

On the logic of induction

Peter A. Flach*

INFOLAB, Tilburg University, PObox 90153, 5000 LE Tilburg, the Netherlands

Abstract.

This paper presents a logical analysis of induction. Contrary to common approaches to inductive logic that treat inductive validity as a real-valued generalisation of deductive validity, we argue that the only logical step in induction lies in hypothesis *formation* rather than evaluation. Inspired by the seminal paper of Kraus, Lehmann & Magidor [18] we analyse the logic of inductive hypothesis formation on the metalevel of consequence relations. Two main forms of induction are considered: explanatory induction, aimed at inducing a general theory explaining given observations, and confirmatory induction, aimed at characterising completely or partly observed models. Several sets of meta-theoretical properties of inductive consequence relations are considered, each of them characterised by a suitable semantics. The approach followed in this paper is extensively motivated by referring to recent and older work in philosophy, logic, and Machine Learning.

1. Introduction

This paper is an attempt to develop a logical account of inductive reasoning, one of the most important ways to synthesize new knowledge. Induction provides an idealized model for empirical sciences, where one aims to develop general theories that account for phenomena observed in controlled experiments. It also provides an idealized model for cognitive processes such as learning concepts from instances. The advent of the computer has suggested new inductive tasks such as program synthesis from examples of input-output behaviour and knowledge discovery in databases, and the application of inductive methods to Artificial Intelligence problems is an active research area, which has displayed considerable progress over the last decades.

On the foundational side, however, our understanding of the essentials of inductive reasoning is fragmentary and confused. Induction is usually defined as inference of general rules from particular observations, but this slogan can hardly count as a definition. Clearly some rules are better than others for given observations, while yet other rules are totally unacceptable. A logical account of induction should shed more light on the relation between observations and hypotheses, much like deductive logic formalises the relation between theories and their deductive consequences.

This is by no means an easy task, and anyone claiming to provide a definitive solution should be approached sceptically. The main contribution of this paper lies in the novel perspective that is obtained by combining older work in philosophy of science with a methodology suggested by recent work in formalising nonmonotonic reasoning. This perspective provides us with a *descriptive* — rather than prescriptive — account of induction, which clearly indicates both the opportunities for and limitations of logical analysis when it comes to modelling induction.

1.1 Problem formulation and approach

I should start by stressing that the study reported on in this paper should be perceived as an application of logical analysis to problems in Artificial Intelligence. Thus, we will take it for granted that there exists a distinct and useful form of reasoning called induction. As a model for this form of reasoning we may take the approaches to learning classification rules from examples that can be found in the Machine Learning literature, or the work on inducing Prolog programs and first-order logical theories from examples in the recently established discipline of Inductive Logic Programming. By taking this position we will avoid the

* E-mail: Peter.Flach@kub.nl.

controversies abounding in philosophy of science as to whether or not science proceeds by inductive methods. This is not to say that I will completely ignore philosophical considerations — in fact, my approach has been partly motivated by works from the philosophers Charles Sanders Peirce and Carl G. Hempel, as I will explain shortly.

The main question addressed in this paper is the following: *Can we develop a logical account of induction that is sufficiently similar to the modern account of deduction?* By ‘the modern account of induction’ I mean the by now standard approach, developed in the first half of this century, of defining a logical language, a semantical notion of deductive consequence, and a proof system of axioms and inference rules operationalising the relation of deductive consequence. By the stipulation that the logical account of induction be ‘sufficiently similar’ to this modern account of deduction I mean that the former should likewise consist of a semantical notion of *inductive consequence*, and a corresponding proof system.

Those perceiving logic as the ‘science of correct reasoning’ will now object that what I am after is a deductive account of induction, and it is known already since Hume that inductive hypotheses are necessarily defeasible. My reply to this objection is that it derives from a too narrow conception of logic. In my view, *logic is the science of reasoning*, and it is the logician’s task to develop formal models of every form of reasoning that can be meaningfully distinguished. In developing such formal models for nondeductive reasoning forms, we should keep in mind that deduction is a highly idealized and restricted reasoning form, and that we must be prepared to give up some of the features of deductive logic if we want to model reasoning forms that are less perfect, such as induction.

The fundamental question then is: which features are inherent to logic *per se*, and which are accidental to deductive logic? To illustrate this point, consider the notion of truth-preservation: whenever the premisses are true, the conclusion is true also. It is clear that truth-preservation must be given up as soon as we step out of the deductive realm. The question then arises whether a logical semantics is mainly a tool for assessing the truth of the conclusion given the truth of the premisses, or whether its main function is rather to define what property is preserved when passing from premisses to conclusion. We will address this and similar fundamental questions in this paper.

Another objection against the approach I propose could be that deductive logic is inherently prescriptive: it clearly demarcates the logical consequences one *should* accept on the basis of given premisses, from the ones one *should not* accept. Clearly, our understanding of induction is much too limited to be able to give a prescriptive account of induction. My reply to this objection is that, while such a demarcation is inherent to logic, its interpretation can be either prescriptive or descriptive. The inductive logics I propose in this paper distinguish between hypotheses one *should not* accept on the basis of given evidence, relative to a certain goal one wants the hypothesis to fulfil, and hypothesis one *might* accept. Put differently, these inductive logics formalise the logic of inductive hypothesis formation rather than hypothesis selection, which I think is the best one can hope to achieve by purely logical means.

The objective pursued in this paper, then, is to develop semantics and proof systems for inductive hypothesis formation. What is new here is not so much this objective, which has been pursued before (see e.g. [4]), but the meta-theoretical viewpoint taken in this paper, which I think greatly benefits our understanding of the main issues. This meta-theoretical viewpoint has been inspired by the seminal paper of Kraus, Lehmann & Magidor [18], where it is employed to unravel the fundamental properties of nonmonotonic reasoning. Readers familiar with the paper of Kraus *et al.* may alternatively view the present paper as a constructive proof of the thesis that their techniques in fact establish a *methodology*, by demonstrating how they can be successfully applied to analyse a rather different form of reasoning.

1.2 Plan of the paper

The paper is structured as follows. In section 2 the philosophical, logical, and Machine Learning backgrounds of this paper are surveyed. Section 3 introduces the main logical tool employed in this paper: the notion of a metalevel consequence relation. Sections 4 and 5 form the technical core of this paper, stating representation theorems characterising sets of metalevel properties of explanatory induction and confirmatory induction, respectively. In section 6 we discuss the implications of the approach taken and results obtained in this paper. Section 7 repeats the main conclusions.

2. Backgrounds

This section reviews a number of related approaches from the philosophical, logical, and Machine Learning literature. With such a complex phenomenon as induction, one cannot hope to give an overview that can be called complete in any sense — I will restrict attention to those approaches that either can be seen as precursors to my approach, or else are considered as potential answers to my objectives but rejected upon closer inspection. We start with the latter.

2.1 Inductive probability

By now it is commonplace to draw a connection between inductive reasoning and probability calculus. Inductive or subjective probability assesses the degree to which an inductive agent is willing to accept a hypothesis on the basis of available evidence. A so-called posterior probability of the hypothesis after observing the evidence is obtained by applying Bayes' theorem to the probability of the hypothesis prior to observation. Rudolf Carnap has advocated the view that inductive probability gives rise to a system of inductive logic [3]. Briefly, Carnap defines a function $c(H,E)$ assigning a *degree of confirmation* (a number between 0 and 1) to a hypothesis H on the basis of evidence E . This function generalises the classical notion of logical entailment — which can be seen as a 'confirmation function' from premisses and conclusion to $\{0,1\}$ — to an inductive notion of 'partial entailment':

'What we call inductive logic is often called the theory of nondemonstrative or nondeductive inference. Since we use the term 'inductive' in the wide sense of 'nondeductive', we might call it the theory of inductive inference... However, it should be noticed that the term 'inference' must here, in inductive logic, not be understood in the same sense as in deductive logic. Deductive and inductive logic are analogous in one respect: both investigate logical relations between sentences; the first studies the relation of [entailment], the second that of degree of confirmation which may be regarded as a numerical measure for a partial [entailment]... The term 'inference' in its customary use implies a transition from given sentences to new sentences or an acquisition of a new sentence on the basis of sentences already possessed. However, only deductive inference is inference in this sense.' [3, §44B, pp.205–6]

This citation succinctly summarises why inductive probability is *not* suitable, in my view, as the cornerstone of a logic of induction. My two main objections are the following.

Inductive probability treats all nondeductive reasoning as inductive. This runs counter to one of the main assumptions of this paper, namely that induction is a reasoning form in its own right, which we want to characterise in terms of properties it enjoys rather than properties it lacks. A more practical objection is that a single logical foundation for all possible forms of nondeductive reasoning is likely to be rather weak. Indeed, I would argue that in many forms of reasoning the goal that is to be fulfilled by the hypothesis, such as explaining the observations, is not reducible to a degree of confirmation.¹

Inductive probability, taken as partial entailment, leads to a degenerated view of logic. This is essentially what Carnap notes when he states that his inductive logic does not establish inference in the same sense as deductive logic (although he would not call it a degeneration). This means that, for instance, the notion of a proof reduces to a calculation of the corresponding degree of confirmation. A possible remedy is to define and axiomatise a *qualitative* relation of confirmation, such as the relation defined by $qc(H,E) \Leftrightarrow c(H,E) > c(H, \mathbf{true})$. However, such a qualitative relation of confirmation can also be postulated without reference to numerical degrees of confirmation, which would give us much more freedom to investigate the relative merits of different axiom systems. In fact, this is the course of action taken by Hempel, as we will see in the next section.

I should like to stress that it is not inductive probability or Bayesian belief measures as such which are criticised here — on the contrary, I believe these to be significant approaches to the important problem of how to update an agent's beliefs in the light of new information. Since belief measures express the agent's subjective estimates of the truth of hypotheses, let us say that inductive probability and related approaches establish a *truth-estimating procedure*. My main point is that such truth-estimating procedures are, generally speaking, *complementary* to logical systems. Truth-estimating procedures

¹Note that degree of confirmation is not a quantity that is simply to be maximised, since this would lead us straight back into deductive logic.

answer a type of question which nondeductive logical systems, in general, cannot answer, namely: how plausible is this hypothesis given this evidence? The fact that deductive logic incorporates such a truth-estimating procedure is accidental to deductive reasoning; the farther one moves away from deduction, the less the logical system has to do with truth-estimation. For instance, the gap between logical systems for nonmonotonic reasoning and truth-estimating procedures is much smaller than the gap between the latter and logical systems for induction. Indeed, one may employ the same truth-estimating procedure for very different forms of reasoning.

2.2 *Confirmation as a qualitative relation*

Carl G. Hempel [15, 16] developed a qualitative account of induction that will form the basis of the logical system for what I call confirmatory induction (section 5). Carnap rejected Hempel's approach, because he considered a quantitative account of confirmation as more fundamental than a qualitative account. However, as explained above I think that the two are conceived for different purposes: a function measuring degrees of confirmation can be used as a truth-estimating procedure, while a qualitative relation of confirmation can be used as the cornerstone for a logical system. I also consider the two as relatively independent: a qualitative confirmation relation that cannot be obtained from a numerical confirmation function is not necessarily ill-conceived, as long as the axioms defining the qualitative relation are intuitively meaningful.

Hempel's objective is to develop a material definition of confirmation. Before doing so he lists a number of adequacy conditions any such definition should satisfy. Such adequacy conditions can be seen as metalevel axioms, and we will discuss them at some length. The following conditions can be found in [16, pp.103–106, 110]; logical consequences of some of the conditions are also stated.

- (H1) *Entailment condition*: any sentence which is entailed by an observation report is confirmed by it.
- (H2) *Consequence condition*: if an observation report confirms every one of a class K of sentences, then it also confirms any sentence which is a logical consequence of K .
 - (H2.1) *Special consequence condition*: if an observation report confirms a hypothesis H , then it also confirms every consequence of H .
 - (H2.2) *Equivalence condition*: if an observation report confirms a hypothesis H , then it also confirms every hypothesis which is logically equivalent with H .
 - (H2.3) *Conjunction condition*: if an observation report confirms each of two hypotheses, then it also confirms their conjunction.
- (H3) *Consistency condition*: every logically consistent observation report is logically compatible with the class of all the hypotheses which it confirms.
 - (H3.1) Unless an observation report is self-contradictory, it does not confirm any hypothesis with which it is not logically compatible.
 - (H3.2) Unless an observation report is self-contradictory, it does not confirm any hypotheses which contradict each other.
- (H4) *Equivalence condition for observations*: if an observation report B confirms a hypothesis H , then any observation report logically equivalent with B also confirms H .

The entailment condition (H1) simply means that entailment 'might be referred to as the special case of *conclusive* confirmation' [16, p.107]. The consequence conditions (H2) and (H2.1) state that the relation of confirmation is closed under weakening of the hypothesis or set of hypotheses (H_1 is weaker than H_2 iff it is logically entailed by the latter). Hempel justifies this condition as follows [16, p.103]: 'an observation report which confirms certain hypotheses would invariably be qualified as confirming any consequence of those hypotheses. Indeed: any such consequence is but an assertion of all or part of the combined content of the original hypotheses and has therefore to be regarded as confirmed by any evidence which confirms the original hypotheses.' Now, this may be reasonable for single hypotheses (H2.1), but much less so for sets of hypotheses, each of which is confirmed separately. The culprit can be identified as (H2.3), which together with (H2.1) implies (H2). A similar point can be made as regards the consistency condition (H3), about which Hempel remarks that it 'will perhaps be felt to embody a too

severe restriction'. (H3.1), on the other hand, seems to be reasonable enough; however, combined with the conjunction condition (H2.3) it implies (H3).

We thus see that Hempel's adequacy conditions are intuitively justifiable, except for the conjunction condition (H2.3) and, *a fortiori*, the general consequence condition (H2). On the other hand, the conjunction condition can be justified by a completeness assumption on the evidence, as will be further discussed in section 5. We close this section by noting that Hempel's material definition of the relation of confirmation of a hypothesis by evidence roughly corresponds to what we would nowadays call 'truth of the hypothesis in the truth-minimal Herbrand model of the evidence'. We will return to material definitions of qualitative confirmation in section 2.5.

2.3 *Abduction*

Predating Hempel's work on confirmation by almost half a century is the work of Charles Sanders Peirce on abduction: the process of forming explanatory hypotheses, which I will briefly discuss in this section.

In a series of lectures on Pragmatism delivered in 1903, Peirce distinguishes three types of reasoning: deduction, induction, and abduction. Induction 'consists in starting from a theory, deducing from it predictions of phenomena, and observing those phenomena in order to see *how nearly* they agree with the theory'. Furthermore,

'The justification for believing that an experiential theory which has been subjected to a number of experimental tests will be in the near future sustained about as well by further such tests as it has hitherto been, is that by steadily pursuing that method we must in the long run find out how the matter really stands.' [13, 5.170]

Note that Peirce claims, like Carnap, that induction evaluates the plausibility of a given theory, rather than constructing that theory from observations. However, inductive hypotheses do not come out of the blue, and this is where abduction comes into play:

'Abduction is the process of forming an explanatory hypothesis. It is the only logical operation which introduces any new idea; for induction does nothing but determine a value, and deduction merely evolves the necessary consequences of a pure hypothesis.

Deduction proves that something *must* be; Induction shows that something *actually is* operative; Abduction merely suggests that something *may* be.

Its only justification is that from its suggestion deduction can draw a prediction which can be tested by induction, and that, if we are ever to learn anything or to understand phenomena at all, it must be by abduction that this is to be brought about.

No reason whatsoever can be given for it, as far as I can discover; and it needs no reason, since it merely offers suggestions.' [13, 5.171]

In other words, *abduction is the process of conjecturing inductive hypotheses*, constrained by the requirement that they should comply with the available observations. Abduction represents the purely logical part of inductive reasoning.²

Peirce proceeds by defining the logical form of abduction. 'It must be remembered', he writes, 'that abduction, although it is very little hampered by logical rules, nevertheless is logical inference, asserting its conclusion only problematically or conjecturally, it is true, but nevertheless having a perfectly definite logical form.' Peirce then defines this logical form, as follows.

²Unfortunately, the term 'abduction' is nowadays used in two different ways. Peirce himself is to blame at least partly for this confusion, sine he first proposed a rather different, syllogistic classification of reasoning forms, which can be summarized as follows. Consider the Aristotelian syllogism *Barbara*: 'All the beans from this bag are white; these beans are from this bag; therefore, these beans are white'. Now there are two ways to exchange the conclusion with one of the premisses, one resulting in the inductive syllogism 'These beans are white; these beans are from this bag; therefore, all the beans from this bag are white', the other in 'All the beans from this bag are white; these beans are white; therefore, these beans are from this bag'. Peirce refers to this latter syllogism as (forming a) *hypothesis*. This syllogistic theory has to a large extent been adopted in the discipline of logic programming, where abduction (ironically, the term was only introduced in Peirce's later theory) is generally perceived as the inference of ground facts from rules and a query that is to be explained. Notice that a logic based on entailment rather than syllogisms is unable to distinguish between the two latter syllogisms, which both embody a form of reversed deduction. See [10].

‘Long before I first classed abduction as an inference it was recognized by logicians that the operation of adopting an explanatory hypothesis — which is just what abduction is — was subject to certain conditions. Namely, the hypothesis cannot be admitted, even as a hypothesis, unless it be supposed that it would account for the facts or some of them. The form of inference, therefore, is this:

The surprising fact, C, is observed;
 But if A were true, C would be a matter of course,
 Hence, there is reason to suspect that A is true.

Thus, A cannot be abductively inferred, or if you prefer the expression, cannot be abductively conjectured until its entire content is already present in the premiss, “If A were true, C would be a matter of course.”’

[13, 5.188]

In short, the view of induction that Peirce offers here is this. Inductive reasoning consists of two steps: (i) formulating a conjecture, and (ii) evaluating the conjecture. Both steps take the available evidence into account, but in quite different ways and with different goals. The first step requires that the conjectured hypothesis *explains* the observations; having a definite logical form, it represents a form of inference. The second step evaluates how well predictions offered by the hypothesis agree with reality; it is not inference, but assigns a numerical value to a hypothesis. In order to avoid terminological problems, I will not use Peirce’s terminology and refer to the first step as explanatory hypothesis formation, and to the second as hypothesis evaluation or validation.

Leaving a few details aside, Peirce’s definition of explanatory hypothesis formation can be formalised as the inference rule

$$\frac{C, \text{---} C}{\text{---}}$$

In this paper I propose to generalise Peirce’s definition by including the relation of ‘is explained by’ as a parameter. This is achieved by lifting the explanatory inference from C to A to the metalevel, as follows:

$$\frac{A \text{---} C}{C \text{---} A}$$

The symbol --- stands for the explanatory consequence relation. Axiom systems for this relation will be considered in section 4.

2.4 Confirmation vs. explanation

We now have encountered two fundamental notions that play a role in inductive hypothesis formation: one is that the hypothesis should be *confirmed* by the evidence, the other that the hypothesis should *explain* the evidence. Couldn’t we try and build the requirement that the hypothesis be explanatory into our definition of confirmed hypothesis?

The problem is that an unthoughtful combination of explanation and confirmation can easily lead into paradox. Let H_1 and H_2 be two theories such that the latter includes the former, in the sense that everything entailed by H_1 is also entailed by H_2 . Suppose E is confirming evidence for H_1 ; shouldn’t we conclude that it confirms H_2 as well? To borrow an example of Hempel: ‘Is it not true, for example, that those experimental findings which confirm Galileo’s law, or Kepler’s laws, are considered also as confirming Newton’s law of gravitation?’ [16, p.104]. This intuition is formalised by the following condition:

- (H5) *Converse consequence condition*: if an observation report confirms a hypothesis H , then it also confirms every formula logically entailing H .

The problem is, however, that this rule is incompatible with the special consequence condition (H2.1). This can be seen as follows: in order to demonstrate that E confirms H for arbitrary E and H , we note that E confirms E by (H1), so by the converse consequence condition E confirms $E \wedge H$; but then E confirms H by (H2.1). Thus, we see that the combination of two intuitively acceptable conditions leads to a collapse of the system into triviality, a clearly paradoxical situation.

Hempel concludes that one cannot have both (H2.1) and (H5), and drops the latter. His justification of

this decision is however unconvincing, which is not surprising since neither is *a priori* better than the other: *they formalise different intuitions*. While (H2.1) formalises a reasonable intuition about confirmation, (H5) formalises an equally reasonable intuition about *explanation*:

- (H5') if an observation report is explained by a hypothesis H , then it is also explained by every formula logically entailing H .

In this paper I defend the position that Hempel was reluctant to take, namely that with respect to inductive or scientific hypothesis formation there is more than one possible primitive notion: the relation 'confirms' between evidence and hypothesis, and the relation 'is explained by'. Each of these primitives gives rise to a specific form of induction. This position is backed up by recent work in Machine Learning, to which we will turn now.

2.5 Inductive Machine Learning

Without doubt, the most frequently studied induction problem in Machine Learning is concept learning from examples. Here, the observations take the form of descriptions of instances (positive examples) and non-instances (negative examples) of an unknown concept, and the goal is to find a definition of the concept that correctly discriminates between instances and non-instances. Notice that this problem statement is much more concrete than the general description of induction as inference from the particular to the universal: once the languages, in which instances and concepts are described, are fixed, the desired relation between evidence and hypothesis is determined. A natural choice is to employ a predicate for the concept to be learned, and to use constants to refer to instances and non-instances. In this way, a classification of an instance can be represented by a truthvalue, which can be obtained by setting up a proof.³ We then obtain the following general problem definition:

- Problem:* Concept learning from examples
- Given:*
- (1) A predicate-logical language
 - (2) A predicate representing the target concept.
 - (3) Two sets P and N of ground literals of this predicate, representing the positive and negative examples.
 - (4) A background theory T containing descriptions of instances.
- Determine:* A hypothesis H within the provided language such that
- (i) for all $p \in P$: $T \cup H \models p$;
 - (ii) for all $n \in N$: $T \cup H \not\models n$.

Notice that condition (ii) is formulated in such a way that the hypothesis only needs to contain *sufficient conditions* for concept membership (since a negative classification is obtained by negation as failure). This suggests an analogy between concept definitions and Horn clauses, which can be articulated by allowing (possibly recursive) logic programs as hypotheses and background knowledge, leading us into the field of Inductive Logic Programming (ILP) [22]. Furthermore, P and N may contain more complicated formulae than ground facts. The general problem statement then becomes: given a partial logic program T , extend it with clauses H such that every formula in P is entailed and none of the formulae in N .

The potential for inductive methods in Artificial Intelligence is however not exhausted by classification-oriented approaches. Indeed, it seems fair to say that most knowledge implicitly represented by extensional databases is non-classificatory. Several researchers have begun to investigate non-classificatory approaches to knowledge discovery in databases. For instance, in previous work I have demonstrated that the problem of inferring the set of functional and multivalued attribute dependencies satisfied by a database relation can be formulated as an induction problem [6, 7, 8]. Furthermore, De Raedt & Bruynooghe have generalized the classificatory ILP-setting in order to induce non-Horn clauses from ground facts [5]. Both approaches essentially employ the following problem statement.

³The alternative is to represent concepts by open formulae, and to operationalize classification by means of subsumption.

Problem: Non-classificatory induction.

Given: (1) A predicate-logical language.
(2) Evidence E .

Determine: A hypothesis H within the provided language such that:
(i) H is true in a model m_0 constructed from E ;
(ii) for all g within the language, if g is true in m_0 then $H \supset g$.

Essentially, the model m_0 employed in the approaches by De Raedt & Bruynooghe and myself is the truth-minimal Herbrand model of the evidence.⁴ The hypothesis is then an axiomatisation of all the statements true in this model, including non-classificatory statements like ‘everybody is male or female’ and ‘nobody is both a father and a mother’.

The relation between the classificatory and non-classificatory approaches to induction is that they both aim at extracting similarities from examples. The classificatory approach to induction achieves this by constructing a single theory that entails all the examples. In contrast, the non-classificatory approach achieves this by treating the examples as a model — a description of the world that may be considered *complete*, at least for the purposes of constructing inductive hypotheses. The approach is justified by the assumption that *the evidence expresses all there is to know about the individuals in the domain*. Such a completeness assumption is reminiscent of the Closed World Assumption familiar from deductive databases, logic programming, and default reasoning — however, in the case of induction its underlying intuition is quite different. As Nicolas Helft, one of the pioneers of the non-classificatory approach, puts it:

‘induction assumes that the similarities between the observed data are representative of the rules governing them (...). This assumption is like the one underlying default reasoning in that a priority is given to the information present in the database. In both cases, some form of “closing-off” the world is needed. However, there is a difference between these: loosely speaking, while in default reasoning the assumption is “what you are not told is false”, in similarity-based induction it is “what you are not told looks like what you are told”.’ [14, p.149]

There is a direct connection between Peirce’s conception of abduction as formation of explanatory hypotheses and the classificatory induction setting, if one is willing to view a theory that correctly classifies the examples as an *explanation* of those examples. In this paper I suggest to draw a similar connection between Hempel’s conception of confirmation as a relation between evidence and potential hypotheses and the non-classificatory induction setting outlined above. Non-classificatory induction aims at constructing hypotheses that are *confirmed* by the evidence, without necessarily explaining them. Rather than studying material definitions of what it means to explain or be confirmed by evidence, as is done in the works referred to above, in the following sections I will be concerned with the logical analysis of the abstract notions of explanation and confirmation.

3. Inductive consequence relations

In the sections to follow I employ the notion of a consequence relation, originating from Tarski [24] and further elaborated by Gabbay [11], Makinson [20], and Kraus, Lehmann & Magidor [18, 19]. In this section I give an introduction to this important metalogical tool that is largely self-contained. The basic definitions are given in section 3.1. In section 3.2 I consider some general properties of consequence relations that arise when modelling the reasoning behaviour of inductive agents. Section 3.3 is devoted to a number of considerations regarding the pragmatics of consequence relations in general, and inductive consequence relations as used in this paper in particular.

3.1 Consequence relations

We distinguish between the language L in which an inductive agent formulates premisses and conclusions of inductive arguments, and the metalanguage in which statements about the reasoning behaviour of the

⁴An alternative approach is to consider the information-minimal partial model of the evidence [8].

inductive agent are expressed in this paper. Let L be a propositional language⁵ over a fixed countable set of proposition symbols, closed under the usual logical connectives. We assume a set of propositional models U , and a satisfaction relation $\models \subseteq U \times L$ that is well-behaved with respect to the logical connectives and compact. As usual, we write $\alpha \models \beta$ for $\forall u \in U. u \models \alpha \Rightarrow u \models \beta$, for arbitrary $\alpha, \beta \in L$. Note that U may be a proper subset of the set of all truth-assignments to proposition symbols in L , which would reflect prior knowledge or *background knowledge* of the inductive agent. Equivalently, we may think of U as the set of models of an implicit background theory T , and let $\alpha \models \beta$ stand for ' α is a logical consequence of T '.

The metalanguage is a restricted predicate language built up from a unary metapredicate \vdash in prefix notation (standing for validity with respect to U in L) and a binary metapredicate \prec in infix notation (standing for inductive consequence). In referring to object-level formulae from L we employ a countable set of metavariables $\alpha, \beta, \gamma, \delta, \dots$, the logical connectives from L (acting like function symbols on the metalevel), and the metaconstants **true** and **false**. Formulae of the metalanguage, usually referred to as *rules* or *properties*, are of the form $P_1, \dots, P_n \vdash Q$ for $n \geq 0$, where P_1, \dots, P_n and Q are literals (atomic formulae or their negation). Intuitively, such a rule is interpreted as an implication with *antecedent* P_1, \dots, P_n (interpreted conjunctively) and *consequent* Q , in which all variables are implicitly universally quantified. An example of such a rule, written in an expanded Gentzen-style notation, is

$$\frac{\alpha \wedge \beta \rightarrow \gamma, \alpha \prec \beta}{\alpha \wedge \neg \gamma \prec \beta}$$

This is a rule with two positive literals in its antecedent, and a negative literal in its consequent. Intuitively, it expresses that an inductive hypothesis β , previously inferred from evidence α , should be withdrawn if the negation of a consequence of α and β together is added to the evidence.

Consequence relations provide the semantics for this metalanguage, by fixing the meaning of the metapredicate \prec . Formally, a consequence relation is a subset of $L \times L$. They will be used to model part or all of the reasoning behaviour of a particular reasoning agent, by listing a number of *arguments* (pairs of *premiss* and *conclusion*) that the agent is assumed to accept⁶. A consequence relation satisfies a rule whenever it satisfies all instances of the rule, and violates it otherwise, where an instance of a rule is obtained by replacing the variables of the rule with formulae from L . A consequence relation satisfies an instance of a rule if, whenever it satisfies the literals in the antecedent of the rule, it also satisfies the consequent. Finally:

- (i) a literal α is satisfied whenever the propositional formula from L denoted by α is true in every model in U ;
- (ii) a literal $\neg \alpha$ is satisfied whenever the propositional formula from L denoted by α is false in some model in U ;
- (iii) a literal $\alpha \prec \beta$ is satisfied whenever the pair of propositional formulae from L denoted by α and β is an element of the consequence relation;
- (iv) a literal $\alpha \not\prec \beta$ is satisfied whenever the pair of propositional formulae from L denoted by α and β is not an element of the consequence relation;

For instance, the consequence relation $\{(p, q), (p \wedge \neg p, q)\}$ violates the rule above. We will often refer to a particular consequence relation as \prec and write $p \prec q$ instead of $\langle p, q \rangle \in \prec$.

A useful analogy with clausal logic is revealed by noting that the object language L establishes a Herbrand universe built up from the proposition symbols in L as constants and the connectives from L as function symbols. Consequence relations then correspond to Herbrand interpretations⁷ of the metalanguage, whose rules can be easily transformed to clausal notation. The following terminology is borrowed from logic programming:

- (i) if all of P_1, \dots, P_n and Q are positive literals the rule is *definite*;
- (ii) if at least one of P_1, \dots, P_n is a negative literal and Q is a positive literal the rule is *indefinite*;
- (iii) if all of P_1, \dots, P_n are positive literals and Q is a negative literal the rule is a *denial*.⁸

⁵One may argue that in induction the distinction between statements about individuals (ground facts) and statements about sets of individuals (universal sentences) is crucial, which calls for a predicate language. This point is discussed in section 6.

⁶This is not to say that the agent actually draws the conclusion when observing the premisses, but merely that she considers it one of the possible conclusions.

⁷These Herbrand interpretations are restricted to the metapredicate \vdash ; \vdash is treated as a built-in predicate.

⁸This exhausts all the possibilities: the case that at least one of P_1, \dots, P_n is a negative literal and Q is a negative

For instance, all the rule systems of [18] are made up of definite rules only. Logic programming theory teaches us that the set of consequence relations satisfying a set of definite rules is closed under intersection. In contrast, the rule system **R** characterised in [19] contains one indefinite rule, *viz.* Rational Monotonicity. Consequently, the set of rational consequence relations is not closed under intersection.

3.2 Some properties of inductive consequence relations

After having stated the main definitions concerning consequence relations I will now list some properties generally obeyed by inductive consequence relations. In this section we will distinguish between explanatory or confirmatory induction, and simply interpret $\alpha \prec \beta$ as 'evidence α '.

The first two rules state that the logical form of evidence and hypothesis is immaterial:

Left Logical Equivalence

$$\frac{\alpha \leftrightarrow \beta, \alpha \prec \gamma}{\beta \prec \gamma}$$

Right Logical Equivalence

$$\frac{\alpha \prec \gamma, \alpha \prec \beta}{\alpha \prec \beta}$$

Left Logical Equivalence states, for instance, that if the evidence is expressed as a conjunction of ground facts, the order in which they occur is immaterial. From the viewpoint of algorithms this may embody a considerable simplification — however, the framework of the present paper is intended to provide a model for inductive reasoning in general rather than particular algorithms.⁹

The following two rules express principles well-known from philosophy of science:

Verification

$$\frac{\alpha \wedge \beta \rightarrow \gamma, \alpha \prec \beta}{\alpha \wedge \gamma \prec \beta}$$

Falsification

$$\frac{\alpha \prec \beta, \alpha \wedge \neg \gamma \prec \beta}{\alpha \prec \beta}$$

In these two rules γ is a *prediction* made on the basis of hypothesis β and evidence α . Verification expresses that if such a prediction is indeed observed, hypothesis β remains a possible hypothesis, while if its negation is observed, β may be considered refuted according to Falsification.¹⁰ One might remark that typically the inductive hypothesis will entail the evidence, so that the first condition in the antecedent of Verification and Falsification may be simplified to $\beta \rightarrow \gamma$. However, this is only the case for certain approaches to explanatory induction; generally speaking inductive hypotheses, in particular those that are confirmed without being an explanation, may not contain all the information conveyed by the evidence. The formulation above represents the general case.

Falsification can be simplified in another sense as shown by the following lemma:

LEMMA 3.1. *In the presence of Left Logical Equivalence, Falsification is equivalent to the following rule:*

Consistency

$$\frac{\alpha \prec \beta, \beta \rightarrow \neg \alpha}{\alpha \prec \beta}$$

Proof. To derive Falsification, suppose $\alpha \wedge \beta \rightarrow \gamma$, i.e. $\beta \rightarrow \neg(\alpha \wedge \neg \gamma)$, then by Consistency $\alpha \wedge \neg \gamma \prec \beta$. To derive Consistency from Falsification, suppose $\alpha \prec \beta$ and $\beta \rightarrow \neg \alpha$, i.e.

literal can be rewritten to (i) or (ii).

⁹Practical algorithms establish a *function* from evidence to hypothesis rather than a relation, i.e. also Right Logical Equivalence would be invalidated by an induction algorithm (of all the logically equivalent hypotheses only one would be output).

¹⁰Notice that, contrary to the previous two rules, Verification and Falsification happen to be meaningful also when modelling the behaviour of an induction algorithm: Verification expresses that the current hypothesis should not be abandoned when the next observation is a predicted one (in the terminology of [1] the algorithm is conservative), while Falsification expresses that the current hypothesis must be abandoned when the next observations runs counter to the predictions of the algorithm (called consistency by [1]). However, in the context of the present paper these are not the intended interpretations of the two rules.

$\alpha \wedge \beta \rightarrow \text{false}$, then by Falsification $\alpha \wedge \neg \text{false} \not\prec \beta$, and by Left Logical Equivalence $\alpha \not\prec \beta$, a contradiction.

Falsification and Consistency rule out inconsistent evidence and hypotheses. The way inconsistent evidence is handled is merely a technicality, and we might have decided to treat it differently — for instance, Hempel’s entailment condition (H1) implies that in his framework inconsistent evidence confirms arbitrary hypotheses. The case of inconsistent hypotheses is different however: it is awkward to say, for instance, that arbitrary evidence induces an inconsistent hypothesis. Furthermore, in inductive concept learning often negative examples are included, that are not to be classified as belonging to the concept, which requires consistency of the induced rule. Also, the adoption of Consistency is the only way to treat explanatory and confirmatory induction in a unified way as regards the consistency of evidence and hypothesis.

In the presence of Consistency a number of other principles have to be formulated carefully. For instance, we have reflexivity only for consistent formulae. In the light of Consistency a formula is consistent if it occurs in an inductive argument, either as evidence or as hypothesis, so we have the following weaker versions of reflexivity:¹¹

Left Reflexivity

Right Reflexivity

$$\frac{\alpha \not\prec \beta}{\alpha \not\prec \alpha}$$

$$\frac{\alpha \not\prec \beta}{\alpha \not\prec \beta}$$

If a consequence relation contains an argument $\alpha \not\prec \alpha$, this signals that α is consistent with the reasoner’s background theory. We will call such an α *admissible* (with respect to the consequence relation), and use conditions of this form whenever we require consistency of evidence

The final rule mentioned in this section is a variant of Verification that allows to add any prediction to the hypothesis rather than the evidence:

Right Extension

$$\frac{\alpha \wedge \beta \rightarrow \gamma, \alpha \not\prec \beta}{\alpha \not\prec \beta \wedge \gamma}$$

Further rules considered in this paper are specific to either explanatory or confirmatory induction, and are therefore to be discussed in later sections.

3.3 *The pragmatics of consequence relations*

Before moving to the technical results of the paper we may spend a few thoughts on the exact nature of consequence relations. As defined above, a consequence relation is an extensional specification of the behaviour of a reasoning agent. The symbol \prec is introduced in order to reason about consequence relations and reasoning behaviour, and functions as a binary predicate in the metalanguage. Rules describing properties of \prec express boundaries of rationality of inductive reasoning: a consequence relation violating such a rule would be considered irrational.

Now suppose \mathbf{X} is a set of such rationality postulates, and let A be a set of inductive arguments. Clearly, by means of \mathbf{X} we could derive additional inductive arguments A' , the significance of which is: if an inductive agent accepts arguments A , and it behaves rationally according to \mathbf{X} , it should also accept arguments A' . For instance, if \mathbf{X} includes the rule of Verification, A contains the argument

`chevy_is_black` \prec `crows_are_black`

and the background knowledge includes

`crows_are_black` \rightarrow `chevy_is_black`

then A' contains the additional inductive argument

¹¹One might argue that induction is inherently non-reflexive if the hypothesis is to generalise the evidence. This point will be taken up in section 6.

¹²Readers with a background in inductive learning may interpret $\alpha \not\prec \alpha$ as ‘hypothesis α does not cover any negative example’.

$\text{chad_is_black} \wedge \text{chevy_is_black} \vdash \text{crows_are_black}$

implying that if the agent would not accept the latter, it would behave irrationally (wrt. \mathbf{X}).

One may now ask: what is the smallest set of arguments containing A and satisfying the rules of \mathbf{X} ? Such a set, if it exists, would represent the *closure* of A under \mathbf{X} , denoted $A^{\mathbf{X}}$. The significance of such a closure is that if two agents start from the same set of arguments A , they cannot possibly disagree about any other argument if they both act rationally according to \mathbf{X} . Clearly, the existence of such a closure operation depends on the rules in \mathbf{X} . As has been remarked before, if all the rules in \mathbf{X} are definite (having only positive literals in their antecedents and consequent) the closure of A under \mathbf{X} is unique for arbitrary A . All the rule systems in [18] consist solely of definite rules.

However, the situation changes drastically if not all rules are definite. In particular, rule systems containing indefinite rules (having at least one negative literal in their antecedent and a positive literal as consequent) will not have an associated closure operation. An indefinite rule represents a rationality principle of the following kind: ‘if you accept this argument, you should accept at least one of those’. One example is Rational Monotonicity as studied in [19]; other examples will be found in this paper. The upshot of such rules is that even if two agents agree on the initial set of arguments they accept, they can disagree about some other arguments without violating the rationality postulates. Put differently, if we have two consequence relations both satisfying the rules of \mathbf{X} , and we take their intersection to find out on what arguments they agree, this intersection itself may not satisfy the rules of \mathbf{X} .

Lehmann & Magidor are not content with this indefiniteness: they define an operation of *rational closure* which selects, among the many supersets of A that satisfy \mathbf{X} , one that has certain desirable properties [19, p.33]. This seems to be motivated by the tacit assumption that *rationality postulates \mathbf{X} should lead to a closure operation $A \rightarrow A^{\mathbf{X}}$* . Such a closure operation can be seen as a consequence operation in the sense of Tarski [24], mapping any set of premisses to the set of its consequences under some inference system. However, the notion of inference represented by such closure operations should be clearly distinguished from the notion of inference represented by consequence relations — if the latter operate on a metalevel, closure operations establish a *meta-metalevel*.

In this paper we will not be concerned with inference on this meta-metalevel. This choice is motivated on methodological rather than technical grounds: I don’t believe that rationality postulates for induction necessarily establish an unequivocal closure operation. If we intersect two inductive consequence relations to see what two inductive agents agree upon, the resulting set of arguments may, as a consequence relation, not satisfy some rationality postulate, but this is, it seems to me, just to be expected.

A further criticism of the meta-metalevel view is that it obscures the status of the consequence relation symbol. The relevant question on the meta-metalevel is: what expressions of the form $\alpha \vdash \beta$ are entailed by a set of expressions of the same form? In other words: on the meta-metalevel \vdash acts as a *connective*.¹³ This is also apparent from Kraus, Lehmann & Magidor’s terminology: they call $\alpha \vdash \beta$ a *conditional assertion*, a set of such conditional assertions is called a *conditional knowledge base*, and they ask themselves the question: *what does a conditional knowledge base entail?*, i.e. what other conditional assertions can be deduced from it? It seems to me that the meta-metalevel perspective is at odds with the metalevel perspective — indeed, it is my conjecture that the theory of entailment of conditional assertions can be developed without reference to an intermediate level of consequence relations.

In this paper the notion of closure will be employed on the metalevel rather than the meta-metalevel, in order to compare consequence relations and rule systems. Given a consequence relation \vdash , its *closure* $C_{\vdash}: L \rightarrow 2^L$ is defined as $C_{\vdash}(\alpha) = \{\beta \mid \alpha \vdash \beta\}$, and $C_{\vdash}(\alpha)$ is referred to as the closure of α under \vdash . A consequence relation \vdash_1 is called (at least) *as restrictive as* another consequence relation \vdash_2 if for every $\alpha \in L$ the closure of α under \vdash_1 is a subset of the closure of α under \vdash_2 — or equivalently, if \vdash_1 is a subset of \vdash_2 — and *more restrictive than* \vdash_2 if in addition $\vdash_1 \neq \vdash_2$. A set of rules \mathbf{X}_1 is (at least) *as restrictive as* another set of rules \mathbf{X}_2 if for every consequence relation \vdash_2 satisfying \mathbf{X}_2 there is a unique least restrictive consequence relation \vdash_1 satisfying \mathbf{X}_1 such that \vdash_1 is as restrictive as \vdash_2 ; we say that \vdash_1 is the \mathbf{X}_1 -restriction of \vdash_2 . In addition the mapping from \vdash_2 to its \mathbf{X}_1 -restriction is required to be a mapping from the set of relations satisfying \mathbf{X}_1 (see [9] for further motivation and analysis of these definitions).

The preceding definitions reflect that rule systems should be compared by comparing the set of conclusions for given premisses, rather than by metalevel entailment. From [KLM90]

¹³This raises the question why conditional assertions cannot be nested, as in $(\alpha \vdash \beta) \vdash \gamma$. Note that the answer to this question is perfectly clear on the metalevel, since this expression makes as little sense as, say, $(\alpha \vdash \beta) \vdash \gamma$.

impression that monotonic consequence is stronger (more restrictive) than preferential consequence because the rule system **M** entails every rule in the rule system **P** (or equivalently, every monotonic consequence relation is preferential). However, this does not work in general: the metalevel axiom $\alpha \vdash \beta$ entails all the rules in **M**, yet defines a very unrestrictive form of reasoning. Furthermore, our criterion also allows to compare rule systems that are not related by metalevel entailment, as we will see below.

I will now proceed with a technical analysis of the process of forming an explanatory hypothesis from evidence *à la* Peirce (section 4) and the process of forming a confirmed hypothesis *à la* Hempel (section 5).

4. Explanatory induction

In this section we will study abstract properties and semantics for explanatory consequence relations. Throughout the section $\alpha \prec \beta$ is to be read as ‘evidence α is explained by hypothesis β ’ or ‘hypothesis β is a possible explanation of evidence α ’. What counts as a possible explanation will initially be left unspecified — the framework of consequence relations allows us to formulate abstract properties of hypothesis formation, without fixing a particular material definition. We will then single out a particular set of properties (the system **EM**) and characterise it semantically by means of strong explanatory structures.

4.1 Properties of explanatory induction

A natural requirement for explanatory induction is that every consistent condition that hypothesis β be evidence counts as a possible explanation. As explained above, the condition that hypothesis β be consistent is expressed by $\beta \prec \beta$, which gives us the following rule:

Admissible Converse Entailment

$$\frac{\neg \alpha, \beta \prec \beta}{\alpha \prec \beta}$$

Another requirement for explanations has been discussed above as (H5’): possible explanations may be logically strengthened, as long as they remain consistent. This is expressed as follows:

Admissible Right Strengthening

$$\frac{\alpha \prec \beta, \beta \prec \gamma}{\alpha \prec \gamma}$$

We may note that Admissible Converse Entailment can be derived from Admissible Right Strengthening if we assume Consistency and the following rule:

Explanatory Reflexivity

$$\frac{\alpha \prec \alpha, \neg \beta \prec \alpha}{\beta \prec \beta}$$

This rule represents a property especially tailored for explanatory induction. It is understood by rewriting it into its contrapositive: from $\alpha \prec \alpha$ and $\beta \prec \beta$ infer $\neg \beta \prec \alpha$, which states that β is inadmissible, i.e. too strong a statement with regard to the background knowledge, its negation is so weak that it is explained by arbitrary admissible hypotheses α .

LEMMA 4.1. *In the presence of Consistency and Explanatory Reflexivity, Admissible Right Strengthening implies Admissible Converse Entailment.*

Proof. Suppose $\neg \alpha$, then by Consistency $\neg \alpha \prec \beta$. Suppose furthermore $\beta \prec \beta$, then by Explanatory Reflexivity $\alpha \prec \alpha$. The desired result follows by Admissible Right Strengthening.

While the rules above express properties of possible explanations, the following two rules concentrate on the evidence. The underlying idea is a basic principle in inductive learning: if the evidence is a set of instances of the concept, we can partition the evidence arbitrarily and find a single hypothesis that is an explanation of each subset of instances. This principle is established by the following two rules:¹⁴

¹⁴In previous work [8] Incrementality was called Additivity, and Convergence was called Incrementality. The

Incrementality

$$\frac{\beta \not\prec \gamma, \beta \prec \gamma}{\alpha \wedge \beta \prec \gamma}$$

Convergence

$$\frac{\alpha \rightarrow \beta, \alpha \prec \gamma}{\beta \prec \gamma}$$

LEMMA 4.2. *If \prec is a consequence relation satisfying Incrementality and Convergence, then $\alpha \wedge \beta \prec \gamma$ iff $\alpha \prec \gamma$ and $\beta \prec \gamma$.*

Proof. The *if* part is Incrementality, and the *only-if* part follows from Convergence.

Incrementality and Convergence are of considerable importance for computational induction, since they allow for an incremental approach. Incrementality 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. Under the rule of Consistency this set is consistent, which yields the following principle:

Left Consistency

$$\frac{\alpha \prec \beta}{\neg \alpha \not\prec \beta}$$

LEMMA 4.3. *In the presence of Right Reflexivity and Admissible Converse Entailment, Left Consistency implies Consistency.*

Proof. Suppose $\alpha \not\prec \neg \alpha$. Now, either $\beta \prec \beta$ or $\beta \not\prec \beta$; in the former case, $\neg \alpha \prec \beta$ by Admissible Converse Entailment, and we conclude by Left Consistency. In the latter case, we have $\delta \not\prec \beta$ for any δ by Right Reflexivity.

It follows that Left Consistency and Consistency are equivalent in the presence of Right Reflexivity, Admissible Converse Entailment, and Incrementality.

Convergence states a *monotonicity* property of induction, which can again best be understood by considering its contrapositive: a hypothesis that is rejected on the basis of evidence β cannot become feasible again when stronger evidence α is available. In other words: the process of rejecting a hypothesis is not defeasible (i.e. based on assumptions), but based on the evidence only. This is the analogue of the monotonicity property of deduction (note that the latter can be obtained by reversing the implication in the first condition of Convergence).

LEMMA 4.4. *The combination of Verification and Predictive Convergence is equivalent with the following rule:*

Predictive Convergence

$$\frac{\alpha \wedge \gamma \rightarrow \beta, \alpha \prec \gamma}{\alpha \wedge \beta \rightarrow \gamma}$$

Proof. To derive Predictive Convergence, suppose $\alpha \wedge \gamma \rightarrow \beta$ and $\alpha \prec \gamma$, then by Verification $\alpha \wedge \beta \prec \gamma$, and by Convergence $\beta \prec \gamma$.

Predictive Convergence implies Convergence, since $\alpha \rightarrow \beta$ implies $\alpha \wedge \gamma \rightarrow \beta$.

Predictive Convergence implies Verification, since $\alpha \wedge \beta \rightarrow \gamma$ implies $\alpha \wedge \beta \rightarrow \alpha \wedge \gamma$.

Predictive Convergence can be seen as a strengthening of Convergence, in the sense that β is not merely a weakening of evidence α , but can be any set of *predictions*. Note that Right Reflexivity is an instance of Predictive Convergence (put $\gamma = \beta$).

The final postulate we consider in this section 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 included implicitly, β is added to the hypothesis instead.

terminology employed here better reflects the meaning of the rules.

Conditionalisation

$$\frac{\alpha \not\prec \beta \wedge \gamma}{\beta \rightarrow \alpha \not\prec \gamma}$$

After having discussed various abstract properties of formation of explanatory hypotheses we now turn to the question of characterising explanatory induction semantically.

4.2 *Strong explanatory consequence relations*

As we have seen, Peirce’s original idea was to define explanatory hypothesis formation as reversed deduction. I will amend Peirce’s proposal in two ways. First, as explained above it is required that the hypothesis be consistent with respect to the background knowledge. Secondly, I reformulate reversed deduction as inclusion of deductive consequences. The main reason for the latter is that in this way the explanatory consequence relation is defined in terms of a property that is *preserved* by arguments (*viz.* explanatory power).

DEFINITION 4.5. An *explanation mechanism* is some consequence relation \vdash . The *explanatory consequence relation* $\not\prec$ defined by \vdash is defined as $\alpha \not\prec \beta$ iff $C_{\vdash}(\alpha) \subseteq C_{\vdash}(\beta) \subseteq L$. A *strong explanatory consequence relation* is defined by a monotonic explanation mechanism.

Thus, an explanation is required to have at least the same consequences under the explanation as the premiss it is obtained from. It should be noted that, in the general case, the conditions $C_{\vdash}(\alpha) \subseteq C_{\vdash}(\beta)$ and $\beta \vdash \alpha$ are not equivalent.¹⁵ However, for monotonic explanation mechanisms they are, which provides us with the following ‘Peircean’ definition of strong explanatory consequence relations.

DEFINITION 4.6. A *strong explanatory structure* is a set $W \subseteq M$. The consequence relation it defines is denoted by $\not\prec_W$ and is defined by: $\alpha \not\prec_W \beta$ iff (i) there is a $m_0 \in W$ such that $m_0 \vdash \beta$, and (ii) for every $m \in W$, $m \vdash \beta \rightarrow \alpha$.

The following system of rules will be proved to axiomatise strong explanatory structures.

DEFINITION 4.7. The system **EM** consists of the following rules: Admissible Right Strengthening, Explanatory Reflexivity, Incrementality, Predictive Convergence, Left Consistency, and Conditionalisation.

We note the following derived rules of **EM**: Convergence, Admissible Converse Entailment and Reflexivity (instances of Predictive Convergence) and Consistency (Lemma 4.3). The following derived rule will also prove useful.

LEMMA 4.8. *The following rule is a derived rule of EM:*

Consistent Right Strengthening
$$\frac{\alpha \not\prec \gamma, \neg \beta \not\prec \gamma}{\alpha \not\prec \beta \wedge \gamma}$$

Proof. Suppose $\neg \beta \not\prec \gamma$; since $\neg(\beta \wedge \gamma) \wedge \gamma \rightarrow \neg \beta$, we have $\neg(\beta \wedge \gamma) \not\prec \gamma$ by Predictive Convergence. Furthermore, suppose $\alpha \not\prec \gamma$, then by Right Reflexivity $\gamma \not\prec \gamma$, so by Explanatory Reflexivity we have $\beta \wedge \gamma \not\prec \beta \wedge \gamma$. We conclude by Admissible Right Strengthening.

Soundness of **EM** is easily checked.

LEMMA 4.9 (Soundness of **EM**). *Any strong explanatory consequence relation satisfies the rules of EM.*

¹⁵ $C_{\vdash}(\alpha) \subseteq C_{\vdash}(\beta)$ implies $\beta \vdash \alpha$ if \vdash is reflexive; $\beta \vdash \alpha$ implies $C_{\vdash}(\alpha) \subseteq C_{\vdash}(\beta)$ if \vdash is transitive.

Proof. Let $w \in W$ be a strong explanatory structure as defined in Definition 4.6, satisfying the rules of **EM**. Admissible Right Strengthening: if $\alpha \prec \beta$ then some model in W satisfies α . Furthermore, if $m \models \beta \rightarrow \alpha$ and $\alpha \prec \beta$, then $m \models \alpha$. Explanatory Reflexivity: we have that some models in W satisfy α , while not all models in W satisfy $\alpha \rightarrow \neg \beta$, i.e. there is a model in W satisfying $\alpha \wedge \beta$ and hence β . Incrementality: if $m \models \gamma \rightarrow \alpha$ and $m \models \gamma \rightarrow \beta$ then $m \models \gamma \rightarrow (\alpha \wedge \beta)$. Predictive Convergence: if $m \models \gamma \rightarrow \beta$ and $m \models \gamma \rightarrow \alpha$ then $m \models \gamma \rightarrow \beta$. Left Consistency: if some model in W satisfies β while all models in W satisfy $\beta \rightarrow \alpha$, then there is a model in W not satisfying $\beta \rightarrow \alpha$. Conditionalisation: trivial.

In order to prove completeness we build a strong explanatory structure W for a given consequence relation \prec satisfying the rules of **EM**, such that $\alpha \prec \beta$ iff $\alpha \prec_W \beta$. For non-empty explanatory relations the following construction is used:

$$W = \{m \in U \mid \text{for all } \alpha, \beta \text{ such that } \alpha \prec \beta \text{ we have } \beta \rightarrow \alpha\}$$

An empty explanatory relation signals inconsistent background knowledge, and is hence defined as empty explanatory structure.

We need a few intermediate results. The following lemma states that every strong explanatory hypothesis is satisfiable in W .

LEMMA 4.10. *Let \prec be a consequence relation satisfying the rules of **EM**, and let W be defined as above. If $\alpha \prec \beta$ then there is a model $m \in W$ such that $m \models \beta$.*

Proof. Let $\alpha \prec \beta$; we will prove that $\{\beta\} \cup \{\delta \rightarrow \gamma \mid \gamma \prec \delta\}$ is satisfiable. Suppose not, then by compactness there is a finite $\Delta \subseteq \{\delta \rightarrow \gamma \mid \gamma \prec \delta\}$ such that $\Delta \rightarrow \neg \beta$. Furthermore, since $\phi \prec \psi$ for any $\psi \rightarrow \phi \in \Delta$, we have $\psi \rightarrow \phi \prec \psi$ for any $\psi \rightarrow \phi \in \Delta$ by Conditionalisation, $\Delta \prec \mathbf{true}$ by Incrementality, and $\Delta \prec \beta$ by Right Reflexivity and Admissible Right Strengthening. But then by Consistency $\beta \rightarrow \neg \Delta$, a contradiction.

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

LEMMA 4.11. *Let \prec be a non-empty consequence relation satisfying the rules of **EM**, and let W be defined as above. If $\gamma \not\prec \gamma$ then γ is unsatisfiable in W .*

Proof. Let $\alpha \prec \beta$, then $\mathbf{true} \prec \mathbf{true}$ by Convergence and Left Reflexivity. Furthermore, if $\gamma \not\prec \gamma$ then $\neg \gamma \prec \mathbf{true}$ by Explanatory Reflexivity, hence $m \models \mathbf{true} \rightarrow \neg \gamma$ for every $m \in W$.

I will now show that W defines a consequence relation that is included in \prec .

LEMMA 4.12. *Let \prec be a consequence relation satisfying the rules of **EM**, and let W be defined as above. If $\alpha \prec_W \beta$ then $\alpha \prec \beta$.*

Proof. Suppose that $\alpha \prec_W \beta$, we will show that either no model in W satisfies β , or there exists a model in W satisfying α .

First of all, if $\beta \not\prec \beta$ then β is unsatisfiable in W by Lemma 4.11. In the remainder of the proof we will assume that $\beta \prec \beta$. We define $\Gamma_0 = \{\neg \alpha\} \cup \{\delta \mid \delta \prec \beta\}$; we will first show that Γ_0 is satisfiable. Suppose not, then by compactness there is a finite $\Delta \subseteq \{\delta \mid \delta \prec \beta\}$ such that $\Delta \rightarrow \neg \alpha$, i.e. $\Delta \rightarrow \alpha$, by Admissible Right Strengthening $\Delta \rightarrow \alpha \prec \beta$. Recall that $\beta \prec \alpha$. But by Incrementality $\Delta \prec \beta$; using incrementality and Convergence, we have $\Delta \prec \alpha$.

Let $m_0 \models \Gamma_0$; clearly $m_0 \models \alpha$ and since $\beta \in \Gamma_0$, $m_0 \models \beta$. It remains to prove that m_0 is in W ; i.e., that for all ψ, ϕ such that $\psi \rightarrow \phi$ we have $m_0 \models \psi \rightarrow \phi$. Let $\phi \not\prec \psi$; if $\neg \beta \not\prec \psi$, then by Consistent Right Strengthening $\phi \not\prec \psi \wedge \beta$, and by Conditionalisation $\psi \rightarrow \phi \not\prec \beta$; thus $\psi \rightarrow \phi \in \Gamma_0$ and therefore $m_0 \models \psi \rightarrow \phi$. On the other hand, if $\neg \beta \prec \psi$ then by Conditionalisation and Convergence $\beta \rightarrow \neg \psi \prec \mathbf{true}$, by Admissible Right Strengthening $\beta \rightarrow \neg \psi \prec \beta$, and by Incrementality and Convergence $\neg \psi \prec \beta$; thus $\neg \psi \in \Gamma_0$ and therefore $m_0 \models \neg \psi$, hence $m_0 \models \psi \rightarrow \phi$.

Armed with the previous three lemmas we can prove the completeness of **EM**.

THEOREM 4.13 (Representation theorem for strong explanatory consequence relations). *A consequence relation is strong explanatory iff it satisfies the rules of EM.*

Proof. The only-if part is Lemma 4.9. For the if part, let \prec be an arbitrary non-empty consequence relation satisfying the rules of **EM**, and let

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

Suppose $\alpha \prec \beta$, then by the construction of W , $m \models \beta \rightarrow \alpha$ for all $m \in W$. Furthermore, by Lemma 4.10 there is a model in W satisfying β . We may conclude that $\alpha \prec_W \beta$. Conversely, if $\alpha \prec_W \beta$ then Lemma 4.12 proves that $\alpha \prec \beta$. We conclude that W defines a consequence relation that is exactly \prec .

For an empty consequence relation put $W = \emptyset$.

In this section we have studied axioms and semantic characterisations for explanatory induction. I have proposed a novel definition of explanatory hypothesis formation in terms of preservation of explanatory power with respect to an explanation mechanism. A representation theorem has been obtained for the special case of a monotonic explanation mechanism. Characterisation of explanatory induction with respect to other (e.g. preferential) explanation mechanisms is left as an open problem.

5. Confirmatory induction

We will now switch from the explanatory (classification-oriented) viewpoint to the confirmatory (non-classificatory) perspective. Throughout this section $\alpha \prec \beta$ is to be read as ‘evidence α confirms hypothesis β ’. Our goals will be to find reasonable properties of \prec under this interpretation (for which we have a good starting point in Hempel’s adequacy conditions), and to characterise particular sets of properties by a suitable semantics.

5.1 Properties of confirmatory induction

In this section I will translate Hempel’s set of adequacy conditions (section 2.2) into rules for confirmatory consequence relations. The conditions will be slightly modified, in order to keep the treatment of inconsistent evidence and hypothesis in line with the explanatory case: inconsistent evidence does not confirm any hypothesis, and inconsistent hypothesis is confirmed by any evidence.

Entailment condition (H1) is translated into two rules:

Admissible Entailment

$$\frac{\alpha \rightarrow \beta, \alpha \prec \alpha}{\alpha \prec \beta}$$

Confirmatory Reflexivity

$$\frac{\alpha \prec \alpha, \alpha \not\prec \neg\beta}{\beta \prec \beta}$$

Admissible Entailment expresses that admissible evidence (i.e. evidence that is consistent with the background knowledge) confirms any of its consequences. In other words, consistent entailment is a special case of confirmation. Confirmatory Reflexivity is the confirmatory counterpart of Explanatory Reflexivity encountered in the previous section. It is added as a separate rule since, in its original formulation, (H1) includes reflexivity as a special case. As with its explanatory counterpart, Confirmatory Reflexivity is best understood when considering its contrapositive: if α is inadmissible, i.e. too strong a statement with regard to the background knowledge, its negation $\neg\beta$ is so weak that it is confirmed by arbitrary admissible formulae α .

Consequence condition (H2) cannot be translated directly, since in the explanatory case consequence relations as defined here we have no means to refer to a set of confirmed sentences. However, a translation of the special consequence condition (H2.1) and the conjunction condition (H2.3) will suffice:

Right Weakening

$$\frac{\alpha \rightarrow \gamma, \alpha \prec \beta}{\alpha \prec \gamma}$$

Right And

$$\frac{\alpha \vDash \beta, \alpha \vDash \gamma}{\alpha \vDash \beta \wedge \gamma}$$

LEMMA 5.1. *If \vDash is a consequence relation satisfying Right And and Right Weakening, then $\alpha \vDash \beta \wedge \gamma$ iff $\alpha \vDash \beta$ and $\alpha \vDash \gamma$.*

Proof. The *if* part is Right And, and the *only-if* part follows from Right Weakening.

Right Weakening expresses that any hypothesis entailed by a given hypothesis confirmed by α is confirmed by α . Notice that Admissible Entailment is an instance of Right Weakening.

LEMMA 5.2. *The combination of Right Extension and Right Weakening implies the following rule:*

Predictive Right Weakening

$$\frac{\alpha \vDash \beta \rightarrow \gamma, \alpha \vDash \beta}{\alpha \vDash \gamma}$$

Proof. In order to derive Predictive Right Weakening, suppose $\alpha \vDash \beta \rightarrow \gamma$ and $\alpha \vDash \beta$, then by Right Extension $\alpha \vDash \beta \wedge \gamma$, and the result follows by Right Weakening.

Predictive Right Weakening implies Right Weakening, since $\beta \rightarrow \gamma$ implies $\beta \rightarrow \beta$.

Predictive Right Weakening implies Right Extension, since $\beta \wedge \beta \rightarrow \gamma$ implies $\beta \wedge \beta \rightarrow \beta \wedge \gamma$.

In words, Predictive Right Weakening expresses that given a confirmatory argument, any predicted formula is confirmed by the same evidence. Notice that by putting $\gamma = \alpha$ in Predictive Right Weakening we obtain Left Reflexivity.

Right And states that the set of all confirmed hypotheses (interpreted as a conjunction) is itself confirmed. The combination of Right And and Right Weakening implies Hempel's general consequence condition (H2): if E confirms every formula of a set K , then it also confirms the conjunction of the formulae in K (by Right And), and therefore also every consequence of this conjunction (by Right Weakening)¹⁶. It has already been remarked that Right And is probably too strong in the general case, if we have inconclusive evidence that is unable to choose between incompatible hypotheses. In this respect it is perhaps appropriate to point at a certain similarity between Right And and Right Extension: the latter rule requires γ to be predicted rather than being confirmed by α .

Like the general consequence condition (H2), general consistency condition (H3) cannot be translated directly into a rule, since we have no means to refer to the set of confirmed formulae. However, in the light of Right And the conjunction of the formulae in this set is itself confirmed, and therefore it is sufficient to formulate a rule expressing the special consistency condition (H3.2) of Consistency previously encountered:

Consistency

$$\frac{\alpha \vDash \beta}{\beta \rightarrow \neg \alpha}$$

Condition (H3.2) expresses that for any formula β , if β is in the set of confirmed hypotheses then β does not entail $\neg \beta$. This principle is expressed by the following rule:

Right Consistency

$$\frac{\alpha \vDash \beta}{\alpha \vDash \neg \beta}$$

LEMMA 5.3. *In the presence of Admissible Entailment and Left Reflexivity, Right Consistency implies Consistency.*

Proof. Suppose $\alpha \vDash \beta \rightarrow \neg \alpha$, i.e. $\alpha \vDash \beta \rightarrow \neg \beta$. Now, either we have $\alpha \vDash \alpha$, or else $\alpha \not\vDash \alpha$. In the former case, $\alpha \vDash \beta$ by Admissible Entailment, and we conclude by Right Consistency. In the latter case, we have $\alpha \vDash \delta$ for any δ by Left Reflexivity.

Clearly, Consistency implies Right Consistency in the presence of Right And. As a corollary to Lemma 5.3, we have that Right Consistency and Consistency are equivalent in the presence of Left Reflexivity,

¹⁶This holds only for finite K , an assumption that I will make throughout.

Admissible Entailment, and Right And.

Finally, the equivalence condition for observations (H4) is translated into

Left Logical Equivalence

$$\frac{\alpha \leftrightarrow \beta, \alpha \prec \gamma}{\beta \prec \gamma}$$

We now turn to the question of devising a meaningful semantics for Hempel's conditions as re-expressed in our framework of inductive consequence relations.

5.2 Simple confirmatory structures

It has been suggested in section 2.5 that a semantics for confirmatory reasoning be expressed in terms of satisfaction by an appropriately constructed model or set of models. More precisely, a confirmatory semantics is conceived as one in which certain regular models are constructed from the premisses, such that a hypothesis is confirmed if it is true in all such regular models. Since this requires some completeness assumptions regarding the evidence, we may call this a *closed* confirmatory semantics. In section 5.4 I will consider a variant which relaxes the assumptions regarding the evidence (see *open* confirmatory semantics).

DEFINITION 5.4. A *simple confirmatory structure* is a triple $W = \langle S, [\cdot], \cdot \rangle$ where S is a set of semantic objects, and $[\cdot]$ and \cdot are functions mapping formulas to sets of semantic objects. The associated confirmatory consequence relation defined by W is given by:
 $\alpha \prec_W \beta$ iff (i) $[\alpha] \subseteq \alpha$ and (ii) $\alpha \subseteq \beta$.

Intuitively, $[\alpha]$ denotes the set of formulas generated from premisses α , each of which should satisfy hypothesis β . This is the notion of a *pragmatic model* considered by Bell [2], who calls it a *pragmatic model*. There are however some differences between each and mine. First, in order to rule out inconsistent premisses, condition (i) in the definition of \prec_W . Furthermore, it allows the possibility that some formulas may not satisfy the premisses (i.e. $[\alpha] \not\subseteq \alpha$). This possibility is ruled out by condition 6. However, the characterisation of \prec in terms of simple confirmatory structures is a problem, and the results obtained below differ from those reported in [1]. The additional condition (i) in the present section and the next.

Bell proves that the system **CS** of pragmatic models, including the condition $[\alpha] \subseteq \alpha$ are axiomatised by the rules of Reflexivity, Right Weakening, and Right And, and one additionally assumption is needed with respect to the logical connectives and classical entailment:

$$\frac{\alpha \prec \beta, \beta \prec \gamma}{\alpha \prec \gamma}$$

Let us call confirmatory structures which satisfy these conditions, as well as $[\alpha] \subseteq \alpha$ or $\alpha \subseteq [\alpha]$ *simple confirmatory structures*. The (closed) confirmatory consequence relation \prec_W will now demonstrate that simple confirmatory structures are axiomatised by the following rule system:

DEFINITION 5.5. The system **CS** consists of the following rules: Predictive Right Weakening, Right And, and Right Consistency.

Derived rules of **CS** include Right Weakening and Right Extension (Lemma 5.2), Left Reflexivity (an instance of Predictive Right Weakening), Admissible Entailment (an instance of Right Weakening), and Consistency (Lemma 5.3).

The soundness of the rules of **CS** is easily proved.

LEMMA 5.6 (Soundness of **CS**). Any simple closed confirmatory consequence relation satisfies the rules of **CS**.

Proof. For Predictive Right Weakening, first of all we have $[\alpha] \subseteq \alpha$ or $\alpha \subseteq [\alpha]$ hence

$[\alpha] \cap [\beta] \subseteq [\alpha \wedge \beta]$
 $[\alpha] \cap [\beta] \subseteq [\alpha \vee \beta]$
 and we have $[\alpha] \subseteq [\alpha]$
 For Right And $[\alpha \wedge \beta] \subseteq [\alpha]$
 For Right Or $[\alpha \vee \beta] \subseteq [\alpha]$

In order to prove completeness we need to build a confirmatory structure W that defines exactly the consequence relation \vdash .

DEFINITION 5.7. Let \vdash be a confirmatory consequence relation. The model $m \in U$ is said to be *normal for α* iff for all β in L such that $\alpha \vdash \beta$, $m \models \beta$.

For admissible formulae normal models will play the role of the regular models in the confirmatory structure we are building (remember that, given a consequence relation \vdash , a formula α is admissible iff $\alpha \vdash \alpha$).

LEMMA 5.8. *If a consequence relation \vdash satisfies Right Weakening and Right And, and let α be an admissible formula. All normal models for α satisfy β iff $\alpha \vdash \beta$.*

Proof. The *if* part follows from Definition 5.7. For the *only if* part, suppose $\alpha \vdash \alpha$ and $\alpha \not\vdash \beta$; I will show that there is a normal model for α that does not satisfy β . Let $\Gamma_0 = \{\neg\beta\} \cup \{\delta \mid \alpha \vdash \delta\}$; it suffices to show that Γ_0 is satisfiable. Since $\alpha \vdash \alpha$, then by compactness there is a finite $\Delta \subseteq \{\delta \mid \alpha \vdash \delta\}$ such that $\alpha \vdash \beta$, i.e. $\Delta \vdash \beta$ by Right Weakening $\alpha \vdash \Delta \rightarrow \beta$. But by Right And $\alpha \vdash \Delta$; using Right And and Right Weakening we obtain $\alpha \vdash \beta$, a contradiction.

In the standard treatment of consistent premisses they have all formulae as consequences, hence no normal models – since in our treatment inconsistent premisses confirm no hypothesis and thus have all models in U as normal models. We have to treat them as a separate case. Given a consequence relation \vdash , let $W = \langle U, [\cdot], \cdot \rangle$ be defined as follows:

- (1) U is the set of models of L under consideration;
- (2) $[\alpha] = \{m \in U \mid m \text{ is a normal model for } \alpha\}$ if α is admissible, and \emptyset otherwise;
- (3) $\alpha \vdash_W \beta$ iff $\alpha \vdash \beta$.

The following compactness result demonstrates that the consequence relation defined by W coincides with the original one. The latter satisfies the rules of \mathcal{S} .

THEOREM 5.9. *Let \vdash be a simple confirmatory consequence relation. A consequence relation is simple confirmatory iff it satisfies the rules of \mathcal{S} .*

Proof. The only-if part is Lemma 5.4. For the if part, let \vdash be a consequence relation satisfying the rules of \mathcal{S} and let \vdash_W be defined as above. We will prove that $\alpha \vdash \beta$ iff $\alpha \vdash_W \beta$, i.e. \vdash is simple confirmatory.

First suppose that $\alpha \vdash \beta$, then by Left Reflexivity $\alpha \vdash \alpha$ and by Right Negativity $\alpha \vdash \neg\neg\alpha$, so the proof of Lemma 5.8 constructs a model m normal for α that does not satisfy β . Furthermore any normal model for α satisfies β by Definition 5.7, hence $[\alpha] \subseteq [\beta]$. We conclude that $\alpha \vdash_W \beta$. Now suppose that $\alpha \vdash_W \beta$, then $[\alpha] \subseteq [\beta]$ since α is admissible by the definition of W . Furthermore we have $[\alpha] \subseteq [\beta]$, i.e. every model normal for α satisfies β and the conclusion follows by Lemma 5.8.

One may note that two of the rules obtained in the previous section are not needed above, viz. Confirmatory Reflexivity and Left Logical Equivalence. The inclusion of these rules poses additional restrictions on simple confirmatory structures:

- (Left Logical Equivalence) If $\alpha \vdash \beta$, then $[\alpha] = [\beta]$.
- (Confirmatory Reflexivity) If $\beta \in \emptyset$ then $[\alpha] \neq \emptyset$.

Additional properties may be obtained by being more explicit about the construction of regular models. The ILP methods referred to in section 2.5 suggest to take the truth minimal Herbrand model(s) of the evidence as the regular model(s). In the analysis of nonmonotonic logics it is customary to abstract this into a preference ordering on the set of models, such that the regular models are the minimal ones

under this ordering. We will write $\text{th}(s)$ for the *th* function.

5.3 Preferential confirmatory consequence relation

The main result of this section is the axiomatisation of Kraus *et al.*'s preferential semantics, as a confirmatory semantics. The objects of the semantics are those that are minimal with respect to a fixed preference ordering.

DEFINITION 5.10. A *preferential structure* is a triple $W = \langle S, I, < \rangle$ where S is a set of *states*, $I: S \rightarrow U$ is a function that associates with a model, and $<$ is a strict partial order¹⁷ on S , called the *preference ordering* at W .¹⁸ W defines a *preferential confirmatory structure* $\langle W, \cdot \rangle$ and a *preferential confirmatory consequence relation* as follows: $\alpha \vDash \beta$ iff $\{s \in S \mid I(s) \in \alpha\} \cap \text{th}(\alpha) \subseteq \text{th}(\beta)$.

Note that preferential confirmatory structures are simple confirmatory structures, and that they satisfy the conditions assumed in [18]. Logical Equivalence and Confirmatory Reflexivity are satisfied, and the smoothness condition, if $s \in \alpha$ then there is $t \in \alpha$ such that $t < s$ — hence $\alpha \neq \emptyset$ implies $[\alpha] \neq \emptyset$.

In comparison with the preferential semantics of [18], the only difference is that in a preferential confirmatory argument the evidence is required to be satisfiable, in order to guarantee the validity of Consistency. The intermediate semantics *states* is mainly needed for technical reasons, and can be interpreted as the set of models the reasoning agent considers possible in that epistemic state.

The following set of rules will be proved to axiomatise preferential confirmatory consequence relations.

DEFINITION 5.11. The system **CP** consists of the rules of **CS**, Confirmatory Reflexivity, Left Logical Equivalence, plus the following rules:

Left Or

$$\frac{\alpha \vDash \gamma, \beta \vDash \gamma}{\alpha \vee \beta \vDash \gamma}$$

Strong Verification

$$\frac{\alpha \vDash \gamma, \alpha \vDash \beta}{\alpha \wedge \gamma \vDash \beta}$$

The first rule is a variant of Confirmatory Reflexivity. The general case: if we weaken the evidence with a disjunction, then the evidence confirms the hypothesis. However, for strong verification, we only require that the hypothesis is confirmed by taking their disjunction. This is a restriction of the general case, which is satisfied if the formula γ is a disjunction of formulas β_i such that $\alpha \vDash \beta_i$ for each i . This is also allowed when the confirmatory structure is a preferential structure. This is very similar to the way in which the strong verification rule is satisfied in the evidence strength logic, where the hypothesis is a disjunction of formulas.

LEMMA 5.11. Soundness of **CP** for preferential confirmatory structures. *CP* satisfies the rules of **CP**.

Proof. The proof only needs to be carried out for the two new rules. For Left Or, note that if $\alpha \vDash \gamma$ and $\beta \vDash \gamma$, then $\alpha \vee \beta \vDash \gamma$ is satisfied, since $\text{th}(\alpha \vee \beta) = \text{th}(\alpha) \cup \text{th}(\beta)$. For Strong Verification, note that if $\alpha \vDash \gamma$ and $\alpha \vDash \beta$, then $\alpha \wedge \gamma \vDash \beta$ is satisfied, since $\text{th}(\alpha \wedge \gamma) = \text{th}(\alpha) \cap \text{th}(\gamma)$.

¹⁷I.e., $<$ is irreflexive and transitive.
¹⁸I.e. for any $S' \subseteq S$ and for any $s \in S'$, either s is minimal in S' , or there is $s' \in S'$ such that $s' < s$. This condition is satisfied if $<$ does not allow infinite descending chains.

¹⁹In the context of nonmonotonic reasoning, Strong Verification is known as Cautious Monotonicity. For the purposes of this paper I prefer to use the first name, which expresses more clearly the underlying intuition in the present context.

non-empty, hence $[\alpha \wedge \gamma] \neq \emptyset$. Now $\langle n, \gamma \rangle \in [\alpha]$ when $s \in \alpha$ and $\langle n, \gamma \rangle \in [\alpha]$.
 Suppose not, then there is a $t \in \alpha$ such that $t \prec \langle n, \gamma \rangle$ and $t \in [\alpha] \subseteq [\alpha \wedge \gamma]$. But this
 contradicts the minimality of s in $[\alpha]$.

In order to prove completeness, we need to build a preferential structure W that verifies complete
 relation \prec satisfying the rules of **CP**, such that $\alpha \prec \beta$ iff $\alpha \prec \beta$. This structure is defined as follows:

- (1) $S = \{\langle m, \alpha \rangle \mid \alpha \text{ is an admissible formula, and } m \text{ is a normal model for } \alpha\}$;
- (2) $l(\langle m, \alpha \rangle) = m$;
- (3) $\langle m, \alpha \rangle \prec \langle n, \beta \rangle$ iff $\alpha \vee \beta \prec \alpha$ and $m \models \beta$.

Thus, states are pairs of admissible formulae and normal models. The label l simply maps a state to the model it contains. Condition (3) defines a total ordering between states: note that α is a special case of $\alpha \vee \beta \prec \alpha$ by means of Left Or, and in fact that α is admissible. The condition is added to make the ordering irreflexive; note that in a sequence any $\langle m, \alpha \rangle$ is minimal in $[\alpha]$.

The main difference between the preferential consequence relations of Kraus *et al.* and my preferential confirmatory consequence relations is the way unsatisfiable formulae are treated. In Kraus *et al.*'s state **P** unsatisfiable formulae are characterised by the fact that they have every formula in L as a plausible consequence, which means that they don't have any normal models. In my framework, unsatisfiable formulae confirm no hypotheses, and have all models in L as normal models. In both cases, the structure W that is used to prove completeness contains only satisfiable formulae in its states. This means that we replicate most of Kraus *et al.*'s results about the structure of preferential consequence relations.

- PROPOSITION 5.13. (1) [18, Lemma 5.3] *The relation \prec is a strict total order.*
 (2) [18, Lemma 5.15] *The relation \prec is such that for any $s \in \alpha$ with s minimal in $[\alpha]$ there exists a state $t \prec s$ minimal in $[\alpha]$.*
 (3) [18, Lemma 5.11] *If $\alpha \vee \beta \prec \alpha$ and m is a normal model for α and β , then m is a normal model for β .*
 (4) [18, Lemma 5.14] *$\langle m, \alpha \rangle \in [\beta]$ iff $m \models \beta$ and $\alpha \prec \beta$.*

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

LEMMA 5.14. *Let \prec be a consequence relation satisfying the rules of CP, and let W be defined as above. If $\alpha \prec \beta$ then $\alpha \prec_W \beta$.*
 Proof. Suppose that $\alpha \prec \beta$, then by Left Reflexivity $\alpha \prec \alpha$ and by Right Consistency $\alpha \prec \neg \alpha$. The proof of Lemma 5.8 constructs a model m normal for α hence $[\alpha] \neq \emptyset$. Furthermore suppose $s = \langle n, \gamma \rangle \in [\alpha]$, then γ is an admissible formula, and a normal model n for γ that satisfies α and $\gamma \vee \alpha \prec \gamma$. By Proposition 5.13 (1) n is a normal model for α and thus satisfies β by Definition 5.1. Hence $\langle n, \gamma \rangle \in [\beta]$. We conclude that $\emptyset \subset [\alpha] \subseteq [\beta]$ when $\alpha \prec \beta$.

The following lemma proves the converse of the previous lemma, completing the proof of the representation theorem.

LEMMA 5.15. *Let \prec be a consequence relation satisfying the rules of CP, and let W be defined as above. If $\alpha \prec_W \beta$ then $\alpha \prec \beta$.*
 Proof. Suppose $\alpha \prec_W \beta$, i.e. $\emptyset \subset [\alpha] \subseteq [\beta]$. We will first prove that α is admissible. Let $\langle n, \gamma \rangle \in [\alpha]$, then $\gamma \prec \gamma$ and n is a normal model for γ . If $\alpha \not\prec \alpha$, then by Confirmatory Reflexivity $\gamma \prec \neg \alpha$ and thus $n \models \neg \alpha$, contradicting the assumption that $\langle n, \gamma \rangle \in [\alpha]$ – s is admissible. Furthermore, given any model m normal for α , $\langle m, \alpha \rangle \in [\alpha] \subseteq [\beta]$ then m satisfies β , and the conclusion follows by Lemma 5.13 (4).

We may now summarise.

THEOREM 5.16 (Representation theorem for preferential confirmatory consequence relations). *A consequence relation is preferential confirmatory iff it satisfies the rules of CP.*

Proof. The only-if part is LEMMA 5.12. For the if part, let \prec be a consequence relation

satisfying the rules of **CP** and let W be defined as above. Lemma 5.14 and Lemma 5.15 prove that $\alpha \prec \beta$ iff $\alpha \prec_W \beta$, i.e. \prec is preferential confirmatory.

We end this section by noting that **CP** represents a more restrictive form of reasoning than Kraus *et al.*'s system **P**, even though neither system entails the other (Reflexivity, a rule of **P**, is invalid in **CP**, while Consistency, a rule of **CP**, is invalid in **P**). This is so because from every preferential consequence relation \prec we can construct a preferential confirmatory consequence relation by removing all arguments $\{\alpha \prec \beta \mid \alpha \prec \delta \text{ for all } \delta \in L\}$, i.e. all arguments with a left-hand side that is inconsistent with respect to the background knowledge. The resulting preferential confirmatory relation is the largest one contained in \prec ; furthermore, all preferential confirmatory relations can be constructed in this way from a preferential relation. Thus, **CP** is more restrictive than **P**, as defined at the end of section 3.3.

5.4 Weak confirmatory consequence relations

Any semantics that is to obey rules like Right Anticipation and Strong Verification must be based on complete assumptions with regard to the evidence. On the other hand, such strong assumptions cannot be made in induction tasks. It seems reasonable, then, to investigate also an alternative approach in which a hypothesized hypothesis is required to be true in *some* of the regular models.²⁰ It is not difficult to see that alternative semantics based on some notion of completeness would invalidate both Right Anticipation and Strong Verification. On the other hand, by not making completeness assumptions one may agree with the property of Convergence.

DEFINITION 5.18. Let $W = \langle S, [\cdot], \cdot \rangle$ be a confirmatory structure. The open confirmatory consequence relation defined by W is given by: $\alpha \prec_W \beta$ iff $[\alpha] \cap \beta \neq \emptyset$.

We will henceforth use the term *open confirmatory relations*, which holds when $[\cdot]$ is identified with \cdot , which is \cdot -behaved with respect to the connectives and containment.

DEFINITION 5.19. A *classical confirmatory structure* is a simple confirmatory structure $\langle S, [\cdot], \cdot \rangle$. A consequence relation is called *weak confirmatory* iff it is the open consequence relation of a classical confirmatory structure.

From this definition it is clear that weak confirmatory consequence relations satisfy both Right Weakening and Right Strengthening (i.e. Convergence), as well as Consistency. One additional rule is needed.

DEFINITION 5.19. The system **CW** consists of the following rules: Predictive Convergence, Predictive Right Weakening, Consistency, and

Disjunctive Rationality

$$\frac{\alpha \vee \beta \prec \gamma, \beta \prec \gamma}{\alpha \prec \gamma}$$

Disjunctive Rationality has not been considered before. The rule can be borrowed from Kraus *et al.*, who identify it as a valid principle of plausible reasoning. In the context of confirmatory reasoning, Disjunctive Rationality is a rather strong rule, which states that a hypothesis is confirmed by disjunctive observations if it is confirmed by at least one of the disjuncts.

The following theorem proves the equivalence of weak confirmatory structures and the system **CW**.

THEOREM 5.20 (Representation theorem for weak confirmatory consequence relations). A consequence relation \prec is weak confirmatory iff it satisfies the rules of **CW**.

Proof. The only-if part involves demonstrating that a weak confirmatory structure satisfies the rules of **CW**, which is trivial.

For the if part, let \prec be an arbitrary consequence relation satisfying the rules of **CW**, and consider the weak confirmatory structure $W = \langle S, [\cdot], \cdot \rangle$.

²⁰This is sometimes called *credulous* inference, in contrast with *skeptical* inference, which requires truth in all regular models.

Some $m \in S$ for $\alpha, \beta \in L$ such that $m \models \alpha \wedge \beta$: $\alpha \prec_W \beta$.
 $\{m \mid m \models \alpha \wedge \beta\} \neq \emptyset$.
 We will prove $\alpha \prec_W \beta$ if $\alpha \prec_W \beta$. The if part follows directly from the construction of W . Suppose $\alpha \prec_W \beta$ and let $m_0 \in S$ be a model $m_0 \in S$.
 Define $\Gamma_0 = \{\alpha \mid \alpha \prec_W \beta\}$; we will first show $\Gamma_0 \vdash \beta$ by
 by compactness there is a finite $\Delta \subseteq \{\delta \mid \neg \delta \prec \beta\}$ such that $\neg(\Delta \wedge \alpha)$, i.e. $\neg(\Delta \wedge \alpha)$; by
 Consistency $\Delta \wedge \alpha \notin \Gamma_0$. Furthermore, since $\neg \delta \prec \beta$ for $\delta \in \Delta$, we have $\Delta \prec \beta$ by Disjunctive
 Rationality and $\neg \Delta \wedge \alpha \prec \beta$ by Convergence. Combining this with $\Delta \prec \beta$ we have
 $(\neg \Delta \wedge \alpha) \wedge \Delta \prec \beta$ by Disjunctive Rationality and $\alpha \prec \beta$ by incrementality. Contradiction,
 so Γ_0 is satisfiable.
 Let $m_0 \in \Gamma_0$; clearly $m_0 \models \alpha$ and, since by Consistency $\beta \in \Gamma_0$, $m_0 \models \beta$. It remains to prove
 that $m_0 \in S$; i.e., that for all ϕ, ψ such that $m_0 \models \phi \wedge \psi$ we have $\phi \prec \psi$. Let $m_0 \models \phi \wedge \psi$, then
 $\neg(\phi \wedge \psi) \notin \Gamma_0$, hence $\phi \wedge \psi \prec \beta$; by Predictive Right Weakening $\phi \wedge \psi \prec \psi \wedge \beta$, by
 Convergence $\phi \prec \psi \wedge \beta$, and by Right Weakening $\phi \prec \psi$.

The system **CW** thus provides an axiomatisation of the relation of logical compatibility.

In this section we have studied abstract properties and semantics for confirmatory induction. In the first part of this analysis I have demonstrated that Hempel's original conditions axiomatise a rather general form of confirmatory inference, characterised by truth in the regular models of the evidence. A close link with nonmonotonic or plausible reasoning is obtained by identifying the regular models with the minimal models under a preference ordering. Finally, I have proposed a more liberal form of confirmatory reasoning based on some notion of consistency, which is more appropriate if the evidence cannot be considered complete. This more liberal form of confirmatory reasoning invalidates the strong rule of Right And. A new representation theorem has been obtained for the extreme form of logical compatibility. Open problems include dropping the condition that regular models be models of the premisses, and more meaningful forms of open confirmatory reasoning.

6. Discussion

In this paper I have combined and extended old and recent work in philosophy, logic, and Machine Learning, in an attempt to gain more insight in induction as a reasoning process. I believe that the approach followed has implications for each of these three disciplines, which I will address separately below.

6.1 Philosophy

Induction is one of the traditional problems of philosophy. In spite of this it is perhaps also one of the least understood, and certainly one that hasn't been satisfyingly solved. From the perspective from which this paper is written there is not one, but two 'problems of induction': one concerned with the justification of accepted hypotheses, and one concerned with the formation of possible hypotheses. Traditionally philosophers have been concerned with the justification problem:

'Why is a single instance, in some cases, sufficient for a complete induction, while in others, myriads of concurring instances, without a single exception known or presumed, go such a very little way towards establishing an universal proposition? Whoever can answer this question knows more of the philosophy of logic than the wisest of the ancients, and has solved the problem of induction.' [21, Book III, Chapter VIII, p.314]

No attempt has been made, in the present paper, to solve *this* problem of induction. As I have argued in section 2.1, I don't think the justification problem manifests itself exclusively with inductive generalisations, but rather with all forms of nondeductive reasoning.

Furthermore, I don't think that the justification problem is a problem of logic. What I perceive as the *logical problem of induction*, and what has been the central problem of this paper, is the problem of finding a sufficiently accurate description, in logical terms, of the process of inductive hypothesis formation. As Hanson observes, this logical problem of induction has been mostly ignored:

‘Logicians of science have described how one might set out reasons in support of an hypothesis once it is proposed. They have said little about the conceptual considerations pertinent to the initial proposal of an hypothesis. There are two exceptions: Aristotle and Peirce. When they discussed what Peirce called “retroduction”²¹, both recognized that the proposal of an hypothesis is often a reasonable affair. (...)

Neither Aristotle nor Peirce imagined himself to be setting out a manual to help scientists make discoveries. There could be no such manual. Nor were they discussing the psychology of discoverers, or the sociology of discovery. There are such discussions, but they are not logical discussions. Aristotle and Peirce were doing logic. They examined characteristics of the reasoning behind the original suggestion of certain hypotheses.’ [12, pp.1073–4]

The first philosophical contribution of this paper, then, has been to reinforce the point that inductive hypothesis formation is a logical matter (and hypothesis selection is not). The second contribution lies in the distinction I have drawn between explanatory induction and confirmatory induction. More generally, I believe that a logical characterisation of inductive hypothesis formation is impossible unless one takes into account the *goal* which the hypothesis is intended to fulfil. Possible goals include providing classifications like in concept learning, or making implicit regularities explicit like in knowledge discovery in databases. Different goals lead to different forms of induction with different logical characteristics. For this reason I have taken a somewhat relativistic viewpoint, by setting up a general framework in which various logics of induction can be set up, analysed, and compared, instead of fixing a particular logic of induction.

Two main families of logics of induction have been singled out, one portraying induction as explanation-preserving reasoning, the other as inference of confirmed hypotheses. As these families should be taken as starting points for further technical research, let us consider a number of possible improvements here. First of all, it has been remarked earlier that the restriction to a propositional object language L is counterintuitive in the context of induction, where the distinction between statements about individuals and statements about sets of individuals appears to be crucial. Upgrading the results of this paper to predicate logic is certainly an important open problem. On the other hand, the propositional analysis of this paper is not meaningless for the predicate logic case. First of all, in finite domains statements of predicate logic can be encoded in propositional logic.²² Furthermore, I would expect each of the predicate logic rule systems to include the corresponding propositional rules, and the predicate logic semantics to be refinements of the propositional semantics.

A related point concerns the (restricted) reflexivity of the logics proposed here, which again runs counter to the intuition of induction as generalisation from instances to populations. We expect names of individuals to be present in the inductive premisses, and absent from the inductive hypothesis. As a consequence changing the names in the premisses (in a way that does not distort the information they convey) would not validate the hypothesis. For instance, consider the observations

Raven(a) \wedge Black(a) \wedge \neg Raven(c) \wedge Black(c) \wedge \neg Raven(d) \wedge \neg Black(d)

Given these observations $\forall x: \text{Raven}(x) \rightarrow \text{Black}(x)$ is a possible hypothesis. In the confirmatory setting this can be interpreted as constructing a regular model from the evidence, and stipulating that any hypothesis is satisfied by that model. However, from that model we can construct another model by replacing a, c, d with (say) e, f, g . It makes sense to say that this constitutes another regular model that is to satisfy any hypothesis — this would rule out any hypothesis that talks about a, c or d , including the observations in the premisses. In the explanatory setting of the confirmatory structures of section 5.2 this would mean to drop the condition that $[\alpha] \subseteq \alpha$ or $[\alpha] \in D$. In this way the perception of induction as reasoning from the particular to the general would be preserved, not by posing syntactic restrictions (which is logically unattractive), but by wiring it into the semantics. I am currently working on such a logic of generalisation.

6.2 Logic

Mathematical logic has been at the focus of attention of logicians at least since the beginning of the century. This has led to an underappreciation of other forms of reasoning. Work on commonsense

²¹Peirce’s translation of Aristotle’s term $\alpha\pi\alpha\gamma\omega\gamma\eta$ — only later Peirce introduced the term ‘abduction’.

²²This may require a substantial background theory, limiting the practical feasibility of this encoding.

reasoning in Artificial Intelligence has revived the interest in non-standard logics, but still most of the map of reasoning forms is in darkness. The regrettable use of the term ‘nonmonotonic reasoning’ for reasoning with rules that have exceptions²³ is a symptom of this twilight: virtually all nondeductive reasoning is nonmonotonic, yet reasoning with default rules constitutes but one possible mode of nondeductive reasoning. Logic should also be concerned with the systematic study of *reasoning forms* (deduction, induction, plausible reasoning, counterfactual reasoning, and so on).

From a logical point of view this paper can be seen as a contribution to the systematic study of reasoning forms, by trying to nail down the essence of induction in logical terms. By employing consequence relations as the central notion in this analysis, rather than, say, introducing non-truth-functional connectives or modalities, attention is focused on the underlying inference mechanism. Abstract properties of consequence relations can be used to classify reasoning forms: for instance, we could say that a reasoning form is deductive if it satisfies Monotonicity, quasi-deductive if it satisfies Cautious Monotonicity, explanatory if it satisfies (Admissible) Converse Entailment, and so on. These properties can be used to chart the whole map of reasoning forms, just like the properties of [18] chart the map of plausible reasoning. Thus, the present paper can be seen as a constructive proof of the thesis that the method of analysis through consequence relations, pioneered by Gabbay [11], Makinson [20], and Kraus, Lehmann & Magidor [18, 19], constitutes in fact a *methodology*, that can be applied to analyse arbitrary forms of reasoning.

It is important to note that, by pursuing the analysis on the level of consequence relations, one is studying a class of logics, i.e. a reasoning form, rather than a particular logic. The formal analysis of a reasoning form is quite different from the material definition of a particular logic. The traditional picture of the latter process is like this. One starts with the semantics, which is designed to provide a precise meaning for the primitive symbols in the language, and formalises the relevant notion of consequence. Only then the proof-theoretic axiomatisation follows, accompanied by proofs of soundness and completeness. However, the abstract analysis of a class of logics may well proceed along different lines. Usually, a number of material definitions of specific logics are available (e.g. default logic, circumscription, negation as failure), and one tries to understand what these logics have in common. This extraction of commonalities may start on the semantic level (e.g. each of the logics selects among the models of the premisses) but also on the meta-theoretical level (e.g. each of the logics is nonmonotonic, and closed under conjunction on the right-hand side). Semantics does not necessarily come first anymore: the notion of Cautious Monotonicity may add as much to the understanding of plausible reasoning as the idea of a preference ordering on models.

Another tendency that can be observed when moving towards more abstract characterisations of reasoning forms is that semantics concentrates more on a characterisation of the notion of consequence, and less on the meaning of the primitive symbols in the language. This raises the question as to what constitutes a logical semantics. This is by no means a settled issue, but different authors have put forward the notion of *preservation* as playing a central role in semantics. For instance, Jennings, Chan and Dowad ‘argue for a generalisation of inference from the standard account in terms of truth preservation to one which countenances preservation of other desirable metalinguistic properties’:

‘(...) **truth is not the only inferentially preservable property.** A system of inference essentially provides procedures by which a set Σ of sentences (for example, the set of one’s beliefs) having some complex of metalinguistic properties can be unfailingly extended to a larger set Σ' having the same complex of properties. By all means, we may regard *truth* as one of the properties to be preserved, but what other properties are to be preserved can depend upon our interests.’ [17, p.1047]

Furthermore, the view that abduction is the logic of preserving explanations has also been put forward by Zadrozny, who calls an inference rule abductive ‘if it preserves sets of explanations’ [25, p.1].

These points underline that many open problems of contemporary logic are conceptual in nature, rather than just technical. Motivated by problems from Artificial Intelligence, logic is widening its scope to include forms of reasoning that are less and less similar to classical deduction. This development is far from nearing its completion. For instance, it seems commonplace among researchers studying consequence relations to assume that any consequence relation should minimally satisfy Reflexivity, Right Weakening, and Right And. Surely this makes sense if one limits attention to quasi-deductive reasoning, but the explanatory consequence relations studied in section 4 of this paper satisfy none of these properties. This paper makes a case for a much more liberal perception of what is a consequence relation.

²³A better term would be ‘plausible reasoning’.

6.3 Machine Learning

Whereas the main contributions of this paper are logical and philosophical in nature, it also provides a novel perspective on inductive Machine Learning. First of all, the framework of inductive logic elaborated above provides a new logical foundation of inductive learning. As such, viewing inductive learning through logic is of course not new, as is witnessed by the subdiscipline of Inductive Logic Programming (ILP) [22]. However, the usual logical perception of inductive learning differs from the one presented here. These two perceptions can be explained by considering the phrase ‘inductive logic programming’. This phrase is ambiguous: it can be taken to mean — as it is usually done — ‘doing logic programming inductively’, but it can also be parsed as *programming in inductive logic* — which corresponds to the perception elaborated in this paper. Let me explain why these interpretations are fundamentally different.

By identifying ILP with doing logic programming inductively, one effectively says that one’s main goal is logic programming, i.e. answering queries by executing a declarative specification by means of its procedural semantics; however, since this declarative specification is only partly known through a number of examples, we should do some inductive patching before the real work can start. This results in a somewhat subsidiary view of induction as a subproblem that needs to be solved before we can do the main task. This can be seen from the problem specification (see section 2.5), which defines induction as a sort of reversed deduction from positive examples p to hypothesis H . The slogan *ILP = Inductive Logic + Programming* offers a different viewpoint, from which the inference from examples to a logic program is the main step. That is, the examples (and background theory) provide the declarative specification of the induction task, which is executed by applying the inference rules of an inductive logic (e.g. specialisation or generalisation operators). The hypothesis is an inductive consequence of the examples.

An immediate advantage of the latter viewpoint is that it provides an independent definition of induction, instead of defining it in terms of something else (e.g. reversed deduction). The following analogy may clarify this point. In mathematics, many concepts are introduced as inverses of other concepts: division as inverse of multiplication, roots as inverses of powers, integration as inverse of differentiation, and so on. However, once such a concept has been introduced in this way, it usually gets an independent treatment, providing further insight in and justification of the new concept. For instance, the definition of a definite integral as the limit of a Riemann sum formalises the idea that a definite integral calculates the area under a curve. The relationship between the new concept and previously defined concepts is then obtained as a theorem (the fundamental theorem of calculus), rather than a seemingly arbitrary definition.

The framework of inductive consequence relations can be used to obtain soundness and completeness proofs of particular sets of operators employed in a particular induction algorithm. As an illustration, consider the inverse resolution operator \leftarrow . If $p \leftarrow q, r, t$ infer $s \leftarrow p, t$. In our framework, a soundness proof of this induction step is established by the derivation, from an appropriate rule system such as **EM**, of the following statement:

$$\frac{p \leftarrow q, r}{s \leftarrow q, r, t \vdash s \leftarrow p, t}$$

which should be read as follows: if $p \leftarrow q, r$ is known from the background knowledge, then $s \leftarrow p, t$ is an inductive consequence of $s \leftarrow q, r, t$.

The framework also provides a well-defined *vocabulary* for reasoning about induction tasks and algorithms. Important notions like incrementality and convergence are linked to the underlying inductive consequence relation. This is important because, as we have seen, different induction tasks may have different characteristics — articulating these characteristics is a necessary and important first step in understanding the induction task at hand, and choosing or devising the right algorithm. Using such a vocabulary we can construct a taxonomy of induction tasks. Two families of induction tasks have been discerned in this paper: explanatory induction, aimed at obtaining classification rules, and confirmatory induction, directed towards extracting structural regularities from the data. While explanatory induction corresponds to typical classification-oriented inductive Machine Learning tasks such as concept learning from examples, induction tasks belonging to the non-classificatory paradigm have only recently started to attract attention [14, 6, 5]. The contribution of the present paper has been to propose Hempel’s notion of qualitative confirmation as the unifying concept underlying these latter approaches. As Helft observed, inductive conclusions may be obtained by closed-world reasoning if the evidence may be considered

complete; however, above I have proposed to distinguish an open variant of confirmatory reasoning, where the completeness assumptions on the evidence are considerably relaxed. While this may lead to an increase in the set of inductive hypotheses not refuted by given evidence, it has the distinct advantage of invalidating the rather strong property of Right And, and restoring the computationally attractive property of Convergence.

7. Conclusions

This paper has been written in an attempt to increase our understanding of inductive reasoning through logical analysis. What logic can achieve for arbitrary forms of reasoning is no more and no less than a precise definition of what counts as a *possible* conclusion given certain premisses. Selecting the best (most useful, most plausible, etc.) hypothesis is an extra-logical matter. The logic of induction is the logic of inductive hypothesis formation.

There is not a single logic of induction. The logical relationship between evidence and possible inductive hypotheses depends on the task these hypotheses are intended to perform. Induction of classification rules such as concept definitions are based on a notion of explanation; an alternative logical account of induction starts from the qualitative relation of confirmation. Other forms of induction are conceivable. In this paper I have proposed a metalevel framework for characterising and reasoning about different forms of induction. The framework does not fix a material definition of inductive hypothesis formation, but can be used to aggregate knowledge about classes of such material logics of induction.

A number of technical results have been obtained. The system **EM** axiomatises explanation-preserving reasoning with respect to a monotonic explanation mechanism. Characterisation of explanatory induction with respect to weaker (e.g. preferential) explanation mechanisms is left as an open problem. The systems **CS** and **CP** axiomatise the general and preferential forms of closed confirmatory reasoning, conceived as reasoning about selected regular models. They represent variations of earlier, differently motivated, characterisations by Bell [2] and Kraus, Lehmann & Magidor [18], with a different treatment of inconsistent premisses. An important open problem here is the axiomatisation of confirmatory structures where regular models may not be models of the premisses. Finally, the system **CW** represents an extreme form of open confirmatory reasoning (i.e. compatibility of premisses and hypothesis). Finding more realistic forms of open confirmatory reasoning remains an open problem.

Acknowledgements

I like to thank John-Jules Meyer and Daniel Lehmann for stimulating discussions. Part of this work was supported by Esprit IV Long Term Research Project 20237 (Inductive Logic Programming 2), and by PECO Pan-European Scientific Network ILPnet (no. CIPA3510OCT920044).

References

- [1] D. Angluin & C.H. Smith (1983), 'Inductive inference: theory and methods', *Computing Surveys* **15** (3): 238–269.
- [2] J. Bell (1991), 'Pragmatic logics'. In *Proc. Second International Conference on Knowledge Representation and Reasoning KR'91*, Morgan Kaufmann, San Mateo, pp. 50–60.
- [3] R. Carnap (1950), *Logical Foundations of Probability*, Routledge & Kegan Paul, London.
- [4] J.P. Delgrande (1987), 'A formal approach to learning from examples', *Int. J. Man-Machine Studies* **26**, pp. 123-141.
- [5] 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, pp. 1058–1063.
- [6] P.A. Flach (1990), 'Inductive characterisation of database relations'. In *Proc. International Symposium on Methodologies for Intelligent Systems ISMIS'90*, Z.W. Ras, M. Zemankowa & M.L. Emrich (eds.), pp. 371-378, North-Holland, Amsterdam. Full version appeared as ITK Research Report no. 23.
- [7] P.A. Flach (1993), 'Predicate invention in Inductive Data Engineering'. In *Proc. European Conference on Machine Learning ECML'93*, P.B. Brazdil (ed.), Lecture Notes in Artificial Intelligence 667, Springer-Verlag, Berlin, pp. 83–94.

- [8] P.A. Flach (1995), *Conjectures: an inquiry concerning the logic of induction*, PhD thesis, Tilburg University.
- [9] P.A. Flach (1996), 'Comparing consequence relations', manuscript.
- [10] P.A. Flach (forthcoming), 'Logical approaches to abductive reasoning and learning: an overview', forthcoming.
- [11] D.M. 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.
- [12] N.R. Hanson (1958), 'The logic of discovery', *Journal of Philosophy* **55**(25): 1073–1089.
- [13] C. Harstshorne, P. Weiss & A. Burks, eds (1931-58), *Collected Papers of Charles Sanders Peirce*, Harvard University, Cambridge.
- [14] N. Helfft (1989), 'Induction as nonmonotonic inference'. In *Proc. First International Conference on Knowledge Representation and Reasoning KR'89*, Morgan Kaufmann, San Mateo, pp. 149–156.
- [15] C.G. Hempel (1943), 'A purely syntactical definition of confirmation', *Journal of Symbolic Logic* **6**(4): 122–143.
- [16] C.G. Hempel (1945), 'Studies in the logic of confirmation', *Mind* **54**(213): 1–26 (Part I); **54**(214): 97–121 (Part II).
- [17] R.E. Jennings, C.W. Chan & M.J. Dowad (1991), 'Generalised inference and inferential modelling'. In *Proc. Int. Joint Conf. on Artificial Intelligence IJCAI'91*, Morgan Kaufmann: 1046–1051.
- [18] S. Kraus, D. Lehmann & M. Magidor (1990), 'Nonmonotonic reasoning, preferential models and cumulative logics', *Artificial Intelligence* **44**, pp. 167-207.
- [19] D. Lehmann & M. Magidor (1992), 'What does a conditional knowledge base entail?', *Artificial Intelligence* **55**, pp. 1-60.
- [20] 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.
- [21] J.S. Mill (1843), *A System of Logic*. Reprinted in *The Collected Works of John Stuart Mill*, J.M. Robson (ed.), Routledge & Kegan Paul, London.
- [22] S. Muggleton & L. De Raedt (1994), 'Inductive Logic Programming: theory and methods', *Journal of Logic Programming*.
- [23] S. Muggleton & W. Buntine (1988), 'Machine invention of first-order predicates by inverting resolution', in *Proc. Fifth International Conference on Machine Learning*, J. Laird (ed.), Morgan Kaufmann, San Mateo: 339–352.
- [24] A. Tarski (1936), 'Über den Begriff der logischen Folgerung', *Actes du Congrès Int. de Philosophie Scientifique* 7: 1–11. Translated into English as 'On the concept of logical consequence'. In *Logic, Semantics, Metamathematics*, A. Tarski (1956), Clarendon Press, Oxford: 409–420.
- [25] W. Zadrozny (1991), *On rules of abduction*, IBM Research Report, August 1991.