

# Neural Network Models of Conditionals: An Introduction

Hannes Leitgeb

Department of Philosophy, University of Bristol,  
Hannes.Leitgeb@bristol.ac.uk

## Abstract

This “lecture notes style” article gives a brief survey of neural network models of conditionals. After short introductions into the studies of neural networks and conditionals, we turn to the notion of an *interpreted dynamical system* as a unifying concept in the logical investigation of dynamic systems in general, and of neural networks in particular. We explain how conditionals get represented by interpreted dynamical systems, which logical systems these conditionals obey, and what the main open problems in this area are.

Keywords: Conditionals. Neural networks. Dynamical systems. Nonmonotonic logic.

## 1 Introduction

Neural networks are abstract models of brain structures capable of adapting to new information. The learning abilities of artificial neural networks have given rise to successful computer implementations of various cognitive tasks, from the recognition of facial images to the prediction of currency movement.

Logic deals with formal systems of reasoning; in particular, inductive logic studies formal systems of reasoning towards plausible but uncertain conclusions. As evidence accumulates, the degree to which it supports a hypothesis, as measured by the logic, should tend to indicate that the hypothesis is likely to be true.

Although sharing a joint focus on information and reasoning, until recently these two areas developed in opposition to each other: neural networks are quantitative dynamic systems, while logical reasoners must be symbolic systems; networks are described by mathematical equations, whereas logic is subject to normative statements about how we ought to reason; neural networks have been studied by scientists, whilst the “problem of induction” is regarded as belonging to the humanities.

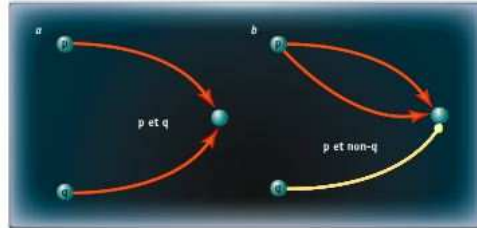
At present this assessment is changing: the emergence of logical formalisms for uncertain reasoning and the discovery that these formalisms apply to neural net processes on the representational level give rise to the expectation that the dynamics of artificial neural networks can be understood in terms of logically valid, and thus rational, rules of inference. As neural networks, commonsense reasoning, and scientific

induction seem to conform to similar logical systems, a joint theoretical framework is in the offing that might lead to new insights into the logical and cognitive basis of both everyday reasoning and science.

## 2 Neural Networks as Models of Reasoning

In their famous article “A Logical Calculus of the Ideas Immanent in Nervous Activity”, McCulloch and Pitts (1943) first introduced artificial neural networks as mathematical abstractions from neural circuits in the brain. A McCulloch-Pitts network consists of a set of nodes and a set of connections between these nodes. Each node can be in one of two possible states: it fires (1), or it does not (0). Each connection is of one of two possible kinds: along inhibitory connections, nodes receive inhibitory signals by which they get deactivated at the next point of time (on a discrete time scale). Via excitatory connections, signals are transferred from one node to another which have a stimulating effect on the target node: if the node does not get inhibited, and if the number of all incoming excitatory signals exceed or are identical to some fixed threshold value that is associated with the node, then the node fires at the next point of time. Despite appearing to be quite simple devices, McCulloch and Pitts were able to prove that in principle every finite automaton can be realized by such a McCulloch-Pitts network. Furthermore, the state transitions which take place in such networks allow for a description in logical terms: if the activity of a node is considered as a truth value, then the node itself may be regarded as an entity which *has* a truth value, i.e., as a formula or proposition. If the “truth value” of a node does not depend on the “truth values” of other nodes (but, say, only some given input), then it is indeed natural to regard such nodes as *atomic* formulas or propositions. Accordingly, if nodes are put together in a network, such that connections between nodes can cause the “truth values” of other nodes to be altered, then the latter nodes may be taken to correspond to *complex* formulas; the semantic dependency of the truth value of a complex formula on the truth values of its component formulas is thus represented by the network topology and the choice of thresholds.

As an example, consider the following two very elementary McCulloch-Pitts networks: In the first network, excitatory connections lead from nodes  $p$  and  $q$  to a third node. If this latter node has a threshold value of 2, then the node is going to fire if and only if both  $p$  and  $q$  were active at the previous point of time. So we can associate the formula  $p \wedge q$  with this node. In the second network, two excitatory lines lead from  $p$  to the output node, whereas  $q$  is connected to the latter by an inhibitory edge. If e.g. the output node has a threshold of 2, it will be activated at the next point of time if and only if  $p$  is set to 1 and  $q$  is set to 0 (and therefore does not have any inhibitory influence). Hence, the third node in the network corresponds to the formula  $p \wedge \neg q$ . Here is a picture<sup>1</sup>:



This way of associating nodes in networks with formulas in the language of classical propositional logic extends to more interesting networks with multiple layers of nodes and with more complex patterns of excitatory and inhibitory connections. E.g., it would be easy to extend the second network in fig. 1 by a node that represents  $\neg(p \wedge \neg q)$ , i.e., a formula which is logically equivalent to the material conditional  $p \supset q$ . If our brains were, at least on some level, similar to neural networks of the McCulloch-Pitts kind, they could thus be understood as collections of simple logical units put together in order to calculate binary truth values from external or internal input. The calculation of the truth values of material conditionals would be a special case of this form of computational processing.

Of course, the McCulloch-Pitts networks are, in several respects, much too simple to be plausible models of actual neural networks in animal or human brains. In particular, they are not yet able to learn. The next decisive step in the development of artificial neural networks was to introduce variable weights which are attached to connections and which encode the degree of influence that nodes can exert on their target nodes via these connections. By sophisticated learning algorithms, these weights can be adjusted in order to map inputs to their intended outputs, e.g., facial images of persons to the names of these persons, or verbs to their correct past tenses. Despite some initial success in the 1950s and 1960s – mainly associated with Frank Rosenblatt's *Perceptrons* – it was only in the 1980s that artificial neural network models of cognition became serious contenders to the dominant symbolic computation paradigm in artificial intelligence. (Rumelhart et al. 1986 is still something like the “bible” of connectionism; Rojas 1993 is a nice introduction to neural networks – have a look at these two for more background information.) As we will explain below, the more recent neural network models do not only differ from the original McCulloch-Pitts networks in terms of complexity and learning abilities, they also differ in terms of the interpretation of their components: instead of assigning meaning – expressed by formulas – to *single* nodes, the modern approach emphasizes that it is rather *patterns* or *sets* of nodes which receive an interpretation.

How does ‘cognition by neural networks’ relate to the traditional ‘cognition by symbolic computation’ paradigm of cognitive science (exemplified by classic Artificial Intelligence)? According to the latter, (i) intelligent cognition demands structurally complex mental representations, such that (ii) cognitive processing is only sensitive to the form of these representations, (iii) cognitive processing conforms to rules, stable over the representations themselves and articulable in the format of a computer program, (iv) (standard) mental representations have syntactic structure with a compo-

sitional semantics, and (iv) cognitive transitions conform to a computable cognitive-transition function (we adopt this characterization essentially from Horgan&Tienson 1996, with slight deviations). Intelligent cognition is supposed to be “systematic” and “productive” (see Fodor&Pylyshyn 1988), i.e., the representational capacities of intelligent agents are supposed to be necessarily closed under various representation-transforming and representation-generating operations (e.g., if an agent is able to represent that  $aRb$ , it is also able to represent that  $bRa$ , etc.). This capacity is hypothesized to be due to the combinatorial properties of languages of mental symbols based on a recursive grammar. A cognitive agent that conforms to the symbolic computation paradigm has the belief that  $\varphi$  if and only if a corresponding sentence  $\varphi$  is stored in the agent’s symbolic knowledge base. The rules that govern cognitive processes according to the symbolic computation paradigm are either represented within the cognitive agent as symbolic entities themselves, or they are hard-wired. Inference processes are taken to be internalizations of derivation steps within some logical system, and the alleged “systematicity” of inferences (see again Fodor&Pylyshyn 1988) is explained by the internal representation or hard-wiring of rules which are only sensitive to the syntactic form of sentential representations.

Cognition by artificial neural networks, on the other hand, belongs to the so-called dynamical systems paradigm of cognitive science which can be summarized by what van Gelder (1998) calls the “dynamical hypothesis”: “for every kind of cognitive performance exhibited by a natural cognitive agent, there is some quantitative [dynamical] system instantiated by the agent at the highest relevant level of causal organization [i.e., at the level of representations], so that performances of that kind are behaviors of that system” (van Gelder 1998, p.622). A dynamical system may be regarded as a pair of a state space and a set of trajectories, such that each point of the space corresponds to a total cognitive state of the system, and every point of the space lies precisely on one trajectory. If a certain point corresponds to the system’s total cognitive state at some time, the further evolution of the system follows the trajectory emanating at this point. Usually, such systems are either defined by differential equations, or by difference equations, defined over the points of the state space: in the first case one speaks of continuous dynamical systems with continuous time, while in the latter case one speaks of discrete dynamical systems with discrete time. In the discrete case, the set of trajectories may be replaced by a state-transition mapping, such that each trajectory is generated by the iterated application of the mapping. A *cognitive* dynamical system is a dynamical system with representations, i.e., where states and state transitions can be ascribed a content or an interpretation. The dynamic systems paradigm assumes that intelligent cognition takes place in the form of state-transitions in quantitative systems, i.e., systems in which a metric structure is associated with the points of the state space, and where the dynamics of the system is systematically related to the distances measured by the metric function. The distances between points may be regarded as a measure of their similarity *qua* total cognitive states. Moreover, the typical dynamical systems that are studied within the dynamical systems paradigm also have a vector space structure, and thus they “support a geometric perspective on system behaviour” (van Gelder 1998, p.619).

Connectionism is the most important movement within the dynamical systems paradigm: Artificial neural networks are the dynamical systems the connectionists

are interested in. Smolensky (1988) characterizes connectionism by the following hypotheses: (i) “The connectionist dynamical system hypothesis: The state of the intuitive processor at any moment is precisely defined by a vector of numerical values (one for each unit). The dynamics of the intuitive processor are governed by a differential equation. The numerical parameters in this equation constitute the processor’s program or knowledge. In learning systems, these parameters change according to another differential equation.” (ii) “The subconceptual unit hypothesis: The entities in the intuitive processor with the semantics of conscious concepts of the task domain are complex patterns of activity over many units. Each unit participates in many such patterns.” (iii) “The subconceptual level hypothesis: Complete, formal, and precise descriptions of the intuitive processor are generally tractable not at the conceptual level, but only at the subconceptual level.” The subconceptual level is the level of analysis that is preferred by the connectionist paradigm, or, as Smolensky expresses it, the *subsymbolic* paradigm; it lies “below” the conceptual level which is preferred by the symbolic computation paradigm, but “above” the neural level preferred by neuroscience.

(i) proves connectionism to belong to the dynamical systems paradigm. The subconceptual unit hypothesis (ii) and the subconceptual level hypothesis (iii) highlight the main differences between the old McCulloch&Pitts approach presented above and modern day connectionism: by (ii), single nodes or single connections in a neural network are normally not supposed to carry any meaning at all; the representing units are distributed patterns of activation which involve a great number of nodes or even the network topology as a whole (see van Gelder 1999 on “Distributed versus local representation”). In more metaphorical terms: there is no “grandmother cell”, i.e., no single neuron which would correspond to a very complex formula which describes your grandmother and which would fire if and only if your grandmother were perceived, but rather your grandmother’s being perceived is represented by some complex pattern of activation which spreads throughout parts of the network at the time of perception. Furthermore, by (iii), if symbols can be attached to the activation patterns of nodes or to other “global” aspects neural networks at all, the transitions from one representing item – one pattern – to another will no longer be effected on the level of these representing items themselves but rather on the sub-symbolic level of nodes and edges. Therefore, it seems to be impossible to translate the computations on the sub-symbolic level into sequences of rules on the symbolic level, let alone into logical rules which apply to complex symbolic expressions. Thus, McCulloch&Pitts’ *logical* approach to neural networks has to be given up, or so it seems. Instead of analyzing cognition in terms of localized representations of formulas – “hard constraints” – Smolensky (1988), p.18, suggests that connectionist cognition proceeds by means of “soft constraints”: “Formalizing knowledge in soft constraints rather than hard rules has important consequences. Hard constraints have consequences singly; they are rules that can be applied separately and sequentially – the operation of each proceeding independently of whatever other rules may exist. But soft constraints have no implications singly; any one can be overridden by the others. It is only the entire set of soft constraints that has any implications. Inference must be a cooperative process [...] Furthermore, adding additional soft constraints can repeal conclusions that were formerly valid: Subsymbolic inference is fundamentally nonmonotonic.” If human reasoning is as connectionists describe it, then McCulloch&Pitts’ account of reasoning in terms

of neural network implementations of truth functions in classical logic can hardly be adequate.

Even if this very last statement about McCulloch&Pitts' theory is true, this does not yet entail that the symbolic computation paradigm and the dynamical systems paradigm themselves have to be completely mutually exclusive, i.e., significant aspects of the two paradigms could actually turn out to be compatible with each other. As Gärdenfors (1994), pp.67f, suggests, the two paradigms might in fact be complementing each other: "they are best viewed as two different perspectives that can be adopted when describing the activities of various computational devices." New results concerning symbol manipulation in networks (see e.g. Chalmers 1990, Chen&Honavar 1999) and on rule extraction from networks (see e.g. d' Avila Garcez et al. 2001, Hölldobler 1993, Hölldobler et al. 2004) show that there might be continuous paths of transition from the one paradigm to the other. Indeed, hybrid systems consisting of both symbolic and network components have been suggested (see e.g. Legendre et al. 1994). Finally, the analysis of neural networks in terms of *logical laws and rules* has been pursued quite intensively in recent years, which is the topic of this overview article.

Here are some relevant references on logical accounts of neural network cognition (they can also be found in the bibliography – note that this is a *very* incomplete list though!):

- A.S. d'Avila Garcez, D.M. Gabbay, and L.C. Lamb (200?): *Connectionist Non-Classical Logics*, to appear with Springer-Verlag.
- Balkenius, C. and P. Gärdenfors (1991): "Nonmonotonic inferences in neural networks", in: J. Allen, R. Fikes, and E. Sandewall (eds.), *Principles of Knowledge Representation and Reasoning*, San Mateo: Morgan Kaufmann, 32–39.
- Blutner, R. (2004): "Nonmonotonic inferences and neural networks", *Synthese* 142, 143–74.
- Leitgeb, H. (2001): "Nonmonotonic reasoning by inhibition nets", *Artificial Intelligence* 128, 161–201.
- Leitgeb, H. (2004): *Inference on the Low Level. An Investigation into Deduction, Nonmonotonic Reasoning, and the Philosophy of Cognition*, Dordrecht: Kluwer, Applied Logic Series.
- Leitgeb, H. (2005a): "Interpreted dynamical systems and qualitative laws: From inhibition networks to evolutionary systems", *Synthese* 146, 189–202.

(For these lecture notes, material contained in Leitgeb 2004, 2005a, and the popular and non-technical exposition of logic and neural networks in Leitgeb 2005b was used; the approach in Leitgeb's articles is greatly inspired by Balkenius and Gärdenfors 1991.)

The main idea behind these theories is that if classical logic is replaced by a different logical calculus – perhaps some symbolically encoded system of probabilistic or nonmonotonic reasoning that is closer to the commonsense reasoning that our brains are usually involved in – then a logical description or characterization of neural network states and processes might be possible. Some of the authors listed above indeed aim at expressing in terms of logical formulas or rules what is going in a neural

network on the level of *distributed* representation. We will present one of these approaches in section 4. If some of these logical accounts of neural network cognition were ultimately successful, then the gap between the dynamic systems paradigm and the symbolic computation paradigm in cognitive science would be bridged, or at least shortened significantly. This would also constitute an important step in understanding what neural networks actually do; otherwise, we might be stuck with an ingenious technical machinery which maps an input to its desired output but where the processes which lead from the one to the other remain uninterpreted and unexplained. Progress on logical account of neural networks might also lead to new insights in uncertain reasoning, induction, and even the philosophy of science – we will return to this in the final open questions section of this article.

### 3 Conditionals: Natural Language and Logical Reconstruction

Conditionals are sentences of an ‘if... then...’ form; alternatively, conditionals are defined to be the propositions that are expressed by such sentences. So, the logical form of a conditional is an expression of the form

If  $\varphi$ , then  $\psi$

or, more formalized,

$\varphi \Rightarrow \psi$

where  $\varphi$  is called the ‘antecedent’ of the conditional and  $\psi$  its ‘consequent’; both the antecedent and the consequent of a conditional are sentences.

Conditionals are crucial in everyday communication, especially when we want to convey an information that goes beyond the currently present perceptual situation. Conditionals also play a major role in philosophical theories about dispositions, causality, laws, time, conditional norms, probability, belief, belief revision, and so forth. Finally, conditionals are related closely to quantifiers, such as ‘All  $\varphi$  are  $\psi$ ’, ‘There are  $\varphi$  which are  $\psi$ ’, ‘Most  $\varphi$  are  $\psi$ ’, etc. (see van Benthem 1984 for a nice discussion of this relationship; more can be found by consulting the theory of *generalized quantifiers* – see e.g. Peters&Westerstahl 2006). But note that in these latter cases, ‘ $\varphi$ ’ and ‘ $\psi$ ’ are place holders for *open formulas* – formulas with a free variable – rather than sentences. This is sometimes overlooked even when specialists discuss these topics:

**Remark 1 (DIGRESSION)** *In computer science, in particular in the literature on nonmonotonic reasoning, the following two sets of locutions are often not distinguished properly: on the one hand,*

- *if  $\varphi$  then normally  $\psi$*
- *if  $\varphi$  then typically  $\psi$*
- *if  $\varphi$  then it is very likely that  $\psi$*

and, on the other,

- normal  $\varphi$  are  $\psi$
- typical  $\varphi$  are  $\psi$
- by far most of the  $\varphi$  are  $\psi$

In the first set, ' $\varphi$ ' and ' $\psi$ ' are to be replaced by sentences such as 'Tweety is a bird' and 'Tweety is able to fly', whereas in the second set ' $\varphi$ ' and ' $\psi$ ' are to be substituted by generics such as 'birds' and 'flyers' (or, in a more formalized context, by open formulas such as ' $x$  is a bird' and ' $x$  is able to fly'). Accordingly, if a member of the first set e.g. expresses something probabilistic, then the probability measure in question should be a subjective probability measure by which rational degrees of belief are attributed to propositions. However, in the case of the members of the second set, the corresponding probability measure should be a statistical one by which (limit) percentages are attributed to properties. End of DIGRESSION.

Among conditionals in natural language, usually the following distinction is made<sup>2</sup>:

1. If Oswald had not killed Kennedy, then someone else would have.
2. If Oswald did not kill Kennedy, then someone else did.

2 is accepted by almost everyone, whilst we do not seem to know whether 1 is true. This invites the following classification: A conditional such as 2 is traditionally called *indicative*. A conditional like 1 is called *subjunctive*. In conversation, the antecedents of subjunctive conditionals are often assumed or presupposed to be false: in such cases, one speaks of these subjunctive conditionals as *counterfactuals*. Subjunctive and indicative conditionals may have the same antecedents and consequents while differing only in their conditional connectives, i.e., their 'if'-'then' occurrences have different meanings.

The classification into indicative and subjunctive conditionals constitutes a philosophical problem in itself, but roughly one proceeds by the following rules of thumb:

- *Subjunctive*:
  - Semantically: represents a denoted act or state not as fact but as contingent or possible ("subjunctive mood").
  - Syntactically: is of a 'had-would' or, in any case, of a '...-would' form.
- *Indicative*:
  - Semantically: represents the denoted act or state as an objective fact ("indicative mood").
  - Syntactically: is of a 'did-did' or a 'does-will' form.

(But note there are always exceptions in natural language!)

When logic developed into a serious philosophical and mathematical discipline in the late 19th and the early 20th century, logicians quickly came up with two suggestions of how to formalize conditionals, whether indicative or subjunctive:



- $A \supset B$ : Formalization by means of material conditionals (material implications).
- $A \rightarrow B$ : Formalization by means of strict conditionals (strict implications).

From an axiomatic point of view, the meaning of the former is given by any of the typical deductive systems for classical propositional logic. The logical systems for the latter were investigated intensively by C.I. Lewis, however it was only after the axiomatic systems of normal modal logic had been developed by S. Kripke that the analysis of  $A \rightarrow B$  in terms of  $\Box(A \supset B)$  emerged as a standard (where  $\Box$  is the necessity operator studied by modal logicians). On the semantic side, the meaning of  $\supset$  is given by its well-known truth table, whereas the semantics of  $\rightarrow$  can be stated on the basis of the usual Kripkean possible worlds semantics of  $\Box$ .

These formalizations of the ‘if... then...’ in classical logic proved to be enormously successful, especially in the formalization of mathematical theories and of fragments of empirical theories. However, there was still a problem: both  $\supset$  and  $\rightarrow$

are *monotonic*, i.e., the rule  $\frac{\varphi \Rightarrow \psi}{\varphi \wedge \rho \Rightarrow \psi}$  is logically valid if ‘ $\Rightarrow$ ’ is replaced by either

of the two connectives. On the other hand, there seem to be many instances of indicative and subjunctive conditionals in natural language which are *nonmonotonic*, i.e., for

which the rule  $\frac{\varphi \Rightarrow \psi}{\varphi \wedge \rho \Rightarrow \psi}$  should not assumed to be valid. E.g., ‘If it rains, I will

give you an umbrella’ does not seem to logically imply ‘If it rains and I am in prison, I will give you an umbrella’, nor does ‘If it rained, I would give you an umbrella’ seem to logically imply ‘If it rained and I were in prison, I would give you an umbrella’. Accordingly, add e.g. ‘...and Kennedy in fact survived all attacks on his life’ to the antecedent of ‘If Oswald did not kill Kennedy, then someone else did’ and the resulting conditional does not seem acceptable anymore. Therefore, philosophical logicians started to investigate new logical systems in which monotonicity (or *strengthening of the antecedent*) would not turn out to be logically valid. Lewis (1973) is the classic treatise on counterfactuals as nonmonotonic conditionals. Since the nonmonotonicity phenomenon was already well known in probability theory – a conditional probability  $P(Y|X)$  being high does not entail the conditional probability  $P(Y \cap Z|X)$  being high – it is not surprising that some of the modern accounts of conditionals also turned out to have a probabilistic semantics: indeed, Adams’ (1975) famously developed a probabilistic theory of indicative conditionals (see Adams 1998 for a more general overview of probability logic). For a good state-of-the-art overview on the philosophical literature on indicative and subjunctive conditionals, see Bennett (2003).

At the same time, or actually a little later (note: philosophy was *first!*), theoretical computer scientists made a similar discovery, though from a quite different point of view: The logical resources which had proved to be so successful for the description of the precise and unchanging world of mathematics, and which had continued to be useful for describing the closed universe of e.g. a chess-playing piece of software, turned out to be insufficient for the replication of commonsense reasoning. Assume you want to describe what happens to your car when you turn the ignition key: well, you might say, the car starts, so ‘if the ignition key is turned in my car, then the car starts’ seems to be the proper description of the situation. But what if the gas tank is empty? You better improve your description by saying ‘if the ignition key is turned

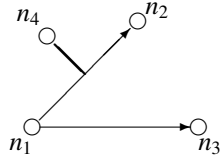
in my car and the gas tank is not empty, then the car starts’. However, this could still be contradicted by a potato that is clogging the tail pipe, or by a failure of the battery, or by an extra-terrestrial blocking your engine, or... The possible exceptions to ‘if the ignition key is turned in my car, then the car starts’ are countless, heterogeneous, and unclear. Nevertheless, we seem to be able to communicate information with such simple conditional sentences, and, which is equally important, we are able to reason with them in a rational way. In order to do so we make use of a little logical “artifice”: we do not really understand ‘if the ignition key is turned in my car, then the car starts’ as expressing that it is not the case that the ignition key is turned and the car does not start – after all, what is negated here might indeed be the case in exceptional circumstances – but rather that *normally*, or with a *high probability*, given the ignition key is turned, the car starts. Instead of trying to enumerate the indefinite class of exceptions in the if-part of a material or strict conditional, we tacitly or explicitly qualify ‘if the ignition key is turned in my car, then the car starts’ as holding only in normal or likely circumstances, whatever these circumstances may look like. As a consequence, the logic of such normality claims again differs from the logic of material or strict conditionals: ‘if Tweety is a bird, then [normally] Tweety is able to fly’ is, presumably, true, but ‘if Tweety is a penguin bird, then [normally] Tweety is able to fly’ is not, and neither is ‘if Tweety is a dead bird, then [normally] Tweety is able to fly’ or ‘if Tweety is a bird with his feet set in concrete, then [normally] Tweety is able to fly’. So computer scientists found themselves in need of describing reality in terms of *nonmonotonic* normality conditionals on the basis of which computers should be able to draw justified inferences about the everyday world while being unaffected by the omnipresence of exceptions. This is the subject of *nonmonotonic reasoning*, one of the most vibrant areas of theoretical computer science in the last 30 years. (See Brewka et al. 1997 for a very nice state-of-the-art overview, Makinson 1994 for a comprehensive logical account, and Schurz&Leitgeb 2005 for a compendium of articles on cognitive aspects of nonmonotonic reasoning. Ginsberg 1987 is outdated but still very useful if one wants to see what nonmonotonic reasoning derives from.)

In the next section we will suggest that so-called *interpreted dynamical systems* may be used to yield a semantics for nonmonotonic conditionals; the logical systems which turn out to be sound and complete with respect to such semantics are standard systems of conditional logic which have been studied both in philosophical logic and nonmonotonic reasoning. Interpreted artificial neural networks will be shown to be the paradigm case examples of interpreted dynamical systems. Although the conditionals that are satisfied by such interpreted artificial neural networks may be regarded to be represented distributedly by these networks, the logical rules they obey are the rules of systems which had been developed in order to make computers cope with the real world by means of symbolic computation, and which had been investigated even before by philosophers who intended to give a proper logical analysis of indicative and subjunctive conditionals. Since the dynamics of state changes in interpreted neural networks can be described correctly and completely by sets of conditionals which are closed under the rules of such logical systems, neural networks can be understood as nonmonotonic reasoners who, when they evolve under an input towards a state of minimal energy, draw conclusions that follow from premises in all minimally abnormal cases.

## 4 From Dynamical Systems to Conditionals: Interpreted Dynamical Systems

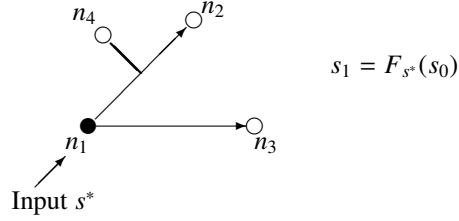
Following Gärdenfors' proposal mentioned above, we will study cognitive dynamical systems from two complementary perspectives. On the one hand, cognitive dynamical systems such as neural networks can be described in terms of differential or difference equations, i.e., as *dynamical systems*. On the other hand, they seem to exemplify cognitive states and processes which can be ascribed propositional contents which may in turn be expressed by sentences or formulas; so they are *cognitive agents* or *reasoners*.

Here is an example. For the sake of simplicity, let us forget about the weights again which are attached to the edges of a typical neural network, and let us also assume that the activation functions which are defined for the nodes in such a network are as straight-forward and simple as in the case of the McCulloch-Pitts networks. Then we might e.g. end up with a simple qualitative neural networks looking like this<sup>3</sup>:

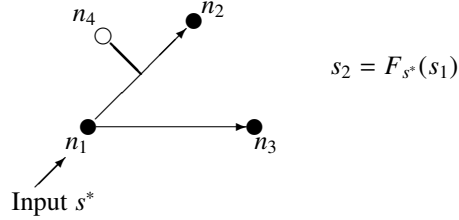


This is a network with four nodes.  $n_1$  is connected to  $n_2$  and  $n_3$  by excitatory connections. In contrast with traditional McCulloch-Pitts networks, there is also an inhibitory connection which leads from  $n_4$  to the *excitatory connection* from  $n_1$  to  $n_2$ . So, if  $n_4$  is active, then this is not going to directly inhibit the activation of some other node at the next point of time, but instead any activity by  $n_4$  will only have the effect that no excitatory stimulus will be able to pass the edge from  $n_1$  to  $n_2$  at the next point of time.

Now, say, node  $n_1$  gets activated by some external stimulus, e.g., by some sensory signal coming from outside of the network. We will assume that such inputs always remain constant for sufficiently long, so in our example one should think of  $n_1$  as being activated from the outside until the computational process that we are interested in has delivered its final output. Formally, we can describe what is going on in the following way: the network is in an initial state  $s_0$ , e.g., the state in which no node fires. This state  $s_0$  may be regarded as a mapping from the set of nodes into the set  $\{0, 1\}$ , such that each node is mapped to 0. Furthermore, the network is committed to an input  $s^*$  which makes  $n_0$  fire but which activates no other node: it is useful to identify such an input with the network state that the input would generate just by means of external influences on the network. Thus, in our case,  $s^*$  will be the state in which the node  $n_0$  is mapped to 1 and in which all other nodes are mapped to 0. The dynamics of the network can now be described in terms of a state transition mapping  $F_{s^*}$  which is given relative to the (constant) incoming input –  $s^*$  – and which is applied to the initial state  $s_0$  in order to determine the next state  $s_1$  of the network. Since no node is active in  $s_0$ , the only nodes which will be active in  $s_1$  will be those activated by means of the input itself, i.e.,  $n_1$ . So we have:

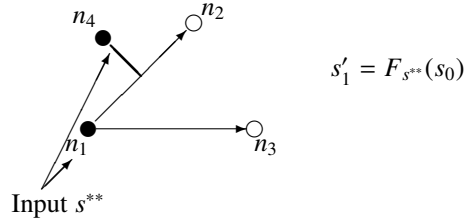


Accordingly, in order to determine the next state  $s_2$  of the network, the state transition mapping  $F_{s^*}$  is applied again. The state transition will be such that the activity of  $n_1$  in  $s_1$  spreads to  $n_2$  and  $n_3$ , which yields:

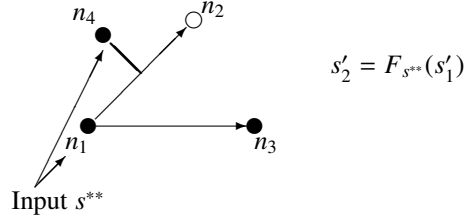


If the state transition mapping is applied again, then nothing is going to happen anymore (until the input to the network changes): hence,  $s_3 = F_{s^*}(s_2) = s_2$ . Connectionists regard such a *stable* or *equilibrium* state as a network's "answer" to the "question" posed by the input. So,  $s_2$  – the state in which only  $n_1, n_2, n_3$  fire – is the output that belongs to the input  $s^*$ . As we will also say,  $s_2$  is an  $s^*$ -stable state,

What would happen if we used a different input but the same initial state: Let  $s^{**}$  be the state in which both  $n_1$  and  $n_4$  fire, i.e., the external input now causes these two nodes to become active. Then we have, by the same token as before:



But now the state transition will be such that the activity of  $n_4$  in  $s_1$  blocks the excitation of  $n_2$  by  $n_1$ . In other words:



Once again, a stable state is reached after two computation cycles, and this time the output to the input state  $s^{**}$  is the state in which  $n_1, n_3, n_4$  fire, i.e.,  $s'_2$  is an  $s^{**}$ -stable state

This was a typical description of (simplified) network processes in the language of the theory of dynamical systems. Our goal is now to complement this description by one according to which cognitive dynamical systems have beliefs, draw inferences, and so forth. So if  $x$  is a neural network, we want to say things like

- $x$  believes that  $\neg\varphi$
- $x$  infers that  $\varphi \vee \psi$  from  $\varphi$
- $\vdots$

where  $\varphi$  and  $\psi$  are *sentences*.

Our task is thus to associate *states* of cognitive dynamical systems with *sentences* or *propositions*: the states of such dynamic systems ought to carry information that can be expressed linguistically.

Let us make this idea more precise. In order to do so, we first have to abstract from the overly simplified dynamical systems which are given by the qualitative neural networks sketched above. Indeed, we want to leave open at this point what our dynamical systems will be like – whether artificial neural networks or not – as long as they satisfy a few abstract requirements.

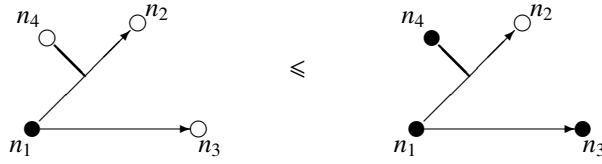
Here is what we will presuppose: We deal with discrete dynamical systems with a set  $S$  of states. On  $S$  a partial order<sup>4</sup>  $\leq$  is defined, which we will interpret as an ordering of the amount of information that is carried by states; so  $s \leq s'$  will be read as:  $s'$  carries at least as much information as  $s$  does. We will also assume that  $\leq$  is “nice” in so far as for every two states  $s$  and  $s'$  there is a uniquely determined state  $\sup(s, s')$  which (i) carries at least as much information as  $s$ , which (ii) carries at least as much information as  $s'$ , and which (iii) is the state with the least amount of information among all those states for which (i) and (ii) hold. Formally, such a state  $\sup(s, s')$  is the *supremum* of  $s$  and  $s'$  in the partial order  $\leq$ . Finally, an internal next-state function is defined for the dynamical system, where this next-state function is like the state transition mapping described above except that – for the moment – we will disregard possible inputs to the system. So in the examples above, an application of the corresponding next-state mapping e.g. would lead to the transmission of the activity of  $n_1$  to  $n_3$  once  $n_1$  gets activated, but it will never lead to any activation of  $n_1$  itself since  $n_1$  can only be activated by external input.

Summing up, we get what is called an ‘ordered discrete dynamical system’ in Leitgeb (2005):

**Definition 2**  $\mathcal{S} = \langle S, ns, \leq \rangle$  is an ordered discrete dynamical system :iff

1.  $S$  is a non-empty set (the set of states).
2.  $ns : S \rightarrow S$  (the internal next-state function).
3.  $\leq \subseteq S \times S$  is a partial order (the information ordering) on  $S$ ,  
such that for all  $s, s' \in S$  there is a supremum  $\sup(s, s') \in S$  with respect to  $\leq$ .

In the example networks above, we had  $S = \{s \mid s : N \rightarrow \{0, 1\}\}$  with  $N = \{n_1, n_2, n_3, n_4\}$  being the set of nodes. In order to define a suitable information ordering  $\leq$  on  $S$ , we can e.g. use the following idea: the more nodes are activated in a state, the more information the state carries. Then we would e.g. have:



If  $\leq$  is defined in this way, then  $\sup(s, s')$  turns out to be the union of the activation patterns that correspond to  $s$  and  $s'$ ; in such a case one may also speak of  $\sup(s, s')$  as the “superposition of the states  $s$  and  $s'$ ”. The internal dynamics of the network is captured by the next-state mapping  $ns$  that is determined by the pattern of excitatory and inhibitory edges in the network.

Now, just as in the example above, we consider an input which is regarded to be represented by a state  $s^* \in S$  and which is supposed to be held fixed for a sufficiently long duration of time. The state transition mapping  $F_{s^*}$  can then be defined by taking both the internal next-state mapping and the input  $s^*$  into account: The next state of the system is given by the superposition of  $s^*$  with the next internal state  $ns(s)$ , i.e.:

$$F_{s^*}(s) := \sup(s^*, ns(s))$$

The dynamics of our dynamical systems is thus determined by iteratively applying  $F_{s^*}$  to the initial state. Fixed points  $s_{stab}$  of  $F_{s^*}$  are again regarded to be the “answers” which the system gives to  $s^*$ . Note that in general there may be *more than just one stable state* for the state transition mapping  $F_{s^*}$  that is determined by the input  $s^*$  (and by the given dynamical system), and there may also be *no stable state* at all for  $F_{s^*}$ : in the former case, there is more than just one “answer” to the input, in the latter case there is no “answer” at all. The different stable states may be reached by starting the computation in different initial states of the system.

Now we are ready to assign formulas to the states of ordered discrete dynamical system. These formulas are supposed to express the content of the information that

is represented by these states. For this purpose, we fix a propositional language  $\mathcal{L}$  which includes (i) finitely many propositional variables  $p, q, r, \dots$ , (ii) and which is closed under the application of the standard classical propositional connectives, i.e.,  $\neg, \wedge, \vee, \supset, \top, \perp$ , where  $\top$  is the *logical verum* (a tautology) and  $\perp$  is the *logical falsum* (a contradiction). The formulas of  $\mathcal{L}$  do not yet include any of the nonmonotonic conditional signs  $\Rightarrow$  that we are interested in. The assignment of formulas to states is achieved by an interpretation mapping  $\mathfrak{I}$ . If  $\varphi$  is a formula in  $\mathcal{L}$ , then  $\mathfrak{I}(\varphi)$  is the state that carries exactly the information that is expressed by  $\varphi$ , i.e., not less or more than what is expressed by  $\varphi$ . So we presuppose that for every formula in  $\mathcal{L}$  there is a uniquely determined state the total information of which is expressed by that formula. If expressed in terms of belief, we can say that in the state  $\mathfrak{I}(\varphi)$  *all the system believes is that  $\varphi$* , i.e., the system only believes  $\varphi$  and all the propositions which are contained in  $\varphi$  from the viewpoint of the system (compare Levesque 1990 on the modal logic of the ‘all I know’ operator). We will not demand that every state necessarily receives an interpretation but just that every formula in  $\mathcal{L}$  will be the interpretation of some state. Furthermore, not just any assignment of states to formulas will do, but we will additionally assume certain postulates to be satisfied which will guarantee that  $\mathfrak{I}$  is compatible with the information ordering that was imposed on the states of the system beforehand. An ordered discrete dynamical system together with such an interpretation mapping is called an ‘interpreted ordered system’ (cf. Leitgeb 2005a). This is the definition stated in detail:

**Definition 3**  $\mathcal{S}_{\mathfrak{I}} = \langle S, ns, \leq, \mathfrak{I} \rangle$  is an interpreted ordered system :iff

1.  $\langle S, ns, \leq \rangle$  is an ordered discrete dynamical system.
2.  $\mathfrak{I} : \mathcal{L} \rightarrow S$  (the interpretation mapping) is such that the following postulates are satisfied:
  - (a) Let  $\mathcal{TH}_{\mathfrak{I}} = \{\varphi \in \mathcal{L} \mid \text{for all } \psi \in \mathcal{L}: \mathfrak{I}(\varphi) \leq \mathfrak{I}(\psi)\}$ :  
then it is assumed that for all  $\varphi, \psi \in \mathcal{L}$ : if  $\mathcal{TH}_{\mathfrak{I}} \vdash \varphi \supset \psi$ , then  $\mathfrak{I}(\psi) \leq \mathfrak{I}(\varphi)$ .
  - (b) For all  $\varphi, \psi \in \mathcal{L}$ :  $\mathfrak{I}(\varphi \wedge \psi) = \sup(\mathfrak{I}(\varphi), \mathfrak{I}(\psi))$ .
  - (c) For every  $\varphi \in \mathcal{L}$ : there is an  $\mathfrak{I}(\varphi)$ -stable state.
  - (d) There is an  $\mathfrak{I}(\top)$ -stable state  $s_{stab}$ , such that  $\mathfrak{I}(\perp) \not\leq s_{stab}$ .

We say that  $\mathcal{S}_{\mathfrak{I}}$  satisfies the uniqueness condition :iff  
for every  $\varphi \in \mathcal{L}$  there is precisely one  $\mathfrak{I}(\varphi)$ -stable state.

How can these postulates be justified? First of all,  $\mathcal{TH}_{\mathfrak{I}}$  is the set of formulas which are the interpretation of states which carry less information than, or an equal amount of information as, any other state with an interpretation. Hence, if  $\varphi \in \mathcal{TH}_{\mathfrak{I}}$ , then the information expressed by  $\varphi$  is contained in every interpreted state of the system. If this is spelled out in terms of belief, then we can say: if  $\varphi \in \mathcal{TH}_{\mathfrak{I}}$ , then  $\varphi$  is believed by the system in every state that has an interpretation. For the same reason, such a belief cannot be revised by the system – it is “built” into the interpreted ordered system independent of its current input or state, as long as the state that it is in has an interpretation at all. In more traditional philosophical terms, we might say that every such

formula is believed *a priori* by the system. So if a material conditional  $\varphi \supset \psi$  follows logically from  $\mathcal{TH}_\mathfrak{I}$ , then – since (rational) belief is closed under logical deduction – also  $\varphi \supset \psi$  must be (rationally) believed by the system in every interpreted state whatsoever; indeed we may think of such a conditional as a strict *a priori* conditional: it is a material conditional which is epistemically necessary in the sense of being entailed by  $\mathcal{TH}_\mathfrak{I}$ , so if  $\Box$  expresses entailment by  $\mathcal{TH}_\mathfrak{I}$ , then for every conditional  $\varphi \supset \psi$  that is derivable from  $\mathcal{TH}_\mathfrak{I}$  it holds that  $\Box(\varphi \supset \psi)$ . But if this so, then the system must regard the propositional information that is expressed by  $\psi$  to be included in the propositional information that is expressed by  $\varphi$  – from the viewpoint of the system,  $\varphi$  must express a stronger proposition than  $\psi$ . In this case, with respect to the information ordering of the system, the state that belongs to  $\psi$  should be “below” the state that is associated with  $\varphi$  or at worst the two states should be equal in the information ordering. In other words,  $\mathfrak{I}(\psi) \leq \mathfrak{I}(\varphi)$  ought to be the case. This is exactly what is expressed by postulate 2a.

Postulate 2b is more easily to explain and justify: the state that belongs to a conjunctive formula  $\varphi \wedge \psi$  should be the supremum of the two states that are associated with the two conjuncts  $\varphi$  and  $\psi$ , just as the proposition expressed by a conjunctive sentence is the supremum of the propositions expressed by its two conjuncts in the partial order of logical entailment.

Postulate 2c makes sure that we are dealing with systems which have at least one “answer” – whether right or wrong – to every “question” posed to the system.

Postulate 2d only allows for interpreted ordered systems which do not end up believing a contradiction when they receive a trivial or empty information (i.e.,  $\top$ ) as an input.

Finally, we are in the position to define what it means for a *nonmonotonic* conditional to be satisfied by an interpreted ordered system. Consider an arbitrary conditional  $\varphi \Rightarrow \psi$  where  $\varphi$  and  $\psi$  are members of our language  $\mathcal{L}$  from above, and where  $\Rightarrow$  is a new nonmonotonic conditional sign. Then we say that a system satisfies  $\varphi \Rightarrow \psi$  if and only if whenever the state that is associated with  $\varphi$  is fed into the system as an input, i.e., whenever the input represents a total belief in  $\varphi$ , the system will eventually end up believing  $\psi$  in its “answer states”, i.e., the state that is associated with  $\psi$  is contained in all the states which are stable with respect to this input. If we collect all such conditionals  $\varphi \Rightarrow \psi$  which are satisfied by the system, then we get what we call the ‘conditional theory’ corresponding to the system. In formal terms:

**Definition 4** Let  $\mathcal{S}_\mathfrak{I} = \langle S, ns, \leq, \mathfrak{I} \rangle$  be an interpreted ordered system:

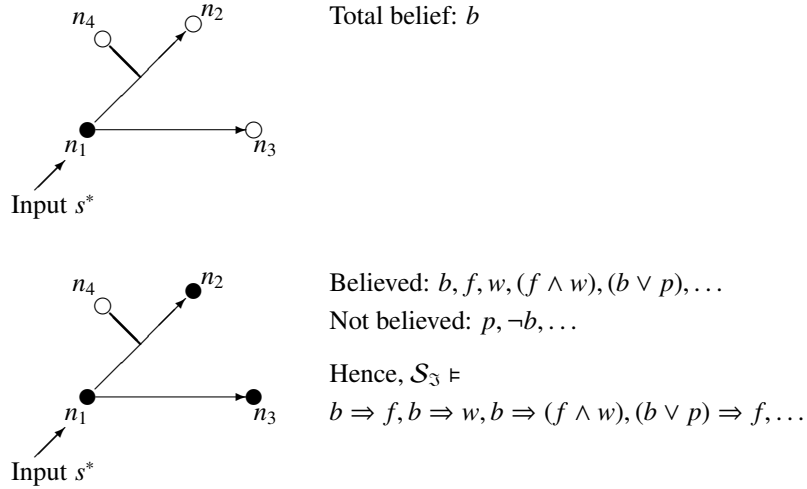
1.  $\mathcal{S}_\mathfrak{I} \models \varphi \Rightarrow \psi$  : iff for every  $\mathfrak{I}(\varphi)$ -stable state  $s_{stab}$ :  $\mathfrak{I}(\psi) \leq s_{stab}$ .
2.  $\mathcal{TH}_\Rightarrow(\mathcal{S}_\mathfrak{I}) = \{\varphi \Rightarrow \psi \mid \mathcal{S}_\mathfrak{I} \models \varphi \Rightarrow \psi\}$   
(the conditional theory corresponding to  $\mathcal{S}_\mathfrak{I}$ ).

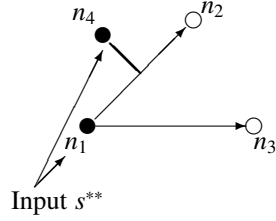
Leitgeb (2007) gives an interpretation of the cognitive states that satisfy conditionals in this way in terms of so-called *conditional beliefs* where conditional beliefs are to be distinguished conceptually from beliefs in conditionals.



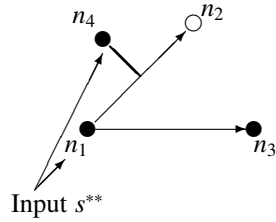
Here is an example: consider again the simple qualitative network which we presented as a discrete ordered dynamical system above. In order to turn it into an *interpreted* ordered system, we have to equip it with an interpretation mapping  $\mathfrak{I}$  that is defined on a propositional language  $\mathcal{L}$ . Let e.g.  $\mathcal{L}$  be determined by the set  $\{b, f, w, p\}$  of propositional variables (for ‘Tweety is a bird’, ‘Tweety is able to fly’, ‘Tweety has wings’, ‘Tweety is a penguin’, or, alternatively, ‘ $x$  is a bird’, ‘ $x$  is able to fly’, ‘ $x$  has wings’, ‘ $x$  is a penguin’). We choose the following interpretation mapping: let  $\mathfrak{I}(b) = \{n_1\}$ ,  $\mathfrak{I}(f) = \{n_1, n_2\}$ ,  $\mathfrak{I}(w) = \{n_1, n_3\}$ ,  $\mathfrak{I}(p) = \{n_1, n_4\}$ , and  $\mathfrak{I}(\neg\varphi) = 1 - \mathfrak{I}(\varphi)$ , where the latter is to be understood in the way that whenever a node is active in  $\mathfrak{I}(\varphi)$  then the same node is inactive in  $\mathfrak{I}(\neg\varphi)$  and vice versa.<sup>5</sup> One can show that there is one and only one interpretation which has these properties and which also satisfies the postulates in definition 3. Note that we have assumed  $\mathfrak{I}(\neg\varphi) = 1 - \mathfrak{I}(\varphi)$  just for convenience, as it becomes easier then to pin down an interpretation for our example. It is not *implied* at all by definition 3 that the pattern of active nodes that is associated with a negation formula  $\neg\varphi$  is actually identical to the complement of the pattern of active nodes that belongs to the formula  $\varphi$ ; this is merely the way in which we set up our example. One consequence of this choice of  $\mathfrak{I}$  is that e.g. the following material conditionals turn out to be members of  $\mathcal{TH}_3$ :  $p \supset b$ ,  $(p \wedge w) \supset b$ ,  $\neg b \supset \neg p$ , and so forth.

Reconsidering our example from above, the dynamics of the system which we studied back then now turns out to have the following symbolic counterparts:





Total belief:  $b \wedge p$



Believed:  $b, p, w, (b \wedge p \wedge w), \neg f, \dots$

Not believed:  $f, \neg b, (f \vee \neg b), \dots$

Hence,  $\mathcal{S}_3 \models$

$(b \wedge p) \Rightarrow \neg f, p \Rightarrow b, p \Rightarrow \neg f, \dots$

By the way: obviously, there will be lots of if-then “laws” about birds and penguins which this interpreted ordered system will get wrong. After all, it would be very surprising indeed if a little network with just four nodes were able to represent all of the systematic relationships between birds and penguins and flying and wings faithfully. But the example should suffice to give a clear picture of how the definitions above are to be applied.

So we find that in this case  $\mathcal{TH}_{\Rightarrow}(\mathcal{S}_3)$  contains e.g.  $b \Rightarrow f, b \Rightarrow w, b \Rightarrow (f \wedge w), (b \vee p) \Rightarrow f, (b \wedge p) \Rightarrow \neg f, p \Rightarrow b, p \Rightarrow \neg f$  without containing e.g.  $b \Rightarrow p, (b \vee p) \Rightarrow p, (b \wedge p) \Rightarrow f$ . In particular, we see that  $b \Rightarrow f \in \mathcal{TH}_{\Rightarrow}(\mathcal{S}_3)$  while  $(b \wedge p) \Rightarrow f \notin \mathcal{TH}_{\Rightarrow}(\mathcal{S}_3)$ .

What can be said in general terms about the conditional theories  $\mathcal{TH}_{\Rightarrow}$  corresponding to interpreted dynamical systems? Here is the answer from the logical point of view:

**Theorem 5** (*Soundness of C*)

Let  $\mathcal{S}_3 = \langle \mathcal{S}, ns, \leq, \mathfrak{I} \rangle$  be an interpreted ordered system:

Then  $\mathcal{TH}_{\Rightarrow}(\mathcal{S}_3)$  is sound with respect to the rules of the system C of nonmonotonic conditional logic (see Kraus et al. 1990 for details on this system), i.e.:

1. For all  $\varphi \in \mathcal{L}$ :  $\varphi \Rightarrow \varphi \in \mathcal{TH}_{\Rightarrow}(\mathcal{S}_3)$  (Reflexivity)
2.  $\mathcal{TH}_{\Rightarrow}(\mathcal{S}_3)$  is closed under the following rules: for  $\varphi, \psi, \rho \in \mathcal{L}$ ,

$$\frac{\mathcal{TH}_3 \vdash \varphi \leftrightarrow \psi, \varphi \Rightarrow \rho}{\psi \Rightarrow \rho} \text{ (Left Equivalence)}$$

$$\frac{\varphi \Rightarrow \psi, \mathcal{TH}_3 \vdash \psi \rightarrow \rho}{\varphi \Rightarrow \rho} \text{ (Right Weakening)}$$

$$\frac{\varphi \Rightarrow \psi, \varphi \wedge \psi \Rightarrow \rho}{\varphi \Rightarrow \rho} \text{ (Cautious Cut)}$$

3. If  $\mathcal{S}_3$  satisfies the uniqueness condition (remember definition 3), then  $\mathcal{TH}_\Rightarrow(\mathcal{S}_3)$  is also closed under

$$\frac{\varphi \Rightarrow \psi, \varphi \Rightarrow \rho}{\varphi \wedge \psi \Rightarrow \rho} \text{ (Cautious Monotonicity)}$$

4.  $\mathcal{TH}_\Rightarrow(\mathcal{S}_3)$  is consistent, i.e.,  $\top \Rightarrow \perp \notin \mathcal{TH}_\Rightarrow(\mathcal{S}_3)$ .

So given the uniqueness assumption – an interpreted ordered system has a unique answer to each interpreted input – the class of conditionals it satisfies is closed under a well-known and important system of nonmonotonic conditional logic, namely the system C of *cumulative reasoning* which is given by the rules listed above. Note that monotonicity, or strengthening of the antecedent, is *not* a valid rule for interpreted systems: as our example from above has shown, there may be formulas  $\varphi, \psi, \rho$  in  $\mathcal{L}$ , such that the conditional  $\varphi \Rightarrow \psi$  is satisfied by a system but  $\varphi \wedge \rho \Rightarrow \psi$  is not.

One can also show a corresponding completeness theorem for the system C with respect to our interpreted ordered systems semantics for  $\Rightarrow$ :

**Theorem 6** (Completeness of C)

Let  $\mathcal{TH}_\Rightarrow$  be a consistent theory of conditionals closed under the rules of C while extending a given classical theory  $\mathcal{TH}$  as expressed by the Left Equivalence and the Right Weakening rules:

It follows that there is an interpreted ordered system  $\mathcal{S}_3 = \langle S, ns, \leq, \mathfrak{I} \rangle$ , such that  $\mathcal{TH}_\Rightarrow(\mathcal{S}_3) = \mathcal{TH}_\Rightarrow$ ,  $\mathcal{TH}_3 \supseteq \mathcal{TH}$ , and  $\mathcal{S}_3$  satisfies the uniqueness condition.

This means that whatever conditional theory you might be interested in, as long as it is closed under the rules of the system C it is possible to find an interpreted ordered system which satisfies precisely the conditionals contained in that theory (and no other conditionals).

It is also possible to extend these results into various directions. In particular, some interpreted ordered systems can be shown to have the property that each of their states  $s$  may be decomposed into a set of substates  $s_i$  which can be ordered in a way such that the dynamics for each substate  $s_i$  is determined by the dynamics for the substates  $s_1, s_2, \dots, s_{i-1}$  at the previous point of time. Such systems are called ‘hierarchical’ in Leitgeb (2005a). We will not go into any details, but one can prove further soundness and completeness theorems for such *hierarchical* interpreted systems and the system  $CL = C + \text{Loop}$  of nonmonotonic conditional logic, where Loop is the following rule:

$$\frac{\varphi_0 \Rightarrow \varphi_1, \varphi_1 \Rightarrow \varphi_2, \dots, \varphi_{j-1} \Rightarrow \varphi_j, \varphi_j \Rightarrow \varphi_0}{\varphi_0 \Rightarrow \varphi_j} \text{ (Loop)}$$

Note that Loop is a weakened version of transitivity, whereas standard transitivity is *not* valid, just as the rule of cautious monotonicity above is a weakened version of monotonicity without standard monotonicity being valid. (Consult Kraus et al. 1990 for more information on CL.)

In Leitgeb (2003, 2004) further soundness and completeness theorems can be found for more restricted classes of interpreted dynamical systems and even stronger logical systems for nonmonotonic conditionals. E.g., the important system P of so-called *preferential reasoning*, where P results from adding the rule

$$\frac{\varphi \Rightarrow \rho, \psi \Rightarrow \rho}{(\varphi \vee \psi) \Rightarrow \rho} \text{ (Or)}$$

to the system CL, is sound and complete with respect to a class of interpreted dynamical systems. P is exactly Adams' (1975) logical system for indicative conditionals as well as the "flat" fragment of Lewis' (1973) logic for subjunctive conditionals ('flat' means: if iterations and other compositions of subjunctive conditionals are ignored). Moreover, various semantics for nonmonotonic reasoning have been found to "converge" on system P as their logical calculus.

As it turns out, if artificial neural networks are extended by an information ordering and an interpretation mapping along the lines explained above, then they are special cases of interpreted ordered systems; moreover, if the underlying artificial neural network consists of layers of nodes, such that the layers are arranged hierarchically and all connections between nodes are only from one layer to the next one, then the interpreted ordered system is indeed a hierarchical one.

In more formal detail:  $\langle U, W, A, O, NET, ex \rangle$  is an artificial neural network :iff

1.  $U$  is a finite and nonempty set of nodes.
2.  $W : U \times U \rightarrow \mathbb{R}$  assigns a weight to each edge between nodes.
3.  $A$  maps each node  $u \in U$  to an activation mapping  $A_u : \mathbb{R}^3 \rightarrow \mathbb{R}$  such that the activation state  $a_u(t+1)$  of  $u$  at time  $t+1$  depends on the previous activation state  $a_u(t)$  of  $u$ , the current net input  $net_u(t+1)$  of  $u$ , and the external input  $ex(u)$  fed into  $u$ , i.e.  $a_u(t+1) = A_u(a_u(t), net_u(t+1), ex(u))$ .
4.  $O$  maps each node  $u \in U$  to an output mapping  $O_u : \mathbb{R} \rightarrow \mathbb{R}$  such that the output state  $o_u(t+1)$  of  $u$  at time  $t+1$  is solely dependent on the activation state  $a_u(t+1)$  of  $u$ , i.e.  $o_u(t+1) = O_u(a_u(t+1))$ .
5.  $NET$  maps every node  $u \in U$  to a net input (or propagation) mapping  $NET_u : (\mathbb{R} \times \mathbb{R})^U \rightarrow \mathbb{R}$  such that the net input  $net_u(t+1)$  of  $u$  at time  $t+1$  depends on the weights of the edges leading from nodes  $u'$  to  $u$ , and on the previous output states of the nodes  $u'$ , i.e.  $net_u(t+1) = NET_u(\lambda u'. \langle W(u', u), o_{u'}(t) \rangle)$ .<sup>6</sup>
6.  $ex : U \rightarrow \mathbb{R}$  is the external input function.

We can view such networks as ordered dynamical systems if we define:

1.  $S = \{s \mid s : U \rightarrow \mathbb{R}\}$ .
2.  $ns : S \rightarrow S$  with  $ns(s)(u) := A_u(s(u), NET_u(\lambda u'. \langle W(u', u), O_{u'}(s(u')) \rangle), 0)$   
(so in the case of the internal next-state function  $ex(u)$  is set to 0).
3.  $\leq \subseteq S \times S$  with  $s \leq s'$  iff for all  $u \in U$ :  $s(u) \leq s'(u)$ .  
( $sup(s, s')$  is thus simply  $max(s, s')$ .)

$\langle S, ns, \leq \rangle$  is an ordered discrete dynamical system, such that  $F_{s^*}(s) = \sup(s^*, ns(s)) = \max(s^*, ns(s))$  which entails that  $F_{s^*}(s)(u) = \max(s^*(u), ns(s)(u)) = \max(s^*(u), A_u(s(u), NET_u(\lambda u'. \langle W(u', u), O_{u'}(s(u')) \rangle), 0))$ , which corresponds to the assumption that the external input to a network interacts with the current activation state of the network by taking the maximum of both. Given this assumption, the dynamics of artificial neural networks and the dynamics of the corresponding ordered dynamical systems coincide. If the network is layered, then the corresponding ordered system is hierarchical. Stable states are regarded as the relevant “answer” states just as in the standard treatment of neural networks. If such networks are equipped with a corresponding interpretation mapping  $\mathfrak{I}$  as defined above, they satisfy conditional theories which are closed under the rules of well-establish systems of logic for nonmonotonic conditionals.

Furthermore, on the level of representation or interpretation we have:

- In interpreted ordered systems, propositional formulas are represented by total states  $s$  of the system; in particular, in interpreted neural networks, propositional formulas are represented by patterns of activity distributed over the nodes of the network.
- In interpreted ordered systems, nonmonotonic conditionals are represented by the overall dynamics of the system; in particular, in interpreted neural networks, nonmonotonic conditionals are represented by the network topology and by the way weights are distributed over the connections of the network. It is not single edges which correspond to conditionals, but the conditional theory that belongs to an interpreted network is a set of soft constraints that is represented by the network as a whole.

Thus, in contrast with the old McCulloch-Pitts idea, the representation of formulas in interpreted dynamical systems is distributed, as suggested by connectionists. At the same time, the set of conditionals satisfied by an interpreted dynamical system is closed under the rules of systems of nonmonotonic conditional logic which were introduced, and which have been studied intensively, by researchers in the tradition of the symbolic computation paradigm of cognitive science. Subsymbolic inference may be fundamentally nonmonotonic, as claimed by Smolensky, but this does not mean that it cannot be formalized in logical terms – it only means that the formalization has to be given in terms of systems of nonmonotonic reasoning.

The dynamical systems paradigm and the symbolic computation paradigm may thus be regarded as yielding complementary perspectives on the one and the same cognitive system. Moreover, since nonmonotonic conditionals have been shown to have interpretations in terms of (i) conditional probability measures, and (ii) orderings of possible worlds by degrees of similarity to the actual world or by degrees of normality or plausibility, the nonmonotonic conditionals that are satisfied by interpreted dynamical systems may be taken to represent aspects of either of these important semantic structures which are also used to analyze human communication and reasoning by means of conditionals.

## 5 Some Open Questions

Here is an (incomplete) to-do-list in this area of research:

*Extending soundness/completeness results:* How can the logical systems discussed by Kraus et al. (1990) and Lehmann&Magidor (1992) be characterized in terms of connectionistically plausible and elegant constraints on interpreted dynamical systems? (So far there are only partial answers to this question, sometimes relying on very restricted classes of dynamic systems.) Which logical systems do we get if we drop the uniqueness assumption (see definition 3)? How can full-fledged systems of conditional logic for subjunctive conditionals, for which nesting of conditionals and the application of propositional connectives to such conditionals is well-defined, be represented by means of interpreted dynamical systems?

*Characterizing learning in neural networks by logical rules:* As we have seen, state transitions in a fixed (possibly, trained) neural network can be described in terms of conditionals. However, it is as yet unknown how learning processes in networks – by which the weights in a network change under the influence of a learning algorithm and training data – can be represented by logical rules. Learning schemes such as Hebbian learning or backpropagation might translate into particular systems of inductive logic in which inferences can be drawn from both factual training data and conditionals to learned conditionals. In order to facilitate this study, computer implementations of interpreted networks and their learning algorithms will be crucial.

*Applying the theory to open problems in uncertain reasoning:* The results achieved by the previous tasks are expected to feed back on open problems in uncertain reasoning. E.g.: Belief revision (see Gärdenfors 1988 for the classic reference, and Hansson 1999 for a textbook) was created as a theory for the “one-shot” revision of beliefs by a single piece of evidence. Attempts of extending the theory to iterated occurrences of evidence led to a multitude of suggestions lacking clear philosophical interpretation. By means of the results achieved in this area, it might be possible to understand evidence-induced changes of networks as iterated belief revisions. We hypothesise that different schemes of iterated revision correspond to, and can be understood as, different learning algorithms for neural networks.

*Applying the theory in philosophy of science:* In philosophy of science, it was realized early on that new empirical evidence can have the effect that previous hypotheses must be withdrawn, since an agent might learn that what she had regarded likely is actually not. As Flach (2000) argues, the same logics that govern valid commonsense inferences can be interpreted as logics for scientific induction, i.e., for data constituting incomplete und uncertain evidence for empirical hypotheses. Schurz (2002) demonstrates that scientific laws are subject to normality or *ceteris paribus* restrictions that obey the logic of nonmonotonic reasoning. At the same time, the study of neural networks is expected to transform our philosophical understanding of science: Churchland (1989) presents networks as models of scientific theories and regards prototype representations in networks as a system’s explanatory understanding of its inputs. Bechtel (1996) explains scientific model building in terms of the satisfaction of soft constraints represented in networks. Bird (2002) observes: “The time is ripe for a reassessment of Kuhn’s earlier work in the light of connectionist and neural-net research”. Is it possible to throw some new light on these trends in philosophy of science on the basis of new findings on logical accounts of neural network reasoning and learning?

## Notes

<sup>1</sup>This image is taken from Leitgeb (2005b).

<sup>2</sup>The following famous example is by Ernest Adams.

<sup>3</sup>Such networks are called ‘inhibition networks’ in Leitgeb (2001).

<sup>4</sup>A partial order  $\leq$  (on  $S$ ) is a reflexive, antisymmetric, and transitive binary relation, i.e.: for all  $s \in S$ :  $s \leq s$ ; for all  $s, s' \in S$ : if  $s \leq s'$  and  $s' \leq s$  then  $s = s'$ ; for all  $s_1, s_2, s_3 \in S$ : if  $s_1 \leq s_2$  and  $s_2 \leq s_3$  then  $s_1 \leq s_3$ .

<sup>5</sup>So the 1 here is actually the constant 1-function, i.e., the function that maps each node to the activation value 1.

<sup>6</sup> $\lambda u'. \langle W(u', u), o_{u'}(t) \rangle$  is the function that maps  $u'$  to the pair  $\langle W(u', u), o_{u'}(t) \rangle$ .

## References

- [1] Adams, E. (1975): *The Logic of Conditionals: An Application of Probability to Deductive Logic*, Synthese Library 86, Dordrecht: Reidel.
- [2] Adams, E. (1998): *A Primer of Probability Logic*, CSLI Lecture Notes.
- [3] d’Avila Garcez, A.S., D.M. Gabbay, and L.C. Lamb (200?): *Connectionist Non-Classical Logics: Distributed Reasoning and Learning in Neural Networks*, to appear with Springer.
- [4] d’Avila Garcez, A.S., K. Broda, and D.M. Gabbay (2001): “Symbolic knowledge extraction from trained neural networks: a sound approach”, *Artificial Intelligence* 125, 153–205.
- [5] Balkenius, C., P. Gärdenfors (1991): “Nonmonotonic inferences in neural networks,” in: J.A. Allen, R. Fikes, E. Sandewall (eds.), *Principles of Knowledge Representation and Reasoning*, San Mateo: Morgan Kaufmann, 32–9.
- [6] Bechtel, W. (1996): “What should a connectionist philosophy of science look like?”, in: R.M. McCauley (ed.), *The Churchlands and Their Critics*, Basil Blackwell, 121–44.
- [7] Bennett, J. (2003): *A Philosophical Guide to Conditionals*, Oxford: Clarendon Press.
- [8] van Benthem, J. (1984): “Foundations of conditional logic”, *Journal of Philosophical Logic* 13/3, 303–49.
- [9] Bird, A. (2002): “What is in a paradigm?”, *Richmond Journal of Philosophy* 1.ii, 11–20.
- [10] Blutner, R. (2004): “Nonmonotonic inferences and neural networks”, *Synthese* 142, 143–74.
- [11] Brewka, G., J. Dix, K. Konolige (1997): *Nonmonotonic Reasoning. An Overview*, Stanford: CSLI Lecture Notes 73.
- [12] Chalmers, D.J. (1990): “Syntactic transformations on distributed representations,” *Connection Science*, vol.2, nos.1&2, 53–62.

- [13] Chen, C.-H., and V. Honavar (1999): “A neural-network architecture for syntax analysis,” *IEEE Transactions on Neural Networks*, vol.10, no.1, 94–114.
- [14] Churchland, P.M (1989): *A Neurocomputational Perspective: The Nature of Mind and the Structure of Science*, MIT Press.
- [15] Flach, P.A. (2000): “Logical characterizations of inductive learning”, in: D.M. Gabbay, R. Kruse (eds.), *Handbook of Defeasible Reasoning and Uncertainty Management Systems*, Vol. 4, Kluwer, 155–96.
- [16] Fodor, J., Z. Pylyshyn (1988): “Connectionism and cognitive architecture: A critical analysis,” *Cognition* 28, 3–71.
- [17] Gärdenfors, P. (1988): *Knowledge in Flux*, Cambridge, Mass.: The MIT Press.
- [18] Gärdenfors, P. (1994): “How logic emerges from the dynamics of information,” in: J. Van Eijck, A. Visser (Eds.), *Logic and Information Flow*, Cambridge: The MIT Press, 49–77.
- [19] Van Gelder, T.J. (1998): “The dynamical hypothesis in cognitive science,” *Behavioral and Brain Sciences* 21, 615–65.
- [20] Van Gelder, T.J. (1999): “Distributed versus local representation,” in: R. Wilson, F. Keil (Eds.), *The MIT Encyclopedia of Cognitive Sciences*, Cambridge: The MIT Press, 236–38.
- [21] Ginsberg, M.L. (ed.) (1987), *Readings in Nonmonotonic Reasoning*, Los Altos: Morgan Kaufmann, 1–23.
- [22] Hansson, S.O. (1999): *A Textbook of Belief Dynamics*, Dordrecht: Kluwer.
- [23] Hölldobler, S. (1993): *Automated Inferencing and Connectionist Models*, Post-Doctoral Thesis.
- [24] Hölldobler, S., P. Hitzler, and A.K. Seda (2004): “Logic programs and connectionist networks”, *Journal of Applied Logic* 2, 245–72.
- [25] Horgan, T., J. Tienson (1996): *Connectionism and the Philosophy of Psychology*, Cambridge: The MIT Press.
- [26] Kraus, S., D. Lehmann, M. Magidor (1990): “Nonmonotonic reasoning, preferential models and cumulative logics,” *Artificial Intelligence* 44, 167–207.
- [27] Legendre, G., Y. Miyata, P. Smolensky (1994): *Principles for an Integrated Connectionist/Symbolic Theory of Higher Cognition*, Hillsdale: L. Erlbaum.
- [28] Lehmann, D., and M. Magidor (1992): “What does a conditional knowledge base entail?”, *Artificial Intelligence* 55, 1–60.
- [29] Leitgeb, H. (2001): “Nonmonotonic reasoning by inhibition nets”, *Artificial Intelligence* 128[1–2], 161–201.



- [30] Leitgeb, H. (2003): “Nonmonotonic reasoning by inhibition nets II”, *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems* 11, suppl. issue 2, 105–35.
- [31] Leitgeb, H. (2004): *Inference on the Low Level. An Investigation into Deduction, Nonmonotonic Reasoning, and the Philosophy of Cognition*, Dordrecht: Kluwer/Springer, Applied Logic Series.
- [32] Leitgeb, H. (2005a): “Interpreted dynamical systems and qualitative laws: From inhibition networks to evolutionary systems”, *Synthese* 146, 189–202.
- [33] Leitgeb, H. (2005b): “Reseaux de neurones capables de raisonner”, *Dossier Pour la Science* (special issue of the French edition of the *Scientific American*) October/December, 97–101.
- [34] Leitgeb, H. (2007): “Beliefs in conditionals vs. conditional beliefs”, *Topoi* 26/1, 115–32.
- [35] Levesque, H. (1990): “All I know: A study in autoepistemic logic,” *Artificial Intelligence* 42, 263–309.
- [36] Lewis, D. (1973): *Counterfactuals*, Oxford: Blackwell.
- [37] Makinson, D. (1994): “General patterns in nonmonotonic reasoning,” in: D.M. Gabbay, C.J. Hogger, J.A. Robinson (Eds.), *Handbook of Logic in Artificial Intelligence and Logic Programming* 3, Oxford: Clarendon Press, 35–110.
- [38] McCulloch, W.S. and W.H. Pitts (1943): “A logical calculus of the ideas immanent in nervous activity”, *Bulletin of Mathematical Biophysics* 5, 115–33. Reprinted in: W.S. McCulloch, *Embodiments of Mind*, Cambridge, Mass.: The MIT Press, 1965.
- [39] Nauck, D., F. Klawonn, R. Kruse (1994): *Neuronale Netze und Fuzzy-Systeme*, Braunschweig: Vieweg.
- [40] Peters, S. and D. Westerstahl (2006): *Quantifiers in Language and Logic*, Oxford: Oxford University Press.
- [41] Rojas, R. (1993): *Theorie der neuronalen Netze*, Berlin: Springer.
- [42] Rumelhart, D.E., J.L. McClelland, and the PDP Research Group (1986): *Parallel Distributed Processing*, Vol 1 and 2, Cambridge: The MIT Press.
- [43] Schurz, G. (2002): “Ceteris paribus laws: Classification and deconstruction,” *Erkenntnis* 57/3, 351–72.
- [44] Schurz, G. and H. Leitgeb (eds.) (2005): special volume on “Non-monotonic and uncertain reasoning in the focus of paradigms of cognition”, *Synthese* 146/1–2.
- [45] Smolensky, P. (1988): “On the proper treatment of connectionism,” *Behavioral and Brain Sciences* 11, 1–23.