorithm ([21]). If this is in fact the case, when additionally equipped with a rule ex-
80 traction method, the resulting architecture should allow for a basic form of closure of the neural symbolic cycle.

Expanding upon work started by Harder and Besold in [22], in the following we put

these ideas into practice by providing a modified core suitable for supervised learning, implementing and executing supervised learning with the backpropagation algorithm and, finally, constructing a rule extraction method. Section 2 gives an overview of the theory and methods underlying this work, after which Section 3 is used for an in-depth documentation of the implemented approach to learning cores and extracting learned rules. We present the empirical results of the corresponding computational experiments in Section 4, followed by a closing discussion and look ahead at future work in Section 5. The proofs corresponding to the theoretical results, together with the pseudo codes of the extraction algorithm, have been relegated to the appendix.

2. Foundations

As conceptual basis for the work presented in subsequent sections, a number of methods and terminology have to be clarified. The three following subsections will respectively give a short introduction to Łukasiewicz logic programs, Hölldobler's and Kalinke's core method, and the backpropagation algorithm. Some basic familiarity with classical logic and neural networks is assumed.

2.1. Łukasiewicz logic programs

We first introduce three-valued Łukasiewicz logic as a formalism, giving the main definitions and a short account of key properties. Then we provide the relevant information about logic programs and weak completion in the Łukasiewicz context, before finally presenting the Stenning-van-Lambalgen consequence operator.

2.1.1. Three-valued Łukasiewicz logic

Łukasiewicz Logic was proposed in its ternary version in 1917 by Polish philosopher and logician Jan Łukasiewicz ([23]), as a result of his work on modalities in logic. The addition of a third value was meant to introduce a notion of possibility and indeterminism to logical reasoning. Metaphysical import aside, this logic was the first one to break with the true/false dichotomy of classical logic and thus lay the conceptual basis for the development of further many-valued and fuzzy logics thereafter. The

110 connective definitions for three-valued Łukasiewicz logic are provided in Figure 2 below (following the notation used in [18]). A noteworthy property, when compared, for instance, to Kleene’s strong logic of indeterminacy ([24]) is the definition of $u \rightarrow u$ and $u \leftrightarrow u$ as true. This allows for the existence of tautologies, i.e. formulas which are true under all interpretations, and preserves some of those familiar from classical
115 logic. Conventionally, an interpretation assigns one of the three values \top , \perp and u to each atom in the Universe U . In the given context it makes sense to define an interpretation as a tuple $I = \langle I^\top, I^\perp \rangle$, where I^\top is a set containing all atoms assigned the value \top and I^\perp contains all atoms assigned \perp . No atom is in both sets and those assigned u are in neither set. One can speak of an *empty* interpretation when $I^\top \cup I^\perp = \emptyset$ and a
120 *partial* interpretation when $I^\top \cup I^\perp \subsetneq U$. I is a *model* of a formula G , iff $I(G) = \top$.

\wedge	\top	u	\perp	\vee	\top	u	\perp	\neg
\top	\top	u	\perp	\top	\top	\top	\top	\top
u	u	u	\perp	u	\top	u	u	u
\perp	\perp	\perp	\perp	\perp	\top	u	\perp	\perp

\rightarrow	\top	u	\perp	\leftrightarrow	\top	u	\perp
\top	\top	u	\perp	\top	\top	u	\perp
u	\top	\top	u	u	u	\top	u
\perp	\top	\top	\top	\perp	\perp	u	\top

Figure 2: Definitions of the connectives for three-valued Łukasiewicz logic.

As already mentioned in the previous section, according to Hölldobler, Łukasiewicz logic is of interest for modelling empirical observations of human reasoning. Specifically, [19] seeks to provide a logical model for the suppression task experiment ([20]). In the corresponding set of experiments, Byrne analysed what conclusions readers will
125 typically draw from a certain class of natural language statements. As an example, when reading the statement “If Marian has an essay to write, she will study late in the library. She has an essay to write.”, 96% of all subjects concluded that “Marian will study late in the library.”. When presented with the same item, but with the added information that “If the library stays open, she will study late in the library.”, only

130 38% of participants conclude that Marian will indeed study in the library. It appears that the additional information led more than half of the participants to revoke their previous inferences, even though this information was not contradictory. The non-monotonicity of this reasoning suggests that it cannot be modelled by classical logic. Against this backdrop, in [19] it is therefore proposed that three-valued Łukasiewicz
135 logic interpreted under weak completion (as explained in the following subsection) fits the findings best.²

2.1.2. Logic programs

Logic programs are defined as a finite set of clauses of the form $A \leftarrow B_1 \wedge B_2 \wedge \dots \wedge B_n$ where the head of the clause, A , is an atom and the B_i , with $1 \leq i \leq n$, in the body
140 are either literals, \top or \perp . Clauses of the form $A \leftarrow \top$ and $A \leftarrow \perp$ are called *positive* and *negative facts* respectively.

These logic programs are interpreted under weak completion, which takes a logic program and transforms it into one single formula, thereby changing how it is evaluated. Firstly, the bodies of all clauses with the same head are concatenated as a disjunction
145 into one body. After this step, the resulting formulas consist of one implication per head. Subsequently, all \leftarrow are replaced with \leftrightarrow . As a result, atoms which are heads in clauses whose bodies all evaluate as \perp are now \perp as well. Finally, concatenating all clauses into one conjunction creates a single formula representing the weakly completed program.

150 Weak completion adds non-monotonicity to Łukasiewicz logic. Atoms which evaluate as \perp because all associated bodies evaluated as \perp , can become \top when adding another clause without contradiction. Also, weakly completed Łukasiewicz logic programs are never contradictory, always having at least one model ([18]).

²There is an ongoing controversy on whether Łukasiewicz logic under weak completion, or completion semantics based on the three-valued logic used by Fitting [25], is better suitable for modelling human reasoning in general, and the suppression task in particular. While this debate and its eventual solution are of general interest, in this article we stay agnostic concerning the matter. Instead, as stated at the end of Section 1, the aim is to equip the system from [18] with a backpropagation-trainable type of core and a rule extraction mechanism, closing the neural-symbolic cycle.

2.1.3. Consequence operators

155 Models for such a logic program P can be computed through a consequence operator Φ_P . Starting from an empty interpretation I , the immediate consequence $\Phi_P(I)$ is calculated as a new interpretation and this process is iterated, until $I = \Phi_P(I)$ and a fixed point is reached. It can be shown ([18]) that the least models of Łukasiewicz logic under weak completion are identical to the least fixed points of the Stenning-van-
160 Lambalgen consequence operator $\Phi_{SvL,P}$ from [26], which is defined as follows:

$$\begin{aligned} \Phi_{SvL,P}(I) &= \langle J^\top, J^\perp \rangle, \text{ where} \\ J^\top &= \{A \mid \exists (A \leftarrow \text{body}) \in P : I(\text{body}) = \top\} \text{ and} \\ J^\perp &= \{A \mid \exists (A \leftarrow \text{body}) \in P \wedge \forall (A \leftarrow \text{body}) \in P : I(\text{body}) = \perp\} \end{aligned}$$

The next section discusses the algorithm introduced by [18], which translates the
165 $\Phi_{SvL,P}$ consequence operator of a given program into a 3-layer feed-forward network, that computes the same function. Like the consequence operator, this network may then be used on multiple iterations until a fixed point is reached.

2.1.4. Example: consequence operator application

In order to illustrate, how the consequence operator $\Phi_{SvL,P}$ functions, we provide a
170 small example. Consider the following four clauses:

$$\begin{aligned} A &\leftarrow \perp \\ B &\leftarrow \neg A \\ C &\leftarrow \neg B \\ C &\leftarrow B \wedge \neg C \end{aligned}$$

All literals default to u in the beginning. The first application of the consequence operator adds A to J_1^\top following the first clause. The bodies of clauses 2,3 and 4 evaluate as u and nothing else happens. The second application, $\Phi_{SvL,P}(I_1)$, finds clause two satisfied and maps B to \top . On the third iteration, the body of clause 3 evaluates

175 as \perp , but the body of clause 4 is still u , and as a result, a fixed point is reached. The sequence of interpretations is given in the table below.

	I_0	I_1	I_2	I_3
A	u	\perp	\perp	\perp
B	u	u	\top	\top
C	u	u	u	u

2.2. The core method

In the following, a detailed description of how the core architecture is set up will be provided. This account chooses a somewhat different perspective than that of the translation algorithm given in [18]. While the algorithmic description is optimal for
180 implementation, the angle used here will hopefully provide a better understanding of the core structure with regard to the modifications that must be made to it, and to the introduction of sigmoidal activation units in particular.

In both input and output layer of the network, each atom A of the program is represented by two neurons. Activation in the first one indicates $A = \top$, while activation in
185 the second one means $A = \perp$. If neither neuron is active, then $A = u$. The core does not allow for both neurons to be active in the same iteration. The input layer also contains one neuron each, representing \top and \perp , which are always active. Each program clause, or rather each clause body, is represented by two neurons $\langle h^\top, h^\perp \rangle$ in the hidden layer.
190 Whether a clauses body is mapped to \top , \perp or u is encoded in the same way as was used for the atoms.

All connections between the layers of the core serve the function of logic gates. An h^\top neuron is connected to one input layer neuron for each conjunct in the clause body it represents. If the conjunct is an atom A , it connects to A^\top , if it is a negated atom $\neg A$,
195 it connects to A^\perp and if the conjunct is \top , it connects to that unit. Connection weights and activation threshold are set to form an 'and'-gate, requiring activation of all input layer neurons for the clause neuron to fire. In case a conjunct is \perp , no connection is formed, but for sake of the logic gate this is treated as a connection to an inactive unit, preventing the clause neuron from ever firing. Respectively, an h^\perp neuron connects to

200 A^\perp neurons, where A is a conjunct and A^\top neurons, where $\neg A$ is one. If \perp is a conjunct, h^\perp connects to the \perp neuron and if \top is a conjunct, no connection is formed. Weights and threshold are set to form an 'or'-gate, such that h^\perp is activated when one or more input neurons fire. This way, clause bodies are represented as \top if and only if all their conjuncts are mapped to \top and represented as \perp , if and only if one or more conjuncts
205 are \perp .

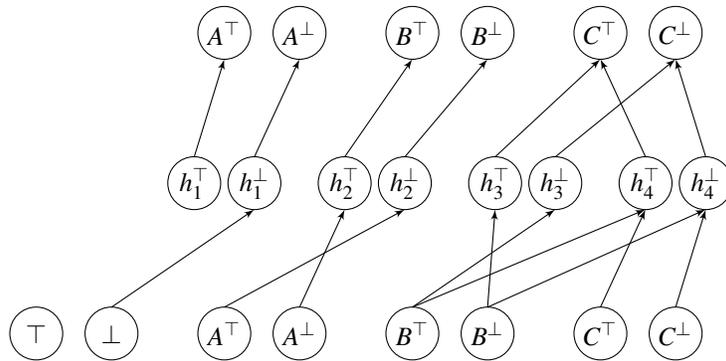
In the output layer, each neuron has one connection for each clause in which the associated atom appears as head. A^\top neurons are connected to the h^\top neurons of the associated clauses, forming an 'or'-gate and A^\perp neurons are connected to the h^\perp neurons, forming an 'and'-gate. Thus atoms are \top when one or more associated clauses
210 are \top , and \perp , when all associated clauses are \perp .

The logic gates are implemented such that all connection weights in the network have the same value $\omega > 0$ and 'or'-gate thresholds are at $0.5 \cdot \omega$, while 'and'-gate thresholds equal to $(l - 0.5) \cdot \omega$, where l is the number of incoming connections. All neurons use the Heaviside activation function, emitting 1, if the received activation
215 meets or exceeds the threshold and 0 otherwise. Given this setup, computing a fixed point merely involves feeding the network's output back into the input layer until it equals the previous output³.

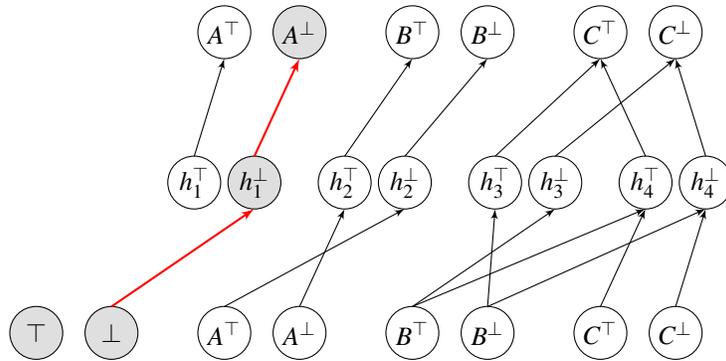
2.2.1. Example: core method computation

For an example of how this core works in practice consider the four clauses used
220 in the previous example: $A \leftarrow \perp$, $B \leftarrow \neg A$, $C \leftarrow \neg B$, $C \leftarrow B \wedge \neg C$. Application of the translation algorithm yields the following multilayer perceptron.

³ The number of iterations necessary for reaching the least fixed point can be shown to be lesser or equal to the number of atoms. The network has no inhibitory connections, so more input always generates equal or more output. Starting from an empty interpretation, each subsequent iteration must therefore activate at least one additional unit, one per atom at maximum, or stop.



The following figures show how activations propagate through the network on each iteration, starting in the input layer at the bottom. Red arrows and grey cells indicate active connections and units. As before, a fixed point is reached after three steps.



The first iteration starts off with \top and \perp , the latter of which activates the h_1^\perp neuron, which, being set as an 'or'-gate, has a threshold of 0.5 and this in turn activates A^\perp in the output, because the 'and'-gate with a single clause also has a threshold of 0.5.

245 for adjusting the weights, usually through gradient descent, so the cost function is minimized. The name backpropagation derives from the order in which the partial derivatives are calculated. This process begins in the output layer and propagates backward, layer by layer, as the calculation in each layer requires the results of its successor.

The derivation of the backpropagation algorithm is fairly general and holds for different cost and activation functions. Given the binary nature of the targets, we choose 250 the logistic cost function⁴, because it encourages binary output. As activation function, the standard sigmoid is used. The specific formulas used in the algorithm depend on the choice of function. A derivation for the algorithm with quadratic loss function, along with a general introduction to the algorithm, can be found in [27] and an analogous 255 derivation for the backpropagation that was used here is provided in the appendix. The implementation uses on-line training, which means that weights are updated every time after calculating the error for one randomly selected sample. The advantage of on-line learning over batch training for our purposes is that the former better accommodates the large variations in sample size that are encountered in different cores. 260 Aside from this, the choice is of no conceptual importance.

Additionally, in our experimental implementation the naïve backpropagation algorithm has been enhanced by using a momentum term, saving the weight adjustment terms in each iteration and adds a fraction of them to the weight adjustment in the next iteration. This tends to speed up convergence by preventing fluctuation of the weights 265 to some extent and also leads to some robustness against small local optima. Since we focused on qualitative rather than quantitative results, this is the only significant modification to the original algorithm. Where a consistent, if limited, level of learning success can be shown with our basic implementation, it is plausible to assume that further attempts with more sophisticated versions of the learning algorithm (also including, 270 for instance, techniques such as regularization, linear-least-squares initialization, or simulated annealing) will yield much better results.

⁴ $J(\vec{w}) = \frac{1}{m} \sum_{i=1}^m [(y^{(i)} \log h_{\vec{w}}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\vec{w}}(x^{(i)}))]$, with \vec{w} the vector of weights and $h_{\vec{w}}(x^{(i)})$ the output produced by the network with sample input $x^{(i)}$ as compared to sample output $y^{(i)}$.

3. Theory and implementation of learning cores

We now document the theoretical groundwork and the actual implementation which have gone into this project, beginning with modifications made to the core architecture, followed by some remarks on the learning algorithm and a thorough discussion of the proposed rule extraction algorithm. The section closes with a list of control measures used in the implementation.

Before supervised learning can be attempted in cores, three problems have to be addressed:

1. The core architecture must reach fixed points to compute results. While the existence of these fixed points has been proven for translated programs, this result does not generalize to cores whose weights have changed over the learning process. The first task, therefore, is to ensure the existence of fixed points throughout the learning process.
2. Following the example given, for instance, in [28], a differentiable activation function must be introduced, while preserving the core's semantics. The backpropagation algorithm relies on the computation of derivatives of the cost function which includes the activation functions of the network. The Heaviside function is not differentiable and must be replaced.
3. One must decide, what kinds of samples will be used for supervised learning. The core in its original form is only capable of computing the least fixed point, when starting from an empty interpretation. If one wants to capture any of the structure of the program, more than one sample is needed for training.

The following subsections address these issues in turn, before drawing all the individual steps and solutions together in a backpropagation algorithm for learning cores. The second to last subsection then introduces the rule extraction algorithm, before the final subsection shortly touches upon two measures introduced in order to assure the proper functioning of the developed methods.

3.1. Ensuring a fixed point with unipolar weights

300 Convergence to a fixed point is essential to the core method. While this property is guaranteed for Hölldobler’s and Kalinke’s discrete cores and will be proven for their sigmoidal analogs below, it is difficult to ensure it throughout the learning period, where the network may change drastically and with little regard for the structure in which it is embedded. Convergence, therefore, should be guaranteed by something
305 other than the initial setting of the weights. A possible solution to this issue, and the one employed here, is to restrict the network to unipolar weights. When limiting all non-bias weights to positive values, there are no inhibitory connections and thus the activation of the network will monotonically increase on every iteration until it plateaus at a fixed point. On the downside, the elimination of inhibitory units of course reduces
310 the modelling capacity of the network. The reason it can nonetheless be done in good conscience here, is that the translation of logic programs into cores itself only uses positive weights and thus ensures that every Łukasiewicz logic program to be learned by a core can be fully modelled, and therefore also learned, using these simpler unipolar networks.

315 A standard activation function used in feed forward networks is the sigmoid function $sig(z) = 1/(1 + e^{-z})$ where $z = \vec{w}^T \vec{x}$, the dot product of the weight vector and the incoming activations. In the implementation of unipolarity is achieved by squaring all but the bias weight in the activation function. So, while the sigmoid function remains the same, z is now computed as $w_0 \cdot x_0 + \sum_{i>0} (\frac{1}{2}(w_i)^2 \cdot x_i)$. The values stored in
320 the weight matrix may still be negative, but will effectively be treated as positive. To preserve the previous behaviour of cores, all non-bias weights are replaced by their respective square root after the translation algorithm has been applied. With this measure, the translation algorithm can ignore the modification to the activation function and act as if it was the standard sigmoid, so long as it only sets positive weights. The subsequent argument that semantics are preserved in the sigmoidal core will also assume the
325 standard activation function.

3.2. Preserving core semantics

With the introduction of sigmoidal activation to the network, the range of possible activations for each neuron changes from 0 and 1 to the interval $]0, 1[$, and what
330 it means for a neuron to 'fire' becomes less clear. To ensure compatibility with the core architecture, the network's output is discretized by rounding it half up to 0 and 1. A fixed point is reached when this rounded output is equal to the input of the current iteration. Whithin the network, however, instead of an activation value it makes more sense to define an interval bounded by a certain value $[A^+, 1]$ where all activation val-
335 ues in that interval are considered as firing, and another interval $[0, A^-]$ of activations regarded as not firing.

As these two intervals should be disjoint, it follows that $A^- < A^+$ and because the classification into firing and non-firing is integral to the way the core is built, it must be ensured that no activations in the interval $]A^-, A^+[$ are produced. Given these changes,
340 the approach of using logic gates, which was explained in the previous section, can be maintained, but must deal with the following complications. Because the output of a non-firing neuron is no longer 0, and can take on values up to A^- , an 'or'-gate must ensure that it won't fire, even if all connected neurons send an activation of A^- each, while at the same time guaranteeing that it will fire if one neuron sends activation
345 A^+ and all others send nothing. Similarly 'and'-gates should fire when all connected neurons from the previous layer send A^+ , but not if all but one send an activation of 1 and the last one sends A^- . It becomes clear that these constraints of maximal non-activating input and minimal activating input⁵ can only be satisfied with the right choice of A^- and A^+ . In the core, both A^- and A^+ are determined by the value of
350 ω . If ω is large, A^- and A^+ , approach 0 and 1 respectively. For a small ω , both values lie close to $1/2$. It can be shown, that semantics of the network are preserved if $\omega > 2\log(2deg - 1)$, where deg is the maximum in-degree among neurons in the output layer. The formal proof for this can be found in the appendix.

⁵The terms are used here in a different context than later on during rule extraction.

3.3. Fixed point calculation with initial activation

355 The original core architecture serves to compute fixed points for a given logic program and no additional input. Evidently, this one sample of (empty) input and corresponding output does not contain exhaustive information about the program which produced it. To train a network which captures the functionality of the program, more samples are required. Given the context of logic programs, it seems like an obvious
360 choice to generate additional samples for possible interpretations of the atoms. There are 3^n possible interpretations for a set of ternary logic formulas P with n atoms. What additional inferences P allows, based on a partial interpretation, provides information about P , and having this information for all 3^n interpretations specifies P to its semantic equivalence class.⁶

365 3.3.1. C-interpretation

A naïve approach for using such partial interpretations in a core is to enforce them as the base activation while running the core, and see what additional inferences are drawn before reaching a fixed point. This is achieved by adding the interpretation to every starting activation on the first iteration as well as to the output at the end of
370 every iteration. This method must be called naïve because the underlying definition of interpretations, while applicable to Łukasiewicz logic, actually makes very little sense for the weakly completed logic programs at hand. Determining the value of an atom from the outset, while leaving the program as is, takes away both the non-monotonicity and the property of non-contradiction. On the plus side, only few changes to the core
375 are necessary to accommodate this interpretation, which will from now on be referred to as C-interpretation.

3.3.2. Ł-interpretation

In order to preserve the semantics of weakly completed Łukasiewicz logic, an alternative Ł-interpretation is proposed. Here the process will be handled slightly differently, as it seems more adequate to model different interpretations in such a way that
380

⁶ If the interpretation leads to no contradictions, it is a model of P , otherwise it is not. Knowing all possible interpretations of P provides the full extension.

they represent logic programs in their own right. As such, setting an atom A to \top or \perp in the interpretation should have the same effect as adding a positive or negative fact to the program. Doing this preserves the important property, that the Stenning-van-Lambalgen consequence operator always reaches a model. Note that in this choice of
385 interpretation setting atoms to false only has an effect, if they do not occur as heads of clauses in the program, and that setting atoms to true in the interpretation will prevent the consequence operator from inferring these atoms to be false even if all other clauses in the program would lead to this conclusion.

Going with the interpretation as adding clauses to the program, the most intuitive
390 approach would be adding neurons to the hidden layer of the core which represent these rules. Unfortunately this does not seem like a viable option. The addition of neurons would change the in-degrees of some of the core's output units which would in turn necessitate an update of their respective bias weights, in order to maintain Łukasiewicz semantics. For each interpretation there could be changes to the whole network which
395 would not only be computationally costly but also pose problems in the context of learning, where changes to the core network should likely be limited to the learning algorithm itself.

Instead it appears more promising to adjust the way in which the inputs to the network are generated and outputs are interpreted. For negative facts like $A \leftarrow \perp$ this
400 is can be done fairly simply. Given weak completion these clauses only affect the program at all, if there is no other clause with head A in the program. In this case A will be set to \perp and keep this value, as there is no other clause to change it. This can be modelled in the core by checking the in-degree of the atom's associated output unit in the network. If the in-degree is 0, A is set to \perp in the input to the network on every
405 iteration as well as on the final output, which in the neural net means the activation of the neurons corresponding to A is $(0, 1)$. Positive facts of the type $A \leftarrow \top$ have to be treated differently. If such a clause exists, A will be true independently of the rest of the program. This means A should be set to \top on all inputs as well as the final output, i.e. the activation is set to $(1, 0)$. Because the clause is part of the interpretation
410 and not translated into the network, the network will still produce activation in the A^\perp output neuron, whenever all program clauses with A in the head are false. The arising

contradiction must be resolved and the easiest way of doing this is to set activation of the A^\perp neuron in the output to 0 on all inputs and the final output.

3.3.3. C^* -interpretation

415 In addition, a form of \mathbb{L} -interpretation will be tested, which is different only in that it leaves out the contradiction resolution step. The resulting new C^* -interpretation can be viewed as using explicit negation, rather than the negation by failure present under \mathbb{L} -interpretation, with regards to elements of the interpretation. Unfortunately, this makes the logic monotonic and allows for contradictions. C^* -interpretation is nonethe-
420 less of interest because, as will become obvious from the results reported below, it can be trained better than \mathbb{L} -interpretation but still bears some similarities. Training under C - and C^* -interpretations performs so similarly that C -interpretation will not be discussed separately in the empirical results.

It is clear that all three interpretations generate the same output for empty inter-
425 pretations. Furthermore it can be shown, that all non-contradictory models under C^* -interpretation are equivalent to the \mathbb{L} -interpretation under the same input⁷. All three interpretations have been implemented and tested.

3.4. Backpropagation in cores

With the modifications to the core that have been described above it is now possible
430 to create samples and test the core's capacity for learning. In the given set-up, two cores are used. The first core is generated by translating the complete program into it and is subsequently run with all possible inputs computing the desired outputs, the pairs of which will be used as training samples. Core number two is created based on a partial version of the program, where some clauses have been deleted. The learning task now,
435 is to train the second core with samples from the first one and see whether it can learn the missing parts of the program.

It may take the core multiple iterations to reach a fixed point, but only the last one is used for training. For a given sample, the core is run on the sample input and when

⁷A sketch of this proof is provided in the appendix

reaching a fixed point returns activation values of all the networks layers. Note that
440 the activation of the output layer is not the final output of the core, which may contain
modifications from interpretation or contradiction-resolution. The backpropagation al-
gorithm is then applied to the network with that activation. Due to the non-classical
activation function used in the network, the algorithm differs slightly from its more
common form. The derivation of the relevant formulas found in the appendix is done
445 analogous to the proof of standard backpropagation found in [27].

3.5. Rule extraction

The algorithm for extracting information from a core discussed in this section fo-
cuses on C- and C*-interpretation. It has been pointed out previously that through
iterations the core’s activation increases monotonically under these interpretations, due
450 to a lack of inhibitory connections in the network. The same reasoning ensures mono-
tonicity with regard to interpretations. While the number of possible interpretations
rises exponentially with the number of atoms, diminishing hopes for a tractable rule
extraction algorithm, the property of monotonicity allows for heavy pruning, making a
viable solution at least for small sample sizes possible.

455 Our algorithm is inspired by the approach for knowledge extraction and the corre-
sponding algorithm for regular networks from [29]. Still, the method presented here
warrants an independent introduction as well as analysis for soundness and comple-
teness.

3.5.1. The basic extraction algorithm

460 The algorithm extracts all minimal activating and all maximal non-activating inputs
for each output neuron of the network, which can then be used to compute the logical
rules generating this activation. In the following, the set of all inputs to the network
will be looked at as a partial order with the input vector of zeros as bottom element
and the vector of ones as top element. Input vectors are ordered in such a way that
465 $v_1 \geq v_2 \Leftrightarrow \forall i : v_1[i] \geq v_2[i]$.

For each output neuron separately, the algorithm traverses the space of all possible
interpretations by advancing alternately an upper and a lower boundary of interpre-

tations starting from top and bottom element. The new boundaries are generated by computing all *direct successors* of each element of the existing boundary. For an element of the lower boundary a direct successor is a copy of the element in which exactly one activation is changed from 0 to 1. The direct successor of an element in the upper boundary, analogously, has one active input less than that element. All inputs connected through a series of direct successions will be called successors and the definition for predecessors follows from this. An input in the lower boundary is said to be *subsumed* by an activating input if it is a successor of that input and is subsumed by a non-activating input, if it is a predecessor of that input. For subsumption in the upper boundary, successor and predecessor relations are reversed. In either case, the target-activation produced by the subsumed input is equal to, and therefore determined by, the other input. Whenever an activating input is found in the lower boundary, which is not subsumed by an input already stored in the set of minimal activating inputs, it is added to that set. The progression through a lower boundary ensures that all predecessors have already been checked and the one that has been found is in fact minimal. Also, if all direct successors of a non-activating input are activating, then that input is added to the list of maximal non-activating inputs. In the upper boundary, relevant inputs are found in an analogous manner, where activation is the default. The algorithm terminates once the two boundaries have passed by one another, which is not implemented explicitly, but a result of the pruning mechanisms discussed below.

Prior to pruning, the soundness and completeness of the extraction algorithm are evident, but spelled out here for the sake of completeness. All activating inputs found in the lower boundary, which are not greater than previously found ones are minimal, as all smaller activating inputs would be in that set. Complementary, all non-activating inputs whose direct successors are activating are maximal, as all their successors are activating due to monotonicity. The analog holds for the upper bound. Thus the extraction is sound. All minimal activating inputs are found by the algorithm, as they are activating and not subsumed by any other activating inputs. All maximal non-activating inputs are found by the algorithm as the successors of each are activating by definition. Here, too, the analog is true for the upper bound and so the extraction is complete as well. Later it will be shown that both properties are preserved under pruning rules.

3.5.2. Pruning

500 Due to the monotonic nature of the space of possible inputs, once an activating input is found in the lower bound, none of its successors need to be investigated any more, as all of them will be activating as well. An efficient extraction algorithm must therefore limit its exploration of the space of possible inputs to the relevant nodes which may hold new information. The pruning is best explained from the perspective of one
505 of the boundaries. From the perspective of the lower boundary, minimal activating inputs will be called *rules* and maximal non-activating inputs are called *anti-rules* (for lack of a better term). These terms are relative to the boundary, such that rules in the lower boundary are anti-rules in the upper boundary and vice versa. When a rule is discovered, all of its successors should be pruned, as their values will hold no new
510 information. This must be ensured both in the current boundary, where it is a rule, and the opposite boundary, where it is an anti-rule and two pruning mechanisms ensure this.

The pruning mechanisms can not be explained without covering the specifics of the algorithm in some detail. It may help to have a look at the pseudo code provided in the
515 appendix for reference.

To avoid too much confusion, the algorithm refers to inputs as vertex objects, owing to the graph-like feel of the partial order. Alongside its input value and a number of other things, each vertex stores a memory array to keep track of the direct successors it should generate and those that should be pruned. This array has one entry for each
520 neuron, which is 1, if switching the value of this input from 0 to 1 or vice versa will generate a valid successor, 0 if the successor and all subsequent successors are invalid, and -1 if the direct successor is invalid, but later ones may be valid.

The *test* function checks whether a vertex is a rule or subsumed by an anti-rule. If the vertex is a rule, the memory array is set to all zeros, so that no successors are
525 generated and the vertex is added to the set of found rules. If the vertex is subsumed by an anti-rule, some, but not necessarily all successors will be subsumed as well. In all places where switching the input would generate a subsumed direct successor, the memory array is set from 1 to -1.

This information is utilized in the *successors* function. As the name suggests, this
530 function creates the direct successors of a given vertex, but also uses the step to ex-
change pruning information among the successors. Firstly, the input vertex is *tested*
again, in case new anti-rules were discovered while traversing the other boundary.
Then, each valid direct successor of the vertex (as indicated by the memory array) is
looked up and if it does not exist yet, is generated and tested. For each such successor
535 which is a rule, the preceding vertex's memory is set to 0 at the index, which was used
to generate that successor, indicating that this successor should not be investigated fur-
ther. After this has been completed, all the successors are traversed for a second time
and all 0s from the vertex memory are copied into their respective memories as well.
This way, each successor receives information about all vertices, with which it shares
540 a common direct predecessor. Now, when a rule is discovered, all of its direct prede-
cessors will set the index in their memory which generated this rule to 0 and pass this
information on to all their successors. If a vertex subsumed by the rule were to be gen-
erated, it would have to have a direct predecessor which is not subsumed by the rule (or
the problem propagates down until this condition occurs). This predecessor, however,
545 must be a successor of one of the rule's predecessors. Therefore it would not generate
that vertex and it follows that no vertices subsumed by rules are generated. Note that
the same does not hold for anti-rules. Finally, the successors function also serves to
determine, whether the given vertex is an anti-rule. This is the case, if the vertex is not
subsumed by an anti-rule (i.e. no entry in the memory is set to -1) and no generated
550 successor has the same target-activation value as the vertex. As rules trivially share
these properties by having all their successors pruned, they must be filtered out. This
is done by checking for the right target-activation, given the vertex's boundary, before
adding it to the set of rules. In the lower bound an anti-rule must be non-activating,
and activating, if in the upper bound. So now, when the *rules_for_target* function fi-
555 nally creates the two boundaries and traverses the partial order, the pruning ensures no
vertices that are subsumed by known rules are looked at.

The anti-rule related pruning is handled with some care by the algorithm, as it can
lead to problematic cases. In general, it might happen, that all direct successors of a
vertex are subsumed by some anti-rules. In such a case, declaring all direct succes-

sors invalid could hinder the generation of valid ones down the line. If, for example, (1, 1, 1, 1, 0, 0, 0) and (1, 0, 0, 0, 1, 1, 1) were known anti-rules, the lower boundary input (1, 0, 0, 0, 0, 0, 0) would generate no direct successors and unsubsumed successors like (1, 0, 0, 1, 1, 0, 0) would not be reached. In the algorithm this problem is solved by looking only at the first anti-rule found, rather than the whole set of subsuming anti-rules. Apart from the trivial case where the anti-rule is the top or bottom element (which ends the search as either all inputs are activating or none are), a single anti-rule, will not prune all successors of the vertex. The indices at which the anti-rule differs from its predecessors (of which, barring the trivial cases, it has at least one) can be changed in the vertex to generate valid successors.

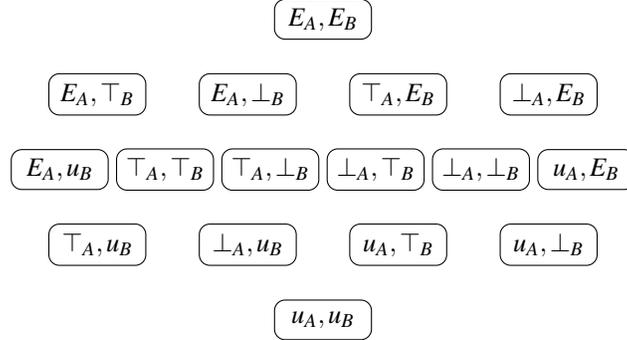
Of course this strategy can, at times, dismiss useful information and generate vertices which are subsumed by known rules at the point of creation. As a stand-in for more elaborate methods it will suffice to generate some first results.

In order to ensure that soundness and completeness are maintained, it must be proven that pruning neither changes the results of examining a particular vertex, nor prevents any rule vertices from being examined. Addressing part one, the test for rules functions in the same way as without pruning and only relies on the vertex's target activation value, which is not affected by pruning. With pruning, the test for anti-rules employs the memory of the given vertex, rather than looking at all successors, but it does so, only to infer the target-activation values of the invalid immediate successors. Assuming that rules have been identified correctly, every 0 in the memory is linked to a successor which has a different target-activation value than the vertex. Each -1 is linked to a successor which has the same value as the vertex. Generating each of these invalid successors would take more time but yield the same result. Therefore all examined vertices are still classified correctly. The second part follows from the soundness of the pruning algorithm.

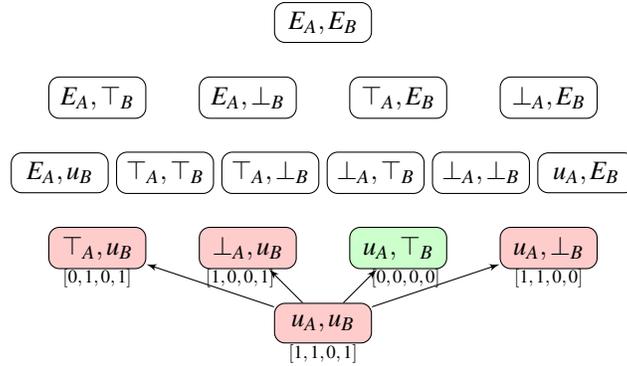
3.5.3. Example: rule extraction

As an illustration of the extraction algorithm, we consider a network with only two literals A and B , which encodes the clause $A \leftarrow B$. The range of possible inputs forms a partial order of 16 elements. In the figure below, E denotes the error input, where both

590 \top and \perp unit of a literal are active.



The extraction algorithm is applied for each of the four units (A^\top , A^\perp , B^\top , B^\perp) separately. We begin by looking at A^\top . Memory cells of highlighted elements are encoded as a 4-tuple, which indicates changes to the units in the same order. So, for example, element (\top_A, u_B) with memory $[0, 1, 0, 1]$ will generate the successors (E_A, u_B) and (\top_A, \perp_B) .



Starting at the bottom element, we find that zero activation in the input does not activate A^\top . Going through the four immediate successors, (\top_A, u_B) and (\perp_A, u_B) both turn out as non-activating inputs and are queued up for the following iteration. Then (u_A, \top_B) is tested and turns out to be a minimal activating input, i.e. a rule, and the memory of the bottom element is adjusted accordingly from $[1, 1, 1, 1]$ to $[1, 1, 0, 1]$. After testing (u_A, \perp_B) and adding it to the queue as well, the memory of the three successors is updated by the bottom element.

3.6. Controls

Verifying the correctness of an implementation as a whole beyond the checking of test cases is usually associated with an enormous effort and has therefore not been in the scope of this project. Checks have been installed at two crucial steps in the program
625 which are worth mentioning:

1. A function has been implemented which can run a core with discrete activation units as in the original algorithm (with the one difference that it squares all non-bias weights to compensate for the fact that they were reduced to their square roots in the new translation). The function is used to run a core with both kinds
630 of activation and for all possible inputs under a chosen interpretation, returning an error if the activations reach different models for the same input. This way, it is possible to verify that the implementation of the translation algorithm with sigmoidal units does preserve the semantics of the cores.
2. As a standard measure to ensure the correctness of the implementation of the
635 backpropagation algorithm, numerical gradient checking has been implemented. As gradient checking is computationally costly, it is only used on a small number of training samples to verify the correctness of backpropagation and disabled for the actual training of the network. Given the more complicated nature of learning in a core as compared to the classical application of backpropagation
640 there are a number of other errors that may occur which prevent learning and are not detected by gradient checking, but the method still serves to eliminate one common source of errors.

4. Empirical results

Test following test results are based on an implementation in Julia.⁸ The source
645 code of our implementation is open source and available for download from GitHub.⁹ As already stated before, a thorough quantitative analysis of training results is not within

⁸ julialang.org

⁹ GitHub ID omitted for blind review.

the scope of this article. Instead, several exemplary cases will be used to highlight consistently observed features of the learning process.

The first such program is displayed below in the format, in which it is read by the
 650 implementation. To keep things simple and in ASCII-encoding, \leftarrow , \wedge , and \neg have been written as `<-`, `&` and `-` respectively. The partial program, consisting of clauses 1, 5 and 6, is translated into a core and then trained on samples generated from the full program.

4.1. Comparison of C*- and L-interpretation

<pre>a <- b & -d & -e b <- a c <- b d <- FF e <- c & d e <- -a & -b</pre>	<pre>a <- b & -d & -e e <- c & d e <- -a & -b</pre>
(a) full program	(b) partial program

Figure 3: P1

The first example program illustrates that there are cases in which the method yields
 655 good results under C*-interpretation. Training was done with learning rate and momentum of 0.05 under C*-interpretation and 0.02 under L-interpretation. During the learning process a number of parameters are measured and reported for intermediate results after every 200 training steps (500 in later examples). *%corr* indicates the percentage of correctly classified training samples, while *eTotal* measures the total number
 660 of errors, i.e. the number of incorrect rounded outputs over all output neurons and all samples. In addition, *avgIt* gives the average number of core iterations over all samples and *costJ* is the total value of the error cost function. If the learning algorithm functions correctly, one would expect a steady decline in the cost function, followed by decrease of the total number of errors, which, in turn, leads to an overall rise in the amount
 665 of correctly classified samples. Since *%corr* does not differentiate between samples with just one error and samples in which every single neuron is wrongly classified, the latter correlation can be quite weak. Especially when the overall number of errors is

still high, $\%corr$ may even increase, as $eTotal$ is reduced, but more evenly distributed across the samples. A similar distribution of errors may also happen in the relationship
 670 between cost function and number of errors.

For the C- and C*-interpretation samples, training of the first test program tends to converge after two to three thousand iterations. Depending on the random initialization, usually one of two optima is reached, the first one being at around 80% correct classification, the second one at 100%.

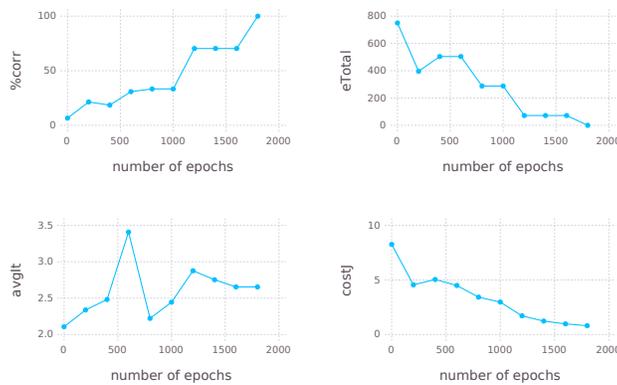


Figure 4: PI training results under C*-interpretation

Under Łukasiewicz interpretation, the cost function can still be seen to generally
 675 decrease, though not always monotonically, before it starts to fluctuate. Choosing
 smaller learning rates remedies the fluctuation to some extent, but in many cases the
 algorithm does not seem to converge even for small learning rates.¹⁰ The graphs also
 show less correlation between the cost function and the total number of errors. This can
 680 be attributed to the conflict resolution mechanism active under Ł-interpretation, which
 may generate errors on the final output that are not accounted for in the cost function.

¹⁰Choosing very small learning rates (in the given example the threshold is at around 0.00003) leads to a steady increase of the cost function. This has been observed across all examined cases but the source has not been found.

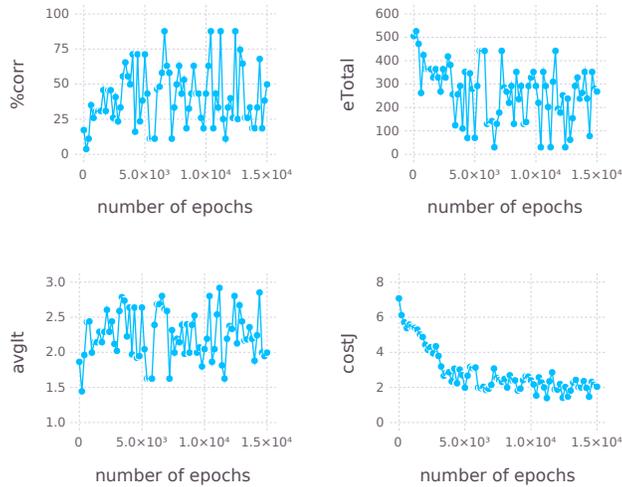


Figure 5: P1 training results under \mathcal{L} -interpretation

4.2. Correlation between error measures

Often times, as displayed by the second example, the learning process quickly gets stuck in local optima under \mathcal{C}^* -interpretation. While the cost function converges, the average number of iterations keeps fluctuating. Still, this does not seem to affect the classification performance. Under \mathcal{L} -interpretation, results are less clear-cut, but the correlation between $costJ$, $eTotal$ and $\%corr$ is still clearly recognizable.

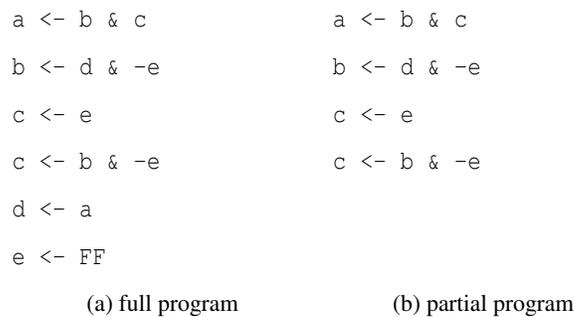


Figure 6: P2

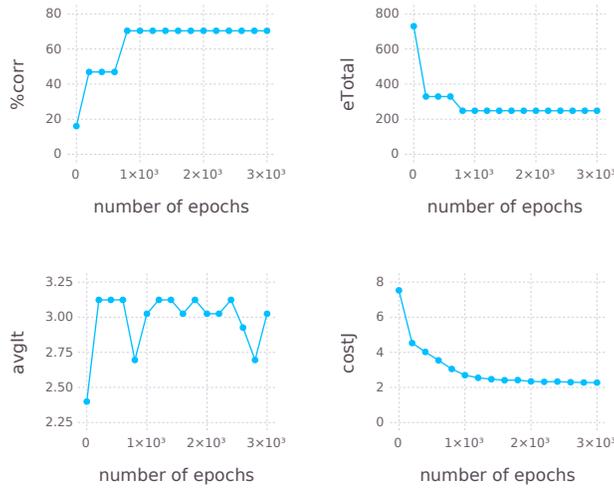


Figure 7: P2 training results under C*-interpretation

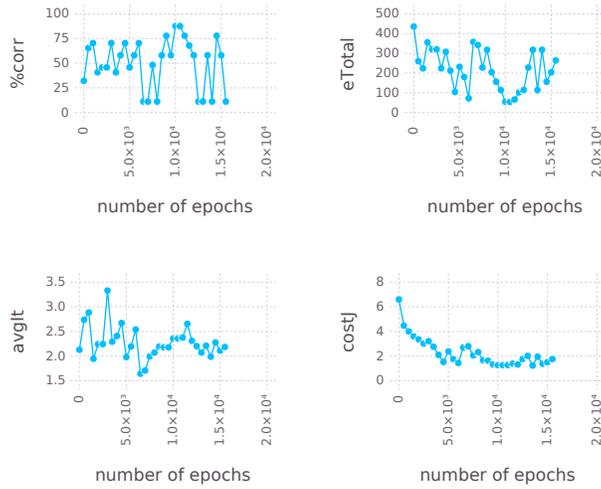


Figure 8: P2 training results under L-interpretation

4.3. *The problem of hidden errors*

The third program to be examined contains eight atoms, three more than the previous programs. This adds six neurons to the network and increases the number of training samples from 243 to 6561. The fluctuation in the plots can in part be explained by this fact. Steps of 500 samples make up less than 10% of the total sample size and may at times lead the algorithm in different directions.

What is interesting about this example, is how the algorithm can be observed plummeting in overall performance in the first 500 training steps. This can be attributed to two factors. Under the logistic cost function, which incentivises many smaller errors over fewer large ones, distributing the error may serve to reduce the overall cost, but increase the total amount. Secondly, the backpropagation algorithm is based on the network output. And as the final output is created only after all facts from the input, backpropagation will perceive all misclassifications which are fixed by this final step. Correcting for these errors will decrease the cost function, but does not increase the core's performance. Under \mathbb{L} -interpretation, the contradiction resolution step contributes to this problem with an additional layer of error correction invisible to the learning algorithm. Moreover, this step cannot be modeled without inhibitory connections, which leaves the algorithm trying—and failing—to correct an error that does not actually exist. There may be multiple reasons for the weak performance under \mathbb{L} -interpretation, but this is certainly one of them.

In this example, these shortcomings are underlined by the fact that the initial performance from partial background knowledge is much better than the local optimum reached through training.

```

a <- TT
b <- a & -c
d <- c
d <- e
f <- e & -b
f <- a & g & -d
h <- FF

```

(a) full program

(b) partial program

Figure 9: P3

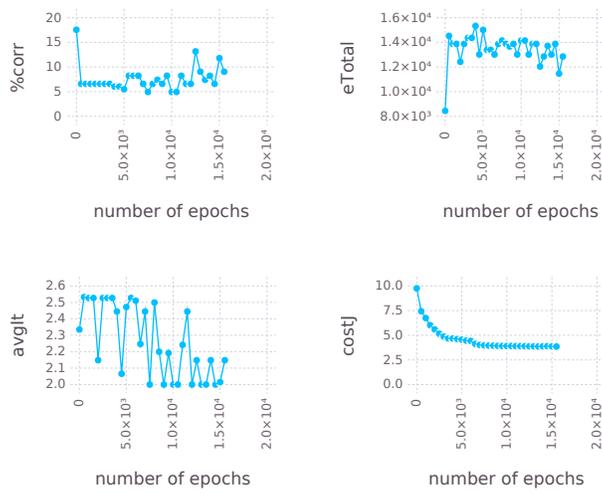


Figure 10: P3 training results under C*-interpretation

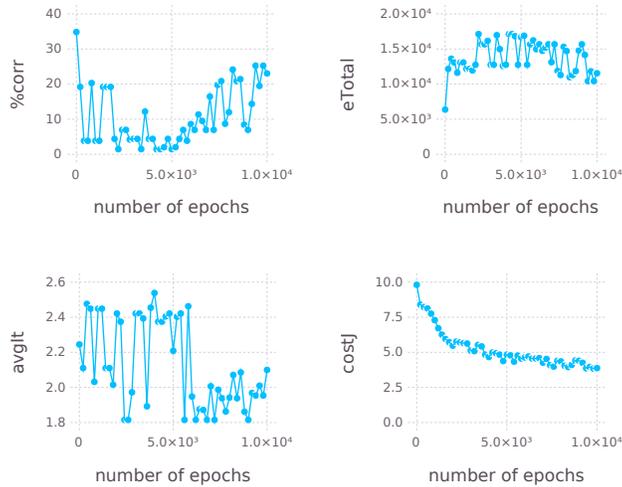


Figure 11: P3 training results under L-interpretation

4.4. Core compression

The final program examined here, is simply a long chain of inferences. The training results are meant to highlight a phenomenon that is less pronounced but observable in most trained cores. In this extreme case the partial program starts with an average number of iterations of 4.14, which immediately drops to around 2, changing very little thereafter. This number includes the last iteration where input and output must be equal. Therefore, in all but very few cases, the inference is compressed into one iteration. In general, trained cores tend to have fewer iterations on average than their translated counterparts. This can be explained by the fact, that the backpropagation algorithm does not take multiple iterations into account and optimizes for correct output after just one iteration. This property does not decrease the performance of trained cores, but it does suggest that they do not fully utilize the core-architecture.

b ← a	b ← a
c ← b	c ← b
d ← c	d ← c
e ← d	e ← d
f ← e	f ← e
g ← f	g ← FF
g ← FF	

(a) full program

(b) partial program

Figure 12: P4

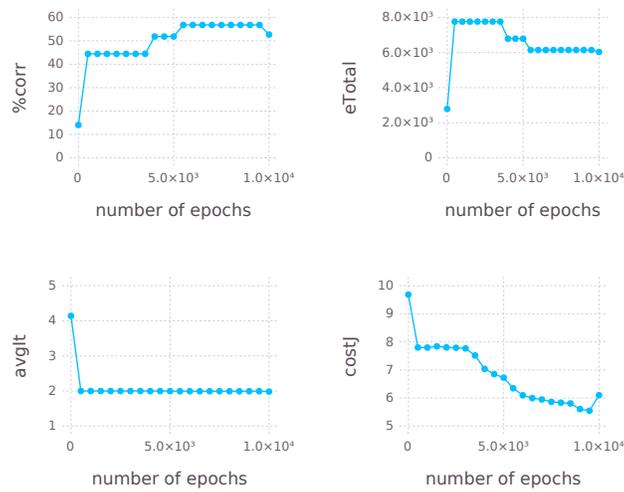


Figure 13: P4 training results under C*-interpretation

4.5. Analysis through rule extraction

The developed rule extraction method is not yet suited to provide a complete picture of the information contained in a core, but it may be used to take a look at individual neurons and their activation rules. This can be done both for translated and trained cores to see, what differences remain after training.

out(d^\top)	a^\top	a^\perp	b^\top	b^\perp	c^\top	c^\perp	d^\top	d^\perp	e^\top	e^\perp	f^\top	f^\perp	g^\top	g^\perp	h^\top	h^\perp
translated 1	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
translated 2	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0
translated 3	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0
trained 1	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
trained 2	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0
trained 3	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0

Figure 14: Extracted minimal activating inputs of d^\top in P3

out(f^\top)	a^\top	a^\perp	b^\top	b^\perp	c^\top	c^\perp	d^\top	d^\perp	e^\top	e^\perp	f^\top	f^\perp	g^\top	g^\perp	h^\top	h^\perp
translated 1	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0
translated 2	0	0	0	0	1	0	0	0	1	0	0	0	0	0	0	0
translated 3	0	0	0	0	0	1	0	0	0	1	0	0	1	0	0	0
trained 1	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0
trained 2	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0

Figure 15: Extracted minimal activating inputs of f^\top in P3

In the third example above, for instance, in which the core was generated from the program P3 with two missing clauses $d \leftarrow c$ and $d \leftarrow e$ and then trained, it turns out that the d^\top neuron's activation rules in the trained core match the full program. While learning was successful with regard to d^\top , d^\perp does not show any activation in the trained core, whereas the core containing the complete translated program has an activation in d^\perp when both c^\perp and e^\perp are active.

These findings alone do not explain, why the trained core classifies less than 10% of samples correctly. Looking further, it can be found that other inference structures have largely broken down. For instance, f^\top has three activation rules in the core containing the complete program. In the trained core, two rules are extracted, one of which is wrong. This does not mean, that the connections to the f^\top output neuron are necessarily wrong. Due to the core's multiple iterations, the activation patterns of different neurons are highly interdependent and errors are hard to localize.

5. Closing discussion and future work

With this project we set out to investigate the applicability of the backpropagation algorithm to Łukasiewicz logic cores. This goal has been met, be it with mixed results. In order to train the cores, a number of modifications had to be made. Sigmoidal
745 activation units were introduced, along with the proposition of unipolar weights, and several options for the generation of training data have been discussed. These steps have been motivated and provided with formal proofs where it was considered necessary, resulting in a trainable core as an intermediate result. Further research on the topic, for example using other supervised learning methods, can be based on this type
750 of core as well. The training process itself is more problematic, as training of cores under \mathbb{L} -interpretation clearly does not perform as well as desired. Whether this weak performance can be sufficiently improved by standard augmentations of the algorithm or whether the contradiction resolution mechanism used in \mathbb{L} -interpretation prevents backpropagation from functioning properly on a more general level, is not clear at the
755 moment. In addition, a rule extraction algorithm has been proposed for cores trained under \mathbb{C} - or \mathbb{C}^* -interpretation. This is very much a work in progress, but should also in its current state provide a starting off point for other ventures in that direction. In the following, a number of problems of the project are discussed, whose resolution would offer ways of tackling the issue of weak performance. They are followed by a list of
760 more specific objectives for future endeavors to improve the project.

5.1. Challenges

In this subsection we shortly outline three open questions which we consider the main challenges for future work continuing from the described stage of development.

5.1.1. Proof for last-iteration backpropagation

765 At this point our work is lacking a proof for the functioning of backpropagation in the way, in which it has been applied. While it is an established fact, that an MLP with a sufficient number of hidden units will be able to model the behaviour of the types of logic programs presented here, and it is evident that a unipolar MLP embedded in the core structure is powerful enough to accomplish this task as well, it is less clear how a

770 unipolar core can be trained to reach this performance. The method of training the core
only on the last iteration, in which the input equals the output, does not take the core's
capacity for multiple iterations into account. It also means that the algorithm is only
guaranteed to work as intended on those inputs which are themselves fixed points. One
indication that backpropagation may not be the method of choice for training cores is
775 the empirical result, that the average number of iterations in the core tends to decline
rapidly over the learning process, before usually settling somewhere around 2, which
is also the minimum possible number for all non-fixpoint inputs. This suggests that
information about the program is compressed in the network, leading to redundancy,
rather than building on itself, as seen in an untrained core.

780 5.1.2. *Missing link between C* and Ł*

Another problem which remains unsolved at the current stage is the disparity be-
tween Ł- and C*-interpretation. As the empirical results indicate, Ł-interpretation is
a lot harder to train and the C- and C*-interpretations have been used, in large part,
as a stand-in, so learning and rule extraction could be explored, in the hope that the
785 gap to Ł-interpretation could be bridged later on. In case this problem cannot be ad-
dressed in a satisfying manner also in the future, another route would be to explore,
how much training a network on C*-interpretation samples can improve performance
on Ł-interpretation test sets. While C*-models are not generally equal to Ł-models,
their similarities might be sufficient to motivate learning on one interpretation in order
790 to improve performance on the other.

5.1.3. *Implausibility of C*-samples*

The problem with this last option and learning on C*-samples in general is that,
while it is easy to produce C*-interpretation samples for training in the course of this
project, any actual application scenario for Łukasiewicz cores will naturally provide
795 Ł-model samples. A method for translating them into their C*-model equivalents has
not yet been put forward and remains an open question for future work.

5.2. Improvements

While the previous subsection focused on open problems, we now want to hint at three fairly straightforward approaches to improving the performance of the described system beyond the current status quo.

5.2.1. Modified backpropagation for better results

As has been outlined in the introductory section, there are a great number of modifications that may be applied to the backpropagation algorithm in order to boost its performance. The version of on-line backpropagation with momentum and standard gradient descent used here is sufficient to provide qualitative results, i.e. whether or not the system improves performance through learning at all. A quantitative analysis of how well a core can be trained would be the consistent next step. Such an analysis should employ some additions to backpropagation likely including a better initialization method for the weights of added neurons and running of multiple initializations. Performance can then be measured by standard means of crossvalidation, partitioning the samples in training, test, and validation set, to compare the results to other approaches with a sample size that seems plausible for application scenarios.

5.2.2. Exploration of other optimization methods

Since backpropagation is only one of many optimization algorithms for neural networks, it is also promising to look at other more general methods. Most problems with the current implementation likely stem from the tight focus of backpropagation on the neural network, which fails to take the rest of the core architecture into account. Genetic algorithms ([30]) are likely a good fit, as they grant more freedom in defining fitness criteria. These could, for example, incentivise a higher number of iterations, preventing the network from condensing all information into one iteration.

5.2.3. Completion of the extraction algorithm

The proposed extraction algorithm is another area where incremental progress may be achieved. In its current state the method is not finished, because it is missing a translation of discovered rules into actual logic program clauses. In trained cores with less

825 than perfect classification this task bears some challenges, as they may not represent a
clean logic program. The relative values of completeness and soundness must then be
weighed against each other in constructing translation methods, a task which warrants
its own thoughtful investigation.

A further goal is to adapt the algorithm to proper Łukasiewicz models. However,
830 doing so may turn out to be so inefficient, due to the loss of monotonicity, that other
extraction algorithms should be considered instead. A property which may still prove
useful for pruning is the monotonicity of Ł-models with regard to the addition of neg-
ative facts.¹¹

5.3. Closing remarks

835 Several problems remain to be solved before the neural-symbolic cycle for three-
valued Łukasiewicz cores can be fully closed. Still, the exploration performed as
part of our project and described in the present article provide reference points for
future work on this topic. Work which will contribute to our general understanding,
of how symbolic and statistical systems interact and hopefully lead to a point, where
840 neural-symbolic methods can be employed to tackle those problems, that are difficult
to address with either paradigm exclusively.

Funding information

This research did not receive any specific grant from funding agencies in the public,
commercial, or not-for-profit sectors.

¹¹This property can be proven as follows: Adding a negative $A \leftarrow \perp$ fact either changes the value of A
from u to \perp or not at all. For each atom B : If $B = \top$, then there is a clause $B \leftarrow body$ with $body = \top$ and
as $body$ is a conjunction of literals it cannot have contained A with value u . If $B = \perp$, then for all clauses
 $B \leftarrow body$, $body = \perp$. There is at least one literal $L = \perp$ in each $body$ and the values of the other conjuncts
do not matter. Therefore each atom which had a value other than u before the addition of $A \leftarrow \perp$ will retain
that value.

845 **References**

- [1] F. Hsu, *Behind Deep Blue: Building the Computer that Defeated the World Chess Champion*, Princeton University Press, 2002.
- [2] J. A. Robinson, A. Voronkov (Eds.), *Handbook of Automated Reasoning*, MIT Press, 2001.
- 850 [3] T. Winograd, *Understanding Natural Language*, Academic Press, 1972.
- [4] J. Haugeland, *Artificial Intelligence: The Very Idea*, MIT Press, Cambridge, MA, USA, 1985.
- [5] W. Dai, S. H. Muggleton, Z. Zhou, Logical vision: Meta-interpretive learning for simple geometrical concepts, in: K. Inoue, H. Ohwada, A. Yamamoto (Eds.),
855 *Late Breaking Papers of the 25th International Conference on Inductive Logic Programming*, Kyoto University, Kyoto, Japan, August 20th to 22nd, 2015., Vol. 1636 of *CEUR Workshop Proceedings*, CEUR-WS.org, 2016, pp. 1–16.
- [6] R. Plamondon, S. N. Srihari, Online and off-line handwriting recognition: a comprehensive survey, *Pattern Analysis and Machine Intelligence*, *IEEE Transactions on* 22 (1) (2000) 63–84.
860
- [7] C. Berger, B. Rumpe, Autonomous driving 5 years after the urban challenge: The anticipatory vehicle as cyper-physical system, in: *2012 Proceedings of the 27th IEEE/ACM International Conference on Automated Software Engineering (ASE 2012)*, 2012, pp. 789–798.
- 865 [8] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, D. Hassabis, Human-level control through deep reinforcement learning, *Nature* 518 (7540) (2015) 529–533.
- 870 [9] A. Garcez, T. R. Besold, L. de Raedt, P. Földiak, P. Hitzler, T. Icard, K.-U. Kühnberger, L. Lamb, R. Miikkulainen, D. Silver, *Neural-Symbolic Learning and*

Reasoning: Contributions and Challenges, in: AAI Spring 2015 Symposium on Knowledge Representation and Reasoning: Integrating Symbolic and Neural Approaches, Vol. SS-15-03 of AAI Technical Reports, AAI Press, 2015.

- 875 [10] A. Garcez, K. Broda, D. M. Gabbay, *Neural-Symbolic Learning Systems: Foundations and Applications*, Perspectives in Neural Computing, Springer, 2002.
- [11] S. Bader, P. Hitzler, Dimensions of neural-symbolic integration - a structured survey, in: S. Artemov (Ed.), *We Will Show Them: Essays in Honour of Dov Gabbay (Volume 1)*, King's College Publications, 2005.
- 880 [12] J. A. Fodor, Z. W. Pylyshyn, Connectionism and cognitive architecture: A critical analysis, *Cognition* 28 (1–2) (1988) 3 – 71.
- [13] L. Valiant, *Probably Approximately Correct: Nature's Algorithms for Learning and Prospering in a Complex World*, Basic Books, 2013.
- [14] H. Gust, K.-U. Kühnberger, P. Geibel, Learning models of predicate logical theories with neural networks based on topos theory, in: *Perspectives of Neural-Symbolic Integration*, Springer, 2007, pp. 233–264.
- 885 [15] S. Bader, P. Hitzler, S. Hölldobler, A. Witzel, A fully connectionist model generator for covered first-order logic programs, in: M. M. Veloso (Ed.), *Proceedings of the 20th International Joint Conference on Artificial Intelligence (IJCAI-07)*, AAI Press, 2007, pp. 666–671.
- 890 [16] S. Hölldobler, Y. Kalinke, Toward a new massively parallel computational model for logic programming, in: *Proceedings of the Workshop on Combining Symbolic and Connectionist Processing*, held at ECAI-94, 1994, pp. 68–77.
- [17] A. d'Avila Garcez, L. Lamb, D. Gabbay, Connectionist modal logic: Representing modalities in neural networks, *Theoretical Computer Science* 371 (2007) 34–53.
- 895 [18] S. Hölldobler, C. D. P. Kencana Ramli, Logics and networks for human reasoning, in: *Artificial Neural Networks - ICANN 2009*, Vol. 5769 of *Lecture Notes in Computer Science*, Springer, 2009, pp. 85–94.

- 900 [19] E.-A. Dietz, S. Hölldobler, M. Ragni, A computational logic approach to the suppression task, in: Proceedings of the 34th Annual Conference of the Cognitive Science Society, 2012, pp. 1500–1505.
- [20] R. M. Byrne, Suppressing valid inferences with conditionals, *Cognition* 31 (1) (1989) 61–83.
- 905 [21] D. E. Rumelhart, G. E. Hinton, R. J. Williams, Learning internal representations by error propagation, in: *Parallel Distributed Processing Vol. 1: Foundations*, MIT Press, 1986, pp. 318–362.
- [22] F. Harder, T. R. Besold, An approach to supervised learning of three-valued lukasiewicz logic in hölldobler’s core method, in: Proceedings of the 4th International Workshop on Artificial Intelligence and Cognition (AIC) 2016, 2016, in
910 press.
- [23] J. Łukasiewicz, On three-valued logic, *The Polish Review* (1968) 43–44.
- [24] S. C. Kleene, *Introduction to Metamathematics*, North Holland, 1952.
- [25] M. Fitting, A kripke-kleene semantics for logic programs, *Journal of Logic Programming* 2 (1985) 295–312.
915
- [26] K. Stenning, M. Van Lambalgen, *Human reasoning and cognitive science*, MIT Press, 2008.
- [27] T. Mitchell, *Machine Learning*, McGraw Hill, 1997.
- [28] A. d’Avila Garcez, G. Zaverucha, The connectionist inductive learning and logic
920 programming system, *Applied Intelligence* 11 (1999) 59–77.
- [29] A. d’Avila Garcez, K. Broda, D. Gabbay, Symbolic knowledge extraction from trained neural networks: A sound approach, *Artificial Intelligence* 125 (2001) 155–207.
- [30] M. Mitchell, *An Introduction to Genetic Algorithms*, MIT Press, 1996.

925 **Appendix A: Proofs**

Backpropagation for logistic cost function with unipolar weights

To follow the notation used by [27], a few terms have to be introduced.

Logistic cost function: $J(\vec{w}) = \frac{1}{m} \sum_{i=1}^m [(y^{(i)} \log h_{\vec{w}}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\vec{w}}(x^{(i)}))]$

Error for one sample d : $E_d(\vec{w}) = \sum_{k \in \text{outputs}} [(t_k \log o_k) + (1 - t_k) \log(1 - o_k)]$

930 **Weight update term:** $\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial w_{ji}}$

Weighted sum of inputs to j : $net_j = x_0 \cdot w_0 + \sum_{k>0} x_{jk} \cdot \frac{1}{2} (w_{jk})^2$

Sigmoid function: $\sigma(x) = 1 / (1 + e^{-x})$

Error term delta: $\delta_j = -\frac{\partial E_d}{\partial net_j}$

The objective of the backpropagation algorithm is to produce weight updates Δw_{ji} for all weights for connections from i to j in the network. The learning rate η will be set independently and therefore the following proof is mostly concerned with a derivation of $\frac{\partial E_d}{\partial w_{ji}}$ for each weight.

Because the weight w_{ji} always enters into E_d in the context of net_j , using the chain rule, the partial derivative can be rewritten as follows:

$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial net_j} \frac{\partial net_j}{\partial w_{ji}}$$

The second factor in the right term yields:

$$\frac{\partial net_j}{\partial w_{ji}} = \begin{cases} i = 0 : x_{ji} \\ i > 0 : w_{ji} \cdot x_{ji} \end{cases}$$

Using above definition, the formula for weight updates follows.

$$\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial w_{ji}} = \begin{cases} i = 0 : \eta \cdot \delta_j \cdot x_{ji} \\ i > 0 : \eta \cdot \delta_j \cdot w_{ji} \cdot x_{ji} \end{cases}$$

940 The way in which δ_j is derived depends on the layer, in which neuron j is located.

δ_j for output layer units

As net_j only affects output o_j , another expansion yields

$$\frac{\partial E_d}{\partial net_j} = \frac{\partial E_d}{\partial o_j} \frac{\partial o_j}{\partial net_j}$$

The two resulting derivatives can be simplified.

$$\frac{\partial E_d}{\partial o_j} = \frac{t_j - o_j}{o_j \cdot (1 - o_j)}$$

$$\frac{\partial o_j}{\partial net_j} = \frac{\partial \sigma(net_j)}{\partial net_j} = o_j \cdot (1 - o_j)$$

This yields the result.

$$\delta_j = -\frac{\partial E_d}{\partial net_j} = -(t_j - o_j)$$

⁹⁴⁵ δ_j for hidden layer units

net_j affects E_d only through first o_j and then the weighted inputs net_k of all neurons that j has outgoing connections to, i.e. that are *downstream* from j (noted as $k \in DS(j)$ below). With this knowledge the following expansions can be made.

$$\frac{\partial E_d}{\partial net_j} = \sum_{k \in DS(j)} \frac{\partial E_d}{\partial net_k} \cdot \frac{\partial net_k}{\partial net_j} = \sum_{k \in DS(j)} -\delta_k \cdot \frac{\partial net_k}{\partial o_j} \cdot \frac{\partial o_j}{\partial net_j}$$

Both partial derivatives in the right term can be transformed. It helps to know that the ⁹⁵⁰ bias-unit has no incoming connection and thus $j > 0$.

$$\frac{\partial net_k}{\partial o_j} = \frac{\partial net_k}{\partial x_{kj}} = \frac{\partial}{\partial x_{kj}} (w_{k0} \cdot x_{k0} + \sum_{j>0} \frac{1}{2} (w_{kj})^2 \cdot x_{kj}) \stackrel{j>0}{=} \frac{1}{2} (w_{kj})^2$$

$$\frac{\partial o_j}{\partial net_j} = \frac{\sigma(net_j)}{\partial net_j} = \sigma' = o_j \cdot (1 - o_j)$$

The result follows.

$$\delta_j = o_j \cdot (1 - o_j) \cdot \sum_{k \in DS(j)} \delta_k \cdot \frac{1}{2} (w_{kj})^2$$

Preservation of Łukasiewicz semantics in sigmoidal cores

Let $[0, A^-]$ and $[A^+, 1]$ be two disjoint intervals representing firing and non-firing activation respectively. Assume $sig(z) = 1/(1 + e^{-z})$, $\omega > 0$ and $deg \geq 1$, the last
955 being the maximum indegree among neurons of the hidden/output layer.

In the following notation, activations A indicate whether they represent firing or not firing ($+$, $-$), are in hidden or output layer (h, o) and associate with a *true* or *false* neuron (\top, \perp), where needed. Note that activations in the input layer are discrete and therefore written as 0 and 1.

960 Proper functioning of logic gates in the network is guaranteed, as long as the following inequalities hold. For each layer and logic gate type there is one formula specifying the minimal input which must lead to activation and one formula for the maximal possible input that must not activate the gate.

- Logic gates in the hidden layer:

965 – On \top clause neurons ('and'-gate)

$$\begin{aligned} \min(A_{h\top}^+) &= sig(deg \cdot 1 \cdot \omega - (deg \cdot \omega - \frac{\omega}{2})) \geq A^+ \\ \max(A_{h\top}^-) &= sig((deg - 1) \cdot 1 \cdot \omega + 0 \cdot \omega - (deg \cdot \omega - \frac{\omega}{2})) \leq A^- \end{aligned}$$

– On \perp clause neurons ('or'-gate)

970

$$\begin{aligned} \min(A_{h\perp}^+) &= sig(1 \cdot \omega - \frac{\omega}{2}) \geq A^+ \\ \max(A_{h\perp}^-) &= sig(0 \cdot \omega - \frac{\omega}{2}) \leq A^- \end{aligned}$$

- Logic gates in the output layer:

– On \top clause neurons ('or'-gate)

$$\begin{aligned} \min(A_{o\top}^+) &= sig(\min(A_h^+) \cdot \omega - \frac{\omega}{2}) \geq A^+ \\ \max(A_{o\top}^-) &= sig(deg \cdot \max(A_h^-) - \frac{\omega}{2}) \leq A^- \end{aligned}$$

975 – On \perp clause neurons ('and'-gate)

$$\begin{aligned} \min(A_{o\perp}^+) &= sig(deg \cdot \min(A_h^+) \cdot \omega - (deg \cdot \omega - \frac{\omega}{2})) \geq A^+ \\ \max(A_{o\perp}^-) &= sig((deg - 1) \cdot 1 \cdot \omega + \max(A_h^-) \cdot \omega - (deg \cdot \omega - \frac{\omega}{2})) \leq A^- \end{aligned}$$

The hidden layer formulas reduce to:

$$\min(A_{h\top}^+) = sig(\frac{\omega}{2}) \geq A^+$$

$$\begin{aligned}
980 \quad \max(A_{h\top}^-) &= \text{sig}\left(-\frac{\omega}{2}\right) \leq A^- \\
\min(A_{h\perp}^+) &= \text{sig}\left(\frac{\omega}{2}\right) \geq A^+ \\
\max(A_{h\perp}^-) &= \text{sig}\left(-\frac{\omega}{2}\right) \leq A^-
\end{aligned}$$

The distinction between $h\top$ and $h\perp$ can thus be ignored in the following. Also it is clear now that $A^+ > \frac{1}{2} > A^-$. Inequalities regarding the output layer can be transformed into:

$$\begin{aligned}
985 \quad \mathbf{1} \quad \min(A_{o\top}^+) &= \text{sig}\left(\left(\min(A_h^+) - \frac{1}{2}\right) \cdot \omega\right) \geq A^+ \\
\mathbf{2} \quad \max(A_{o\top}^-) &= \text{sig}\left(\left(\text{deg} \cdot \max(A_h^-) - \frac{1}{2}\right) \cdot \omega\right) \leq A^- \\
\mathbf{3} \quad \min(A_{o\perp}^+) &= \text{sig}\left(\left(\text{deg} \cdot \min(A_h^+) - \frac{2\text{deg}-1}{2}\right) \cdot \omega\right) \geq A^+ \\
\mathbf{4} \quad \max(A_{o\perp}^-) &= \text{sig}\left(\left(\max(A_h^-) - \frac{1}{2}\right) \cdot \omega\right) \leq A^-
\end{aligned}$$

It can be seen that formulas 1 and 4 imply the hidden layer inequalities, 2 implies 4 and 3 implies 1.¹² Satisfying formulas 2 and 3 therefore satisfies the other inequalities as well. Now, as it has been established that:

$$\begin{aligned}
\min(A_h^+) &= \text{sig}\left(\frac{\omega}{2}\right) \\
\max(A_h^-) &= \text{sig}\left(-\frac{\omega}{2}\right)
\end{aligned}$$

In the 3 layer network it follows:

$$\begin{aligned}
995 \quad \max(A_{o\top}^-) &= \text{sig}\left(\left(\text{deg} \cdot \text{sig}\left(-\frac{\omega}{2}\right) - \frac{1}{2}\right) \cdot \omega\right) \\
\min(A_{o\perp}^+) &= \text{sig}\left(\left(\text{deg} \cdot \text{sig}\left(\frac{\omega}{2}\right) - \frac{2\text{deg}-1}{2}\right) \cdot \omega\right)
\end{aligned}$$

This results in the final inequalities

$$\begin{aligned}
\text{sig}\left(\left(\text{deg} \cdot \text{sig}\left(-\frac{\omega}{2}\right) - \frac{1}{2}\right) \cdot \omega\right) &\leq A^- \\
\text{sig}\left(\left(\text{deg} \cdot \text{sig}\left(\frac{\omega}{2}\right) - \frac{2\text{deg}-1}{2}\right) \cdot \omega\right) &\geq A^+
\end{aligned}$$

1000 Knowing that $A^+ > \frac{1}{2} > A^-$, these have a solution precisely iff

$$\begin{aligned}
\left(\text{deg} \cdot \text{sig}\left(-\frac{\omega}{2}\right) - \frac{1}{2}\right) &< 0 \\
\left(\text{deg} \cdot \text{sig}\left(\frac{\omega}{2}\right) - \frac{2\text{deg}-1}{2}\right) &> 0
\end{aligned}$$

Either is true iff $\omega > 2 \log(2\text{deg} - 1)$.

¹² $(\text{deg} \cdot \min(A_h^+) - \frac{2\text{deg}-1}{2})$ is maximal with $\text{deg} = 1$, then equal to $(\min(A_h^+) - \frac{1}{2})$

Relation of C-models and L-models*

1005 **Claim:** If an initial activation under C*-interpretation leads to a fixed point with no contradiction, the same fixed point will be reached under L-interpretation. Assuming the core architecture implementation itself is correct, this means that all least C*-models are also least L-models.

1010 **Proof sketch:**

- (1) By design, the one thing differentiating C*-consequence from L-consequence is that in L contradictions are resolved at the end of each iteration.
- (2) Due to monotonicity, every output neuron activated while finding a C* fixed point stays active in further iterations.
- 1015 (3) It follows from (1), that if the fixed point of C* and L on the same input differ, it's because a contradiction was resolved in L.
- (4) It follows from (2), that if a contradiction occurs in C during fixed point generation, it will still be there in the fixed point.
- (5) (3) and (4) imply, that if there is no contradiction in the C* fixed point, none was
1020 resolved in the corresponding L fixed point and thus, the fixed points are equal.

Appendix B: Extraction algorithm

Type definitions

Type **Vertex**

```
lbub      /* true if lower boundary, false if upper boundary */
target    /* index of target output neuron */
isRule    /* true if rule or subsumed by a rule */
isPos     /* true if target is active given activation in val */
decID     /* decimal ID */
val       /* array of input values */
mem       /* memorizes, which successors should be generated */
```

Algorithm 1: Vertex

Type **RuleSet**

```
lbRules   /* list of found minimal activating inputs */
ubRules   /* list of found maximal non-activating inputs */
```

Algorithm 2: RuleSet

Auxiliary functions

```
Function get_id      /* produces key to identify vertices */
```

Input: val

Output: decimal value of val read as a binary number

Algorithm 3: get_id

Function **succ_id** /* computes decID that differs from vertex at
index */

Input: vertex, index

if *vertex.mem[index] = 0* **then**

| newID := -1

else

max := length(vertex.val)

newID := vertex.decID

newID := newID + $2^{(max-index)}$ **if** *vertex.lbub*, newID - $2^{(max-index)}$

otherwise

Output: newID

Algorithm 4: succ_id

Function **make_succ** /* create successor from copy of vertex */

Input: vertex, index

newV := deepcopy(vertex)

newV.val[index] := 1 **if** *vertex.lbub*, 0 **otherwise**

newV.mem := abs(vertex.mem) /* 'forget' temporary -1 blockings */

newV.mem[index] := 0

newV.decID = succ_id(vertex, index)

Output: newV

Algorithm 5: make_succ

Function **in_rules_lb** /* true if vertex subsumed by rule in lower bound */

Input: vertex, ruleset

for $i = 1:\text{length}(\text{ruleset.lbRules})$ **do**

if $\text{minimum}(\text{vertex.val} - \text{ruleset.lbRules}[i]) == 0$ **then**

 isSub = true

 subRule = ruleset.lbRules[i] /* also returns subsuming rule */

Output: isSub, subRule

Output: false, zeros

Function **in_rules_ub** /* analogous function for upper bound */

Algorithm 6: in_rules_lb, in_rules_ub

Central algorithm

```

Function test          /* finds rules and subsumed anti-rules */
Input: core, vertex, ruleset
cOut := run_core(core, vertex.val)          /* get core output */
vertex.isPos := true if cOut[vertex.target] = 1, false otherwise
if vertex.lub then
  if vertex.isPos then
    ruleset.lbRules ← vertex.val
    vertex.mem := zeros
    vertex.isRule := true
    /* lower bound rule found */
  else
    isSub, subRule := in_rules_ub(vertex, ruleset)
    if isSub then set vertex.mem to -1 where subRule = 1 and mem = 1
    /* isSub: subsumed by lower bound anti-rule */
else
  if vertex.isPos then
    isSub, subRule = in_rules_lb(vertex, ruleset)
    if isSub then set vertex.mem to -1 where subRule = 0 and mem = 1
    /* isSub: subsumed by upper bound anti-rule */
  else
    ruleset.ubRules ← vertex.val
    vertex.mem := zeros
    vertex.isRule := true
    /* upper bound rule found */

```

Algorithm 7: test

Function **successors**

Input: core, vertex, vvector, queue

test(core, vertex, ruleset)

noMatchSucc := true **if** *minimum(vertex.mem) ≥ 0*, false **otherwise**

for $i = 1:\text{length}(\text{vertex.mem})$ **do**

```
    if vertex.mem[i] = 1 then           /* traverse valid successors */
    |   succID := succ_id(vertex,i)
    |   if ¬isdefined(vvector, succID) then /* create if necessary */
    |   |   vvector[succID] := make_succ(vertex, i)
    |   |   test(core, vvector[succID], ruleset)
    |   |   queue ← succID
    |   if vvector[succID].isRule then vertex.mem[i] := 0 /* pool successor
    |   |   info */
    |   if vvector[succID].isPos = vertex.isPos then noMatchSucc := false
```

if *noMatchSucc* **then** /* test for anti-rules */

```
    |   if vertex.lbub ∧ ¬vertex.isPos then ruleset.ubRules ← vertex.val
    |   if ¬vertex.lbub ∧ vertex.isPos then ruleset.lbRules ← vertex.val
```

for $i = 1:\text{length}(\text{vertex.mem})$ **do** /* distribute info to all successors

*/

```
    |   if vertex.mem[i] = 1 then
    |   |   succID := succ_id(vertex, i)
    |   |   set vvector[succID].mem to 0 where vertex.mem is 0
```

Algorithm 8: successors

Function rules_for_target

Input: core, target

\perp := vertex for bottom element for core

\top := vertex for top element for core

queue := Queue(integers) /* stores keys of vertices such that */

queue \leftarrow \perp .decID, \top .decID /* the boundaries alternate by layer */

vvector := vector(vertices) with 3^n entries, where n = number of output units in core

vvector[begin] := \perp

vvector[end] := \top

ruleset := empty RuleSet

while \neg empty(queue) **do** /* repeat until all vertices pruned or covered */

```
┌   index := dequeue(queue)
├   v := vvector[index]
└   successors(core, v, ruleset, vvector, queue)
```

Algorithm 9: rules_for_target