Rationality Postulates for Induction

Peter A. Flach[†]

Abstract. I study the process of inductive hypothesis formation from two perspectives: finding general rules that explain given specific evidence, and finding general rules that are confirmed by the evidence. Both forms of hypothesis formation are axiomatised on the metalevel of consequence relations, which provides us with a set of rationality postulates for various forms of induction.

1. Introduction

In this paper I develop a rationality theory of inductive reasoning, applying the methodology of Kraus, Lehmann & Magidor (1990) to views on philosophy of science developed by Peirce and Hempel. This introductory section sketches the motivation for this paper, defines the main concepts and notation, and provides an overview.

1.1 Motivation

Consider a reasoning agent \mathbf{R} that, upon observing several black crows, conjectures that all crows are black. Such an agent is said to reason *inductively*: from specific observations about members of a class it tentatively draws general conclusions about the whole class. Clearly, not any general conclusion will do: for instance, \mathbf{R} couldn't conclude that all crows are white. Thus, we may draw a distinction between rational and irrational behaviour of \mathbf{R} . This paper studies limits of rational behaviour of inductive reasoning agents.

In the ideal case, rationality simply prescribes what conclusions are allowed from what premisses, as in the case of formal deductive reasoning, where rationality is defined by a fixed notion of (deductive) validity. In the general case, however, rationality draws a distinction between reasoning behaviours that is more complex than distinguishing between valid and invalid arguments. For instance, inductive rationality as studied in this paper does not prescribe the number of black crows one must have observed before adopting the hypothesis 'all crows are black', but it does prescribe that once one adopts this hypothesis, observation of yet another black crow cannot be a reason to reject it. This establishes a rationality principle of the form 'if a reasoning agent accepts this inductive argument, it should also accept that inductive argument'. Most of the rationality principles we will consider in this paper are of this form. Other possible forms we will consider are 'if a reasoning agent accepts this inductive

[†]Author's address: INFOLAB, Tilburg University, PObox 90153, 5000 LE Tilburg, the Netherlands, tel. +31 13 466 3119, fax +31 13 466 3069, email Peter.Flach@kub.nl.

Part of this work was supported by Esprit Long Term Research Project 20237 (Inductive Logic

Programming), and by PECO Pan-European Scientific Network ILPnet (no. CIPA35100CT920044).

P.A. Flach (1996), 'Rationality postulates for induction'. In *Proc. Theoretical Aspects of Rationality and Knowledge*, Yoav Shoham (ed.), pp. 267-281, Morgan Kaufmann.

argument, it should not accept that argument' and 'if a reasoning agent does not accept this argument, then it should accept that argument'.

Related to the distinction between rationality and validity is the distinction between the processes of hypothesis formation and hypothesis selection. Hypothesis formation is the process that determines, given certain evidence, the set of possible hypotheses, i.e. conjectures that are not ruled out by the evidence. *Hypothesis selection* is the process of selecting, among the possible hypotheses, one or more that meet given criteria. For instance, the reasoning agent may select only those conclusions that are sufficiently justified by the premisses, typically by assessing the truth of the conclusion given the truth of the premisses. Such an assessment procedure, which can be said to generalise the two-valued function of deductive validity to a conti ously-valued function of inductive validity or confirmation, has for instance bee developed by process Carnap (1953). However, in this paper I will concentrate on the purely ogic: of hypothesis formation. That is, an inductive hypothesis is not ily the ces conclusion adopted by the inductive a nt on basis of the evidence possible conclusion.

1.2 **Preliminaries**

In formalising the logic of inducti language L, built up from a set connectives, a set M of propositiona well-behaved with respect to the log be thought f as a set of truth as assumed co pact. As usual, we wr of all possible propos proper subs implicit bac ground theory; thus, allows to co

sume pro propos ion symbols and the isu ves (for all practical purpos conne propositional variables). if m t for all $m \in M$. Note that M may be a

sitional iogical

L that is *M* can will be

s, representing the set of models of an al mode c can be interpreted as 'the background theory

An in we say th consequenc

ctive consequence relation $\nvDash \subseteq L \times L$ is a set of pairs of formulae; if $\alpha \nvDash \beta$ β is a possible inductive hypothesis given evidence α . Inductive del the behaviour of inductive agents — at this stage we do to fix a particular definition of k, but study rationality postulates limiting possible defitions. For instance, given background knowledge including

rows_are_black->chevy_is_black

and the fact that the reasoner accepts the inductive argument

```
chad_is_black < crows_are_black
```

our rationality postulates should prescribe that the hypothesis 'all crows are black' can still be maintained after observation of Chevy being black, by stipulating that the inductive argument

chad_is_black < crows_are_black < crows_are_black



should likewise be accepted by t inductive agent. This observation gives rise to the following postulate:

(1.1) If $\alpha \in \beta$ and $\rightarrow \gamma$, then $\alpha \land \gamma \in \beta$.

This postulate expresses a principle of *verification*: if hypothesis β is tentatively concluded on the basis of evidence α , and prediction γ drawn from β is subsequently observed, then this counts as a verification of β . In this paper I will investigate a range of such rationality postulates, as well as some of the ways in which they interact.

1.3 Overview of the paper

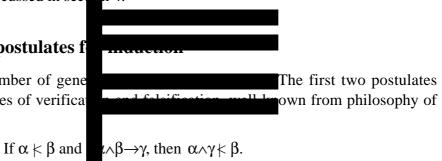
As intuition constitutes the primary source of justification for such rationality postulates, we need to proceed with care — after all, intuitions sometimes turn out to be incompatible, leading to paradoxical situations. One such a paradox was noted by the philosopher Carl G. Hempel, when he investigated the logic of confirmation (Hempel, 1945). As I have demonstrated elsewhere this particular paradox can be dissolved by making a clear distinction between the statements 'this evidence *confirms* that hypothesis' and 'this evidence *is explained by* that hypothesis' (Flach, 1995). In general, the search for confirmed hypotheses may lead to different conjectures than the search for explanatory hypotheses, an observation that is backed up by recent experimental and theoretical work in inductive logic programming (De Raedt & Bruynooghe, 1993). In this paper I will therefore consider two different forms of induction, which I have termed *explanatory induction* (section 2.1) and *confirmatory induction* (section 2.2).

In addition to studying single rationality postulates, I will also study some of the ways in which they interact. Ultimately, this is done by constructing a semantics, and establishing the equivalence of set of postulates and the semantics by means of a representation theorem. In section 3 two such representation theorems will be given, one for each form of induction the significance, limitations and prospects of the approach will be discussed in section 4.

2. Rationality postulates f

We start with a number of gene express the principles of verificar science:

(I1)



β.

(I2) If
$$\alpha \not\models \beta$$
 and $\alpha \land \beta \rightarrow \gamma$, then $\alpha \land \neg \gamma \not\models$

(I1) has been considered above in a slightly less general form. The difference with (1.1) is that predictions may be deduced from a hypothesis together with the evidence on which the hypothesis is based — one might say that, in general, the 'epistemic outcome' of accepting an inductive argument $\alpha \ltimes \beta$ is the belief that $\alpha \land \beta$ rather than

just β . This is relevant in cases required to include all the most from presence by the evidence. Of course, such a prediction may also be ided as a urther hypothesis: (I3) If $\leq \beta$ and $\alpha \beta \rightarrow \gamma$, then $\alpha \neq \beta \land \gamma$.

The postulate of disfication (I2) can be equivalently stated as follows (the proof requires the postulate of eff logical equivalence stated below as (I7)):

(I2') If
$$\rightarrow \neg \alpha$$
, then $\alpha \not < \beta$.

Principle (I2'), obvious as it may be, has a few technical consequences. For instance, we have reflexivity only for consistent formulae. Clearly, in the light of (I2') a formula is consistent if it occurs in an arbitrary inductive argument, so we have the following two weaker versions of reflexivity:

(I4) If
$$\alpha \in \beta$$
, then $\alpha \in \alpha$.
(I5) If $\alpha \in \beta$, then $\beta \in \beta$.

Given an inductive consequence Any admissible formula is consi sidered later the converse of this quires consistency of a formula the

The next two postulates e contents of premiss and conclusion

$$k = \beta$$
.
 α admissible iff $\alpha = \alpha$.
 α of postulates to be con-
tement is also true. Hence, whenever a postulate re-
lition of the form $\alpha \neq \alpha$.
 $\alpha = \alpha + \alpha$.
 $\alpha = \alpha + \alpha + \alpha$.
 $\alpha = \alpha + \alpha + \alpha$.
 $\alpha = \alpha + \alpha + \alpha + \alpha$.

(I6) If
$$\alpha \not\models \beta$$
 and $\varphi \leftrightarrow \gamma$, then $\alpha \not\models$
(I7) If $\alpha \not\models \gamma$ and $\xi \leftrightarrow \beta$, then $\beta \not\models$

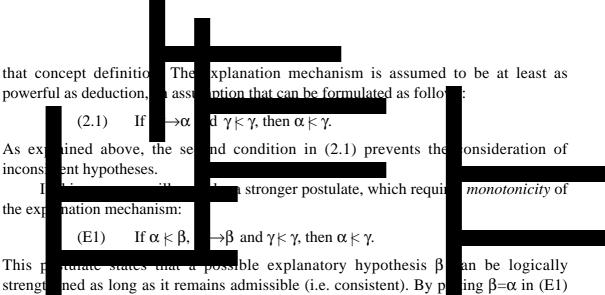
One may argue that, in some cases, the syntactic form of a hypothesis is relevant: for instance, in inductive logic programming the hypotheses take the form of logic programs, and clearly an efficient or at least terminating program is desirable over an inefficient or non-terminating one. However, we deem such considerations to be extralogical. Against (I7) one may raise the objection that many concept learning algorithms perform differently when the examples are re-ordered. However, such behaviour is peculiar to a particular algorithm: an algorithm that would be order-independent would not behave irrationally.

γ.

Clearly, postulate (I1–7) are too weak to distinguish induction from other forms of reasoning. In order to obtain stronger sets of postulates we need to consult our intuitions regarding explanatory and confirmed hypotheses, respectively.

2.1 Rationality postulates for explanatory induction

Throughout this section $\alpha \not\in \beta$ is interpreted as ' β is a possible explanation of α ', which is loosely taken to mean that α can be derived from hypothesis β by means of some *explanation mechanism*. For instance, in inductive concept learning β would be a concept definition, and α would represent the classification of an instance according to



strengt ned as long as it remains admissible (i.e. consistent). By pluing $\beta = \alpha$ in (E1) we obtain (2.1), if in addition we assume that $\alpha \not\models \alpha$ follows from $\rightarrow \alpha$ and $\gamma \not\models \gamma$. Since $\gamma \rightarrow \alpha$ implies $\neg \alpha \not\models \gamma$ by (I2'), this assumption is assured by the following postulate:

(E2) If $\gamma \in \gamma$ and $\neg \alpha \notin \gamma$, then $\alpha \notin \alpha$.

This postulate represents a third weakening (besides (I4) and (I5) discussed above) of reflexivity. The underlying intuition is best explained by considering its contrapositive: if $\alpha \not \ll \alpha$, i.e. α is too strong with respect to the background theory, then its negation $\neg \alpha$ is so weak that it is explained by arbitrary admissible γ .

The next postulate expresses a principle well-known from algorithmic concept learning: if α represents the classification of an instance and β its description, then we may either induce a concept definition from examples of the form $\beta \rightarrow \alpha$, or we may add β to the background theory and induce from α alone. Since in our framework background knowledge is include implicitly, β is added to the hypothesis instead.

(E3) If $\alpha \in \beta \land \gamma$, to $\beta \rightarrow \alpha \in \gamma$.

The converse implication will be

The last two postulates compositionality, namely $\alpha \land \beta \models \beta$ statement of one half of this prin form:

system. explanatory induction establish a principle of ff $\alpha \notin \gamma$ and $\beta \notin \gamma$. The first postulate is simply the planatic state of the second takes a slightly more general

(E4) If $\alpha \vDash \gamma$ and $< \gamma$, then $\alpha \land \beta \vDash \gamma$. (E5) If $\alpha \vDash \gamma$ and $\alpha \rightarrow \beta$, then $\beta \vDash \gamma$.

Both postulates are of considerable importance for computational induction, since they allow for an incremental approach. (E4) states that pieces of evidence can be dealt with in isolation. Another way to say the same thing is that the set of evidence explained by a given hypothesis is conjunctively closed. By the consistency principle (I2') this set is consistent, which yields the following principle:

(2.2) If
$$\alpha \not\in \beta$$
, then $\neg \alpha \not\in \beta$.

(E5) states a monotonicit understood by considering its con evidence β cannot become feasi other words: the process of rej assumptions), but based on the ev property of deduction (which is postulate (E5) can be jointly expr

It is easy to show that, conversely 2.2) implies (I2') in the presence of (2.1) and (I5). property of induction, which can again best be positive: a hypothesis that is rejected on the basis of dence α is available. In ing a hypothesis is not defeasible (i.e. based on ence only. This is the analogue of the monotonicity nplication in the second condition). We may further note that the verification postulate (I1) and the monotonicity sed as

(E5') If
$$\alpha \not\models \beta$$
 and $\alpha \land \beta \rightarrow \gamma$, then $\gamma \not\models \beta$.

Again, this formulation stresses the fact that the epistemic outcome of accepting an inductive argument $\alpha \ltimes \beta$ is $\alpha \land \beta$.

Rationality populates for confirmatory induction 2.2

lecrease

explanat

Throughout this section $\alpha \in \beta$ is interpreted as ' β is confirmed by α ', which is loosely taken to mean that β instance, α could be regularity 'investment: too weak to count as a rate — it is merely bein

The following p confirmation:

(2.3) If
$$\iota \rightarrow \gamma$$
 and

This principle may be generalise consequences of confirmed hypot

(C1) If
$$\alpha \not\in \beta$$
 and $\rightarrow \gamma$

This principle of right weaker strengthening postulate (E1) ob combining (C1) and (E1) leads to hypothesis, Hempel dismissed (same, since (E1) and (C1) are ir Note that (C1) and (I3) can be join y expressed as

ciple ex tive entaiment is a special case of

presses some regularity that is implicitly present in α . For

o the following postulate, which states that logical ses are also confirmed:

hen interests increase'. Note that this hypothesis is

i of the current of investment rate given the interest

omical data, and β could be the

 γ , then $\alpha < \gamma$.

guished from the right ned for explanatory induction. Upon noting that situation in which arbitrary evidence confirms any eel compelled to do the nded to formalise quite different intuitive notions.

```
If \alpha \in \beta and
                                                       \alpha \wedge \beta \rightarrow \gamma, then \alpha < \gamma.
(C1')
```

The next postulate represents a confirmatory variant of reflexivity, analogous to (E2).

> If $\alpha \ltimes \alpha$ and $\alpha \nvDash \neg \beta$, then $\beta \ltimes \beta$. (C2)

¹Stricly speaking, it would be better to say that the hypothesis is *not disconfirmed* by the evidence.

Taken contrapositively, (C2) states that if some β does not confirm itself, it must be too strong a statement — hence, its negation is so weak that it is confirmed by arbitrary admissible α .

The remaining postulates for confirmatory induction we will consider are considerably stronger. The following postulate states that the conjunction of confirmed hypotheses is itself confirmed:

(C3) If
$$\alpha \in \beta$$
 and $\alpha \in \gamma$, then $\alpha \in \beta \land \gamma$.

This postulate can be viewed as a completeness assumption with regard to the evidence, in the sense that α contains all the possible states of affairs that are *relevant* for the domain. For instance, when talking about ravens and their colours such complete evidence would contain some black ravens, but also a few non-ravens, some of which are black and some of which are not. Only under such a completeness assumption one can safely assume that 'what you are not told looks like what you are told' (Helft, 1989, p.149). Hempel considered (C3) a valid principle of confirmation. On the other hand, it is clear that (C3) does not necessarily follow from our intuitions about confirmation — indeed, Schlechta (1995) rejects it. Notice that in the presence of (C3) the consistency principle (I2') is equivalent with

(2.4) If
$$\alpha \ltimes \beta$$
, then $\alpha \nvDash \neg \beta$.

Notice also that (C3) contradicts an incrementality principle like (E5): weakening the evidence typically disconfirms some hypotheses. A considerably weaker 'left weakening' principle that is compatible with (C3) is expressed by the following postulate:

(C4) If $\alpha \in \gamma$ and $\beta \in \gamma$, then $\alpha \lor \beta \in \gamma$.

This postulate states that if both α and β confirm γ , then the knowledge that at least one of them is true should not disconfirm γ .

The final postulate for confirmatory induction we will consider reinvokes the principle of verification (I1), but in a much stronger sense. Recall that verification expresses that if α confirms β , and γ is a prediction drawn from α and β , then subsequent observation that γ is indeed the case does not disconfirm β . The following postulates extends this to arbitrary γ that are also confirmed by α :

(C5) If
$$\alpha \in \beta$$
 and $\alpha \in \gamma$, then $\alpha \land \gamma \in \beta$.

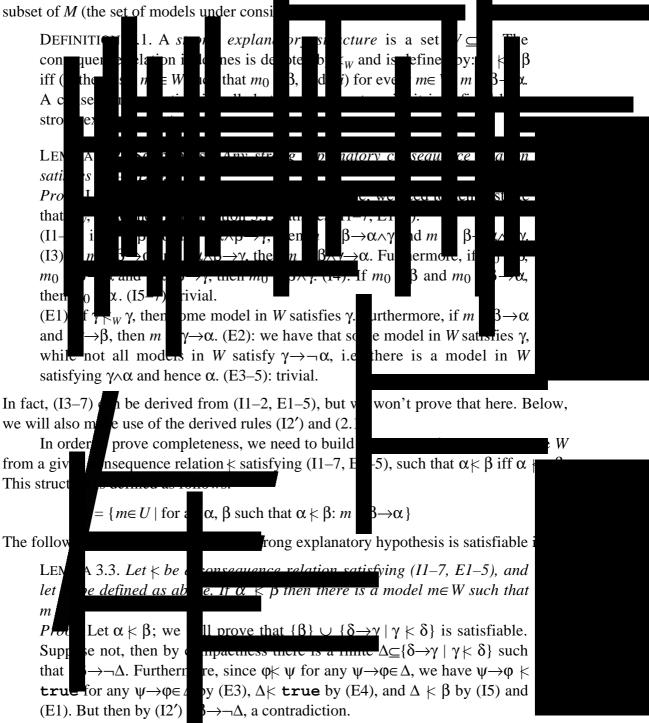
In other words: the assumptions necessary to induce β from α do not contradict the assumptions needed to obtain γ from α — thus (C5) expresses another aspect of the 'completeness' of the evidence.

3. Representation theorems

In this section I present some initial representation results regarding the postulates discussed above.

3.1 Strong explanatory structures

The first result is that postulates (E1-5) characterise explanatory induction as a value also envisage weaker (nonmonotonic) not called *strong* explanatory. The corresponsible subset of *M* (the set of models under consi



e strong enough, together w

ns of explanation, this form of

ng semantic structure is char

n (I1–7), to

induction is

terised by a

Furthermore, we have that every inadmissible formula is unsatisfiable in W.

LEMMA 3.4. Let \nvDash be a non-empted by E1-5, and let W be defined as about Proof. Let $\alpha \nvDash \beta$, then **true** \nvDash **true** by (E2), hence m

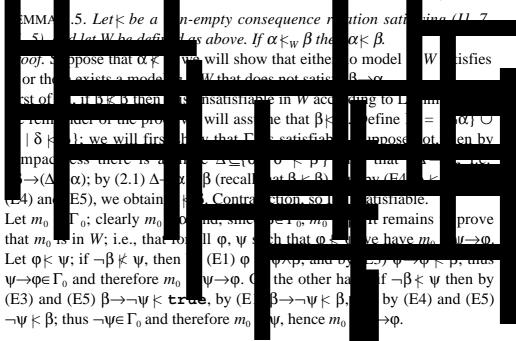
Insequence retation statisfying (11–7, 1. If $\gamma \not\models \gamma$ then γ is unsatisfiable in W. by (E5) and (I4). Further pre, if $\gamma \not\models \gamma$ **rue** $\rightarrow \neg \gamma$ for every $m \in W$

f it satisfies (I1–7,

.nd let B→α}



w show that W defines a consequence relation that is include in $\not\in$.



Armed with the previous three lemmas we can prove the co

THEOREM 3.6 (Representation theorem for strong explanator relations). A consequence relation is strong explanator E1-5).

Proof. The only-if part is Lemma 3.2. For the if par non-empty consequence relation satisfying (I1–7, E1–

 $W = \{ m \in U \mid \text{ for all } \alpha, \beta \text{ such that } \alpha \notin \beta : \}$

Suppose $\alpha \notin \beta$, then by the construction of W, $m \longrightarrow \alpha$ for all $m \in W$. Furthermore, by Lemma 3.3 there is a model in W satisfying β . We may conclude that $\alpha \notin_W \beta$. Conversely, if $\alpha \notin_W \beta$ then Lemma 3.5 proves that $\alpha \notin \beta$. We conclude that W defines a consequence relation that is exactly \notin . For an empty consequence relation put $W = \emptyset$.

3.2 Preferential confirmatory structures

The second result is the characterisation of (C1–5), which is an adaptation of a result by Kraus, Lehmann & Magidor (1990). The main idea is to construct, from the evidence, a set of *regular* interpretations, which exhibit the same regularities as the evidence, and to

define a hypothesis to be confirmed if it is true in every regular interpretation. Notice that such regular interpretations are not necessarily models of the evidence. For instance, in the ravens example we may interchange the names of ravens and nonravens without invalidating the inductive conclusion.

However, in my first investigations I have followed the more traditional approach of assuming a preference ordering on interpretations, and to define the regular interpretations as the most preferred models of the evidence. This idea is motivated by the approaches of Helft (1989) and De Raedt & Bruynooghe (1993), who imp ment the completeness assumption regarding the evidence by the Closed World Assur a result, the characterisation of confirmatory induction stated below as Theor a close variant of Kraus et al.'s characterisation of preferential reasonin familiar with the latter paper may want to skip the present section, which is in form of completeness' sake. The implications of this correspondence between induction and nonmonotonic reasoning will be discussed in section 4 below.

tion. As h 3.14 is readers

DEFINITION 3.7. A preferential confirmatory structure is a triple $\langle S, l, \rangle$, where S is a set of *states*, l: $S \rightarrow M$ is a function that labels every te with a model, and < is a strict partial order on S, called the *prefere* сe ordering, that is smooth². A state $s \in S$ satisfies a formula $\alpha \in L$ iff l(s)α; the set of states satisfying α is denoted by $[\alpha]$, and a minimal element of $[\alpha]$ (wrt. <) will be called a *regular* state for α . The consequence relation defined by W is denoted by k_w and is defined by: $\alpha k_w \beta$ iff (i) there is a state $s \in S$ satisfying α , and (*ii*) every regular state for α satisfies β . A consequence relation is called preferential confirmatory iff it is defined by a preferential confirmatory structure.

In comparison with the preferential models of (Krau that in a confirmatory argument the evidence is re guarantee the validity of (I2). The intermediate needed for technical reasons, and can be interpreted considers possible in that epistemic state.

t al., 1990), the only difference is ired to be satisfiable, in order to hainly s the models the reasoning agent

LEMMA 3.8 (Soundness). Any preferential relation satisfies (I1–7, C1–5).

Proof. The proof is easy and will only be carried

out for (C1-5).

опріттиногу

(C1): if all regular states for α satisfy β and $\rightarrow \gamma$, then all regular states for α satisfy γ . (C2): suppose $[\alpha]$ is non-empty, and not all regular states for α satisfy $\neg\beta$; it follows that some state in S satisfies β . (C3): if all regular states for α satisfy β and γ , then clearly they satisfy $\beta \wedge \gamma$.

(C4): note that $[\alpha \lor \beta] = [\alpha] \cup [\beta]$; thus, if $[\alpha]$ and $[\beta]$ are non-empty then so is $[\alpha \lor \beta]$. Furthermore, the set of regular states for $\alpha \lor \beta$ is a subset of the

²I.e. for any $S' \subseteq S$ and for any $s \in S'$, either s is minimal in S', or there is a $t \in S'$ such that t < s and t is minimal in S'. This condition is satisfied if < does not allow infinite descending chains.

union of the regular states for α and β , since a state cannot be minimal in $[\alpha \lor \beta]$ without being minimal in at least one of $[\alpha]$ and $[\beta]$.

(C5): we need the fact that \langle is a smooth partial order. Suppose that $[\alpha]$ is non-empty, and all regular states for α satisfy β and γ — clearly, $[\alpha \land \beta]$ is non-empty. Now, let s be regular for $\alpha \land \beta$, then $s \in [\alpha]$; I will prove that s is regular for α . Suppose not, then there is a $t \in [\alpha]$ such that t < s and t is regular for α . Now, every state regular for α satisfies β , hence $t \in [\alpha \land \beta]$. But this contradicts the minimality of s in $[\alpha \land \beta]$, hence s is regular for α and thus satisfies γ .

As a matter of fact, (I1) and (I4–6) can be derived from (I2–3, I7, C –5). Another derived rule that will be used below is (2.4).

In order to prove completeness, we need to build a preferentia structure W from a given consequence relation \ltimes satisfying (I1–7, C1 $\alpha \in \beta$ iff $\alpha \in \beta$. As in the case of explanatory structures, such a confirm is built from a specific set of models. These models are selected relation formula, as follows.

), such that ory structure

DEFINITION 3.9. Let k be a conjectural consequence relation. The model $m \in U$ is said to be *normal for* α iff for all β in L such that $\alpha \in \beta$, m Β.

So, a model is normal for a formula if it satisfies every confirmed hypothesis. Thus, given certain evidence the set of normal models decreases when the set of confirmed hypotheses increases. Notice that every model in U is normal for an inadmissible formula, which is therefore not satisfied by some of its normal models. An admissible formula is satisfied by every normal model, however. Notice also that if α is admissible and γ is inadmissible, then by (C2) $\alpha \not\in \neg \gamma$, hence no normal model for α satisfies γ .

The set of models normal for admissi will be used to build e formula preferential confirmatory structure. The follow normal models.

LEMMA 3.10. Suppose a consequence re and let α be an admissible formula. All $\alpha \in \beta$.

Proof. The if part follows from Definition For the only-if part, suppose $\alpha \ltimes \alpha$ and normal model for α that does not satisfy suffices to show that Γ_0 is satisfiable. S there is a finite $\Delta \subseteq \{\delta \mid \alpha < \delta\}$ such that $\alpha \not\in \Delta \rightarrow \beta$. But by (C3) $\alpha \not\in \Delta$; using (C. contradiction.

ig lemma	ates the key result a
tion k sa	fies (C1) and (C3)
ormal mod	s for α satisfy β iff
9.	
кр, 1 wn	now that there is a
Let $\Gamma_0 = \{$	$\beta\} \cup \{\delta \mid \alpha \not\in \delta\}; \text{ it }$
pose not,	en by compactness
$\rightarrow \beta$, i.e.	$\alpha \rightarrow (\Delta \rightarrow \beta)$; by (C1)
-	, we obtain $\alpha \in \beta$, a

Notice from the proof of Lemma 3.10 that normal models exist for any admissible α .

Given a consequence relation \ltimes satisf based on a preferential confirmatory structu

- (1) $S = \{ \langle m, \alpha \rangle \mid \alpha \text{ is an admissible for } \}$
 - (2) $l(\langle m, \alpha \rangle) = m;$
 - (3) $\langle m, \alpha \rangle < \langle n, \beta \rangle$ iff $\alpha \lor \beta \vDash \alpha$ and *m* B.

Thus, states are pairs of admissible formulas and no simply maps a state to the model it contains. Co ordering between states: note that $\beta \ltimes \alpha$ is a special and the fact that α is admissible. The condition m irreflexive; note that as a consequence any $\langle m, \alpha \rangle \in S$ is minimal in [α].

g (I1-C1–5), the completeness proof is

nula, a *m* is a normal model for α };

> napennie runction ai mouen ition (3) defines the preference se of $\alpha \lor \beta \vDash \alpha$ by means of (C4), β is added to make the ordering

The main difference between the preferential consequence relations of Kraus et al. and my preferential confirmatory consequence relations is the way unsatisfiable formulas are treated. In Kraus et al.'s framework unsatisfiable formulas are characterised by the fact that they have every formula in L as a plausible consequence,

which means that they don't have normal models. In V framework, unsatisfiable formulas confirm no hypotheses, and have all models in V as normal models. In both cases, the structure W that is used to prove complete ass contains only satisfiable formulas in its states. This means that we can replicat most of Kraus *et al.*'s results about the structure W (KLM 5.13 refers to (Kraus et al.,

> is a strict partial order. $s \in [\alpha]$, either s is minimal

> > 11–7, C1–5), and

e is a state in S

so by Lemma 3.10

ble formula, *n* is a

ition 3.11 (4). By

ա/ելալ.

Β.

PROPOSITION 3.11. (1) (KLM 5.13) The relation (2) (KLM 5.15) The relation < is smooth: for an in $[\alpha]$ or there exists a state t<s minimal in $[\alpha]$ (3) (KLM 5.11) If $\alpha \lor \beta \ltimes \alpha$ and m is a normal del for α that satisfies β , then m is a normal model for β . (4) (*KLM* 5.14) $\langle m, \alpha \rangle$ is minimal in [β] iff $m \square \beta$ and $\alpha \lor \beta \vDash \alpha$.

The first two propositions express that W is a preferential confirmatory structure. remaining two are used in the proof of the following lemma.

LEMMA 3.12. Let k be a consequence relation satisfyin *let W be defined as above. If* $\alpha \in \beta$ *then* $\alpha \in W$ *become for the constant of a bound of the constant of Proof.* Suppose that $\alpha \in \beta$; we will show that (i) the satisfying α , and (*ii*) every minimal state in $[\alpha]$ satisfies (*i*) By (I4) α is admissible; furthermore, by (2.4) $\alpha \not < \neg$ there exists a model *m* normal for α . We conclude that (*ii*) Suppose $\langle n, \gamma \rangle$ is minimal in $[\alpha]$, then γ is an admis normal model for γ that satisfies α , and $\gamma \vee \alpha \not\in \gamma$ by Prop Proposition 3.11 (3) *n* is a normal model for α , hence *n*

The following lemma proves the converse of Lemma 3.12, and completes the proof of the representation theorem.

LEMMA 3.13. Let $k \in a$ consequence relation satisfying (I1-7, C1-5), and let W be defined as above. If $\alpha \ltimes_W \beta$ then $\alpha \ltimes \beta$.

Proof. Suppose $\alpha \ltimes_W \beta$, then α must be admissible (since no state in *S* satisfies an inadmissible formula). Furthermore, given any model *m* normal for α , $\langle m, \alpha \rangle$ is minimal in [α], hence *m* satisfies β , and the conclusion follows by Lemma 3.10.

We may now summarise.

THEOREM 3.14 (Representation theorem for preferential confirmatory consequence relations). A consequence relation is preferential confirmatory iff it satisfies (I1–7, C1–5).

Proof. The only-if part is Lemma 3.8. For the if part, let \notin be a consequence relation satisfying (I1–7, C1–5) and let *W* be defined as above. Lemmas 3.12 and 3.13 prove that $\alpha \notin \beta$ iff $\alpha \notin_W \beta$, i.e. \notin is preferential confirmatory.

4. Discussion

In the last part of this paper I will discuss the significance of the approach and results presented above and make a comparison with related work.

Although the proof-theoretic approach of section 2 is nicely balanced by a semantic analysis in section 3, I think the former — analysing the logical characteristics of induction on the level of consequence relations — constitutes the main contribution of the paper. This approch has of course been heavily influenced by work of Gabbay and Krass et al. (1990) on nonmonotonic consequence (1985), Makinson (19 ny approach demonstrates that their method is in relations and operation inything, o mention Gärdenfors' work on fact a me nodology. rationality ostulates fo been stated mainly as a starting point for further The ıs ha research atory induction to reversed this area. resting perspective is obtained by defining deduction An equivale U $\alpha \ltimes \beta \Leftrightarrow$ β) \subset *L* $\alpha) \subseteq Cn$

where $Cn(\alpha) = \{\gamma \mid \alpha \mid \gamma\}$ denotes the set of consequence or the *explanatory power* of α . In words: given evidence α , β is an explanatory hypothesis of the explanatory power of β exceeds that of α (without reaching inconsistency). This view of induction as *explanation-preserving* reasoning can be subsequently generalised by considering weaker (nonmonotonic) explanation mechanisms \uparrow .³

³As Daniel Lehmann observed, any inductive consequeynce relation thus defined will satisfy transitivity (Lehmann, personal communication). We are currently studying cumulative, preferential and rational explanation mechanisms (Kraus *et al.*, 1990; Lehmann & Magidor, 1992).

The preferential semantics for modelling the form of confirmatory induction considered by Hempel, Helft and De Raedt has of course been borrowed from work in nonmonotonic reasoning. Readers familiar with that field will have recognised some of the postulates — for instance, our stronger verification principle (C5) is known as cautious monotonicity. On one hand, the correspondence is not so surprising: both forms of reasoning rely on some form of completeness of the premisses, be it 'what you don't know is false' or 'what you don't know resembles what you know'. A pragmatic difference is that in nonmonotonic reasoning conclusions are usually specific ('it flies'), while in induction conclusions are typically general ('all crows are black'). On the other hand, induction often proceeds from incomplete evidence. A formalisation of weaker forms of confirmatory induction requires to relax the completeness assumption regarding the evidence, and thus to drop (C3) — but (C3) will be sound for any semantics that selects a subset of the models of the evidence, all of which should satisfy the hypothesis. A semantics that constructs a set of regular interpretations from the evidence, some of which do not satisfy the evidence (for instance because names have been permuted) seems a promising research direction.

Philosophically speaking, the view of induction as inference of explanations owes much to the ideas of Peirce. In fact, Peirce identified the formation of explanatory hypotheses with reversed deduction (Hartshorne *et al.*, 1931–58, Vol.1, p.117). The framework of rationality postulates provides a finer-grained perspective on the matter. Hempel was the first to analyse the logical relation of confirmation, and even proposed a number of rationality postulates or adequacy conditions, among which (variants of) I2', I6, I7, C1, C3, 2.3 and 2.4 (see Hempel, 1943). He then proceeded to develop a definition of confirmation, which comes surprisingly close to minimal Herbrand model semantics. However, Hempel never gave a complete axiomatisation of his definition.

In summary, we have proposed a new framework for thinking about the logic of induction, a problem that has been bothering philosophers for over twenty centuries, by bringing together old and recent work in philosophy, logic, and artificial intelligence.

References

- R. CARNAP (1950), Logical Foundations of Probability, Routledge & Kegan Paul, London.
- L. DE RAEDT & M. BRUYNOOGHE (1993), 'A theory of clausal discovery', in *Proc. 13th International Conference on Artificial Intelligence IJCAI'93*, Morgan Kaufmann, San Mateo: 1058–1063.
- P.A. FLACH (1995), *Conjectures an inquiry concerning the logic of induction*. PhD thesis, Tilburg University.
- D. GABBAY (1985), 'Theoretical foundations for non-monotonic reasoning in expert systems', in *Logics* and *Models of Concurrent Systems*, K.R. Apt (ed.), Springer-Verlag, Berlin: 439–457.
- P. GÄRDENFORS (1988), *Knowledge in Flux modeling the dynamics of epistemic states*, MIT Press, Cambridge, MA.
- C. HARTSHORNE, P. WEISS & A. BURKS, EDS (1931-58), *Collected Papers of Charles Sanders Peirce*, Harvard University, Cambridge.

- N. HELFT (1989), 'Induction as nonmonotonic inference', in *Proc. First International Conference on Knowledge Representation and Reasoning KR*'89, Morgan Kaufmann, San Mateo: 149–156.
- C.G. HEMPEL (1943), 'A purely syntactical definition of confirmation', *Journal of Symbolic Logic* **6**(4): 122–143.
- C.G. HEMPEL (1945), 'Studies in the logic of confirmation', *Mind* **54**(213): 1–26 (Part I); **54**(214): 97–121 (Part II).
- S. KRAUS, D. LEHMANN & M. MAGIDOR (1990), 'Nonmonotonic reasoning, preferential models and cumulative logics', *Artificial Intelligence* 44: 167-207.
- D. LEHMANN & M. MAGIDOR (1992), 'What does a Conditional Knowledge Base Entail?', *Artificial Intelligence* **55**(1): 1-60.
- D. MAKINSON (1989), 'General theory of cumulative inference', in *Proc. 2nd International Workshop on Non-Monotonic Reasoning*, M. Reinfrank, J. de Kleer, M.L. Ginsberg & E. Sandewall (eds.), Lecture Notes in Artificial Intelligence 346, Springer-Verlag, Berlin: 1–18.
- K. SCHLECHTA (1995), 'A reduction of the theory of confirmation to the notions of distance and measure', Proc. European Conference on Symbolic and Quantitative Approaches to Reasoning and Uncertainty ECSQARU '95, Lecture Notes in Artificial Intelligence 946, Springer-Verlag, Berlin: 387–394.