

**Computer Science Department Technical Report  
University of California  
Los Angeles, CA 90024-1596**

**ON THE LOGIC OF DEFAULTS**

**Hector Geffner**

**July 1988  
CSD-880058**



## On the Logic of Defaults

Hector Geffner  
Cognitive Systems Lab.  
Dept. of Computer Science, UCLA  
LA, CA 90024

Technical Report  
R-110  
June, 1988

### Abstract

We present an alternative interpretation of defaults which draws on probability theory and notions of relevance. The result is a system made up of a body of six rules which appears to overcome some of the weaknesses of other non-monotonic logics proposed in AI. We also analyze several examples and discuss some of the issues that require further research.

## 1 Introduction

A main feature exhibited by commonsense reasoning is the ability to jump to conclusions which additional information might later defeat. The limitation of classical logic to handle this kind of reasoning, has in recent years prompted the development of non-monotonic logics: logics in which the addition of new axioms might render old theorems invalid (see [Ginsberg, 87]).

The usual approach for defining these logics has been to extend classical first order logic by appealing to notions such as consistency [McDermott and Doyle, 80; Reiter, 80] or minimal models [McCarthy, 80;86]. More recently however, these logics have become subject of closer scrutiny and some of their weaknesses have become more apparent (see e.g. [Reiter and Criscuolo 81; Hanks and McDermott, 86; Morris, 87]). These analyses have revealed that the interpretation of defaults provided by these formalisms is weaker than what appears to be the intended interpretation. Conclusions that appear to be implicit in a given set of defaults fail to be sanctioned and, furthermore, as no semantic account of defaults themselves is provided, it is usually not clear where the source of the difficulties lie.

We argue here that there is more to default reasoning than non-monotonicity. We say that defaults represent hard, context-dependent constraints among beliefs and, as such, obey certain laws. Our approach is to uncover such laws and incorporate them into the logic. For that purpose, and following [Geffner and Pearl, 87], we advocate an interpretation of defaults which draws on probability theory and notions of relevance. We show that not only does the resulting system of defeasible inference usually exhibits the intended preferences when dealing with interacting defaults, but that it also provides a perspective from which such preferences can be understood.

The proposed scheme is presented in section 2.1 In the rest of section 2 we analyze several examples and introduce some refinements. In section 3 we discuss some issues that require further research.

## 2 A Logic of Defeasible Inference

### 2.1 Preliminary Definitions

**Conventions.** We use roman capital letters  $A, B, \dots$  as syntactic variables standing for first order wffs, and capital italic letters  $\Gamma, K, E, \dots$  for sets of closed wffs or sentences. Object level formulas are typed in typewriter style, e.g.  $\exists x. \text{block}(x)$ . Tuples of variables are represented by  $x, y, \dots$  while  $a, b, \dots$  stand for tuples of ground terms. The symbols ' $\vdash$ ' and ' $\not\vdash$ ' stand for provability and non-provability in first order logic with equality, respectively. Material implication is represented by the symbol ' $\Rightarrow$ '. For a set  $S$  of formulas, we use  $\phi(S)$  to refer to the formula obtained by conjoining the formulas in  $S$ . When no confusion arises, we omit the  $\phi(\cdot)$  operator and write, for instance,  $\vdash \neg S$ , as a shorthand for  $\vdash \neg\phi(S)$ .

The logic we shall present will be referred as **L** and will be characterized by a body of six rules of inference. The goal of **L** is to sanction as theorems the highly likely consequences that follow from a given context. A context  $E_K$  is built from two sets of wffs: a set  $K$  of sentences presumed to be true in every conceivable situation, called the *background context*, and a set of  $E$  of facts which characterize a particular situation and referred here as the *evidential set*.

Defaults are represented in  $K$  by sentences of the form  $\forall x. A(x) \wedge \neg ab_i(x) \Rightarrow B(x)$ , where  $A$  and  $B$  are wffs with free variables among those of  $x = \{x_1, \dots, x_n\}$ , and with  $ab_i$  playing the role of McCarthy's abnormal predicate. As we assume different defaults to involve different abnormality predicates, we shall sometimes abbreviate such defaults as  $\Delta_i(x)$ . For a particular tuple  $a$  of ground terms, the formula  $A(a) \wedge \neg ab_i(a) \Rightarrow B(a)$  represents a particular default instance, sometimes abbreviated as  $\Delta_i(a)$ .

Abnormality predicates  $ab_i$  receive a special treatment in **L**. For a tuple of ground terms  $a$ , sentences of the form  $\neg ab_i(a)$  are regarded as *candidate assumptions*, i.e. they may be assumed to hold in certain contexts. When the assumption  $\neg ab_i(a)$  holds, we also say that the default instance  $\Delta_i(a)$  holds, and viceversa. A *candidate assumption set* simply refers to a finite set of candidate assumptions.

We say that a candidate assumption set  $AS$  is consistent in context  $\Gamma$ , if  $\Gamma \not\vdash \neg AS$ . A formula  $H$  derivable from a context  $\Gamma$  augmented by a consistent candidate assumption set  $AS$ , will be said to be *arguable* in such a context, and we shall refer to such a derivation as an *argument* for  $H$  in  $\Gamma$ , and to  $AS$  as the *support* of the argument.

**L** defines an irrelevance predicate  $I(\cdot)$ , which is used to certify whether it is legitimate to jump to a defeasible

conclusion in a given context. Roughly, the idea is that if  $H$  represents an assumption believed in context  $\Gamma$ , and  $E'$  represents an additional body of evidence, then belief in  $H$  is authorized to persist as long as  $E'$  does not provide additional support for  $H$ 's negation, or, as we shall say, when  $E'$  is *irrelevant* to  $\neg H$  in context  $\Gamma$ . This is captured by the following definition:

**Definition.** A set of sentences  $E'$  is said to be irrelevant to a sentence  $H$  in context  $\Gamma$ , written  $I(H; E'|\Gamma)$ , iff for any candidate assumption set  $AS$ , such that  $E', \Gamma \vdash \neg AS$  and  $E', \Gamma, AS \vdash H$ , we also have that  $\Gamma, AS \vdash H$ .

This definition of irrelevance possesses a convenient graphical interpretation we shall often exploit. For instance, fig. 1, depicts a background context  $K$  with formulas:

- (1)  $\forall x. B(x) \wedge \neg ab_1(x) \Rightarrow F(x)$
- (2)  $\forall x. P(x) \wedge \neg ab_2(x) \Rightarrow \neg F(x)$
- (3)  $\forall x. P(x) \Rightarrow B(x)$
- (4)  $\forall x. CB(x) \Rightarrow B(x)$

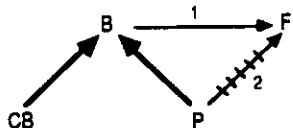


Figure 1:  $B$  separates  $CB$  from  $F$ , i.e.  $I(F(t); CB(t)|K, B(t))$

Paths in this type of graphs<sup>1</sup> correspond to arguments and, irrelevance, to a form of graph separation.<sup>2</sup> For instance, the path  $CB \rightarrow B \rightarrow F$  suggests that for any particular individual  $t$ ,  $F(t)$  is derivable from  $CB(t)$ ,  $K$ , and any support including the assumption  $\neg ab_1(t)$ . Notice that provided any such support, it is easy to verify that  $F(t)$  can be also derived from  $B(t)$  and  $K$ , what amounts to say, considering that there are no more paths from  $CB$  to  $F$ , that  $CB(t)$  is irrelevant to  $F(t)$  in context  $\{B(t)\}_K$ , i.e.  $I(F(t); CB(t)|K, B(t))$ .

Usually we will show a set of sentences  $E'$  to be irrelevant to a sentence  $H$  in a context  $E_K$ , by showing that in the corresponding graph, all the relevant paths that connect nodes corresponding to formulas in  $E'$  to the node that corresponds to  $H$ , are mediated by  $E$ . Clearly in such situations, if from a given support,  $H$  is not derivable from  $E$  and  $K$ ,  $H$  will be certainly not derivable from  $E, K$  and  $E'$ . We should keep in mind, however, that links 'contrapose'. So, a path from  $P$  to  $\neg F$  not only represents an argument for  $\neg F(t)$  given  $P(t)$ , but also an argument for  $\neg P(t)$  given  $F(t)$ . The reader might verify for instance, in

<sup>1</sup>In these graphs, we usually label the link that corresponds to default  $\Delta_i(x)$  with the index  $i$ , in order to facilitate reference.

<sup>2</sup>A similar correspondence between graph separation and conditional independence has been extensively exploited by Judea Pearl in the context of probabilistic networks (see for instance [Pearl and Verma, 87]) We borrow here some of his terminology.

the example above, that, by virtue of the different 'signs' of the links converging to  $F$ ,  $B(t)$  is relevant to  $\neg P(t)$  in  $K$ .

As for the most part the background context will remain fixed, we will find useful to abbreviate  $I(H; E'|\Gamma)$ , with  $\Gamma = E_K$ , as  $I_K(H; E'|E)$ . We also say that  $E'$  is relevant to  $H$  in context  $\Gamma$ , whenever  $I(H; E'|\Gamma)$  does not hold.

## 2.2 The Rules of Inference

The core of  $L$  is given by two sets of inference rules. We write  $\Gamma \vdash H$  to denote that sentence  $H$  is derivable from context  $\Gamma$ . Likewise,  $\Gamma, E' \vdash H$  states that  $H$  is derivable from the context that results from augmenting  $\Gamma$  with  $E'$ . Notice that the provability relation associated with the symbol ' $\vdash$ ' is not monotonic:  $\Gamma, E' \vdash H$  does not always follow from  $\Gamma \vdash H$ .

The first set of rules is given by [Geffner and Pearl, 87]:

### Rule 1 (Logic Theorems)

If  $\Gamma \vdash H$  then  $\Gamma \vdash H$

### Rule 2 (Triangularity)

If  $\Gamma \vdash H'$  and  $\Gamma \vdash H$  then  $\Gamma, H' \vdash H$

### Rule 3 (Bayes)

If  $\Gamma \vdash H'$  and  $\Gamma, H' \vdash H$  then  $\Gamma \vdash H$

### Rule 4 (Disjunction)

If  $\Gamma, H' \vdash H$  and  $\Gamma, H'' \vdash H$  then  $\Gamma, H' \vee H'' \vdash H$

It can be shown [Pearl and Geffner, 88] that the consequences of each rule are guaranteed to be highly likely whenever its premises are. Similar rules were proposed by Adams in his logic of conditionals [Adams, 66].

Hereafter, considering that the background context remains fixed for the most part, we will find useful to abbreviate  $K, E \vdash H$  as  $E \vdash_K H$ .

Rules 1–4 express how conclusions that hold in one context can be carried to a slightly different context provided certain conditions are satisfied. They do not specify however, the contexts under which candidate assumptions in  $K$ , i.e. defaults, can be assumed to hold. In particular, they do not authorize to infer that Tweety flies for instance, given that Tweety is a bird and that typically birds fly. This issue is addressed by another pair of rules, the first of which, specifies the initial context in which a given candidate assumption might be assumed to hold, while the second one uses such assumption to 'jump' to conclusions not refuted by the evidence.

Clearly, if  $\forall x. A(x) \wedge \neg ab_i(x) \Rightarrow B(x)$  is a default in  $K$ , then for a tuple of ground terms  $\mathbf{a}$ , it is reasonable to assume  $\neg ab_i(\mathbf{a})$  to hold when  $A(\mathbf{a})$  is all that is believed. Each default, however, is a belief in itself, not formed in vacuum, but on top of other relevant and irrelevant beliefs. Here we assume  $K$  to partially model such set of beliefs for every default in it, thus, authorizing for a default  $\forall x. A(x) \wedge \neg ab_i(x) \Rightarrow B(x)$ , the following inference

rule:<sup>34</sup>

### Rule 5 (Assumptions)

If  $A(\mathbf{a}), K \not\vdash \text{ab}_i(\mathbf{a})$  then  $A(\mathbf{a}), K \vdash \neg \text{ab}_i(\mathbf{a})$

As we shall see, such assumption turns out to be quite reasonable provided we restrict  $K$  to contain only statements whose truth does not depend on the particular context (e.g. “penguins are birds”), leaving in  $E$ , the context dependent information available (e.g. “Tweety flies”).<sup>5</sup>

Still, the rules above are not sufficient for maintaining derived conclusions in the presence of additional, but irrelevant information. For instance, while rules 1-5 authorize to conclude that Tweety flies, given that it is a bird and that birds typically fly, they fail to preserve such conclusion upon learning, say, Tweety’s color. This issue is addressed by an additional rule which appeals to the notion of irrelevance introduced above. The idea essentially is that a default  $\forall x.A(x) \wedge \neg \text{ab}_i(x) \Rightarrow B(x)$  permits ‘jumping’, say, from  $A(\mathbf{a})$  to  $B(\mathbf{a})$ , whenever we know the relevant assumption  $\neg \text{ab}_i(\mathbf{a})$  to hold, and the new evidence does not provide an argument supporting its negation. More precisely:

### Rule 6 (Irrelevance)

If  $\Gamma, A(\mathbf{a}) \vdash \neg \text{ab}_i(\mathbf{a})$  and  $I(\text{ab}_i(\mathbf{a}); E' | \Gamma, A(\mathbf{a}))$ ,  
then  $\Gamma, E', A(\mathbf{a}) \vdash B(\mathbf{a})$

We argue below that what matters when testing the legitimacy of inferring  $B(\mathbf{a})$  from  $A(\mathbf{a})$  in context  $\Gamma$  when coming to know  $E'$ , is not the existence of arguments in support of  $\text{ab}_i(\mathbf{a})$  but, more accurately, the existence of arguments for  $\text{ab}_i(\mathbf{a})$  in which the new information  $E'$  plays a role. These are precisely the arguments which sanction the relevance of  $E'$  to  $\text{ab}_i(\mathbf{a})$ .

In order to illustrate this last point, consider for instance a candidate assumption  $\neg \text{ab}_i(\mathbf{a})$  believed in context  $\Gamma$ . Usually, in such context there would be different sets of assumptions  $AS_i$  logically inconsistent with  $\neg \text{ab}_i(\mathbf{a})$ . For instance, in a context including information about the flying abilities of penguins and birds, the assumption that corresponds to the default instance “if Tweety is a penguin then it does not fly,” will be logically inconsistent with the assumption that corresponds to the default instance “if Tweety is a bird then it flies,” whenever Tweety is known to be a penguin. In such cases, *independently* of the new information  $E'$ , any argument whose support includes any of the sets of assumptions  $AS_i$  inconsistent with  $\neg \text{ab}_i(\mathbf{a})$  in  $\Gamma$ , will automatically constitute an argument for  $\text{ab}_i(\mathbf{a})$  in the context  $\{\Gamma, E'\}$ . What the definition of irrelevance above simply does, is not to take those arguments into account: for  $E'$  to be relevant to  $\text{ab}_i(\mathbf{a})$  in  $\Gamma$ , there has to

<sup>3</sup>Note that rule 5 permits deriving  $B(\mathbf{a})$  from  $A(\mathbf{a})$ , but not  $\neg A(\mathbf{a})$  from  $\neg B(\mathbf{a})$ . What amounts to say that the two logically equivalent sentences  $\forall x.A(x) \wedge \neg \text{ab}_i(x) \Rightarrow B(x)$  and  $\forall x.\neg B(x) \wedge \neg \text{ab}_i(x) \Rightarrow \neg A(x)$  are interpreted by  $L$  as encoding two different defaults. More about default contraposition in sections 3 and 4.

<sup>4</sup>The consistency test is for discarding from  $K$  some of the default instances otherwise implicit in the default ‘schemas’ in  $K$ , and its role should not be confused with the role consistency plays in other formalisms (e.g. [Reiter, 80; McDermott and Doyle, 80]). That convention allows us to write a ‘unique-name hypothesis’, for instance, as:  $\forall x.\forall y.\neg \text{ab}_i(x, y) \Rightarrow x \neq y$ , without implying those default instances in which  $x = y$ .

<sup>5</sup>More about this distinction in section 3.

be an argument for  $\text{ab}_i(\mathbf{a})$  with a support  $AS'$ , logically consistent with  $\neg \text{ab}_i(\mathbf{a})$  in  $\Gamma$ .

Note that in particular, if  $\neg \text{ab}_i(\mathbf{a})$  represents an assumption believed in  $\Gamma$  and  $\neg \text{ab}_k(\mathbf{b})$  represents an assumption logically inconsistent with  $\neg \text{ab}_i(\mathbf{a})$  in  $\Gamma$ , i.e.  $\Gamma, \neg \text{ab}_i(\mathbf{a}) \vdash \text{ab}_k(\mathbf{b})$ , not only does  $L$  authorize to ‘ignore’ the default instance  $\Delta_k(\mathbf{b})$  corresponding to  $\neg \text{ab}_k(\mathbf{b})$  as long as  $\neg \text{ab}_i(\mathbf{a})$  is believed<sup>6</sup>, but to ignore such default instance even in order to evaluate the relevance of new information to  $\text{ab}_i(\mathbf{a})$ . We say in those cases that  $\neg \text{ab}_i(\mathbf{a})$  *dominates* the assumption  $\neg \text{ab}_k(\mathbf{b})$  in  $\Gamma$ , and thus, the default instance  $\Delta_k(\mathbf{b})$ .

Finally, we summarize a couple of meta-theorems that follow from the rules above, we shall later appeal to:<sup>7</sup>

### Theorem 1 (Logical Closure)

If  $E \vdash_K H$ ,  $E \vdash_K H'$ , and  $H, H' \vdash H''$ , then  $E \vdash_K H''$ .

### Theorem 2 (Exceptions)

If  $E \vdash_K H$  and  $E, H' \vdash_K \neg H$  then  $E \vdash_K \neg H'$ .

## 2.3 Examples

**Example 1.** Let us first consider a background context  $K$  in which it is known that both penguins (P) and circus-birds (CB) are birds (B), and that most birds fly (F), though most penguins do not (Fig. 1):

$$\begin{aligned} \forall x.B(x) \wedge \neg \text{ab}_1(x) &\Rightarrow F(x) \\ \forall x.P(x) \wedge \neg \text{ab}_2(x) &\Rightarrow \neg F(x) \\ \forall x.P(x) &\Rightarrow B(x) \\ \forall x.CB(x) &\Rightarrow B(x) \end{aligned}$$

Let us now say we learn about a penguin called Tim. We can then conclude by means of rule 5 that  $\neg \text{ab}_2(\text{Tim})$  holds in context  $\{P(\text{Tim})\}_K$ , i.e.  $P(\text{Tim}) \vdash_K \neg \text{ab}_2(\text{Tim})$ . Likewise, being  $L$  closed under logical implication (Theorem 1), we can further conclude  $P(\text{Tim}) \vdash_K \neg F(\text{Tim})$ .

Note that extending the context  $\{P(\text{Tim})\}_K$  to include  $B(\text{Tim})$ , does not affect either conclusion since, by means of rule 2 and the fact that  $P(\text{Tim}) \vdash_K B(\text{Tim})$  follows (rule 1), formulas that hold in context  $\{P(\text{Tim})\}_K$ , can also be shown to hold in the enhanced context  $\{P(\text{Tim}), B(\text{Tim})\}_K$ . In particular thus, we obtain

$$P(\text{Tim}), B(\text{Tim}) \vdash_K \neg F(\text{Tim}).$$

$L$  does not authorize reasoning in the opposite direction though. While  $B(\text{Tim}) \vdash_K \neg \text{ab}_1(\text{Tim})$  and, as a consequence,  $B(\text{Tim}) \vdash_K F(\text{Tim})$  can be derived, the conclusion  $B(\text{Tim}), P(\text{Tim}) \vdash_K \neg F(\text{Tim})$  cannot. Nor is  $P(\text{Tim})$  irrelevant to  $\text{ab}_1(\text{Tim})$  in context  $\{B(\text{Tim})\}_K$ , as the presence of an argument for  $\text{ab}_1(\text{Tim})$  in  $\{B(\text{Tim}), P(\text{Tim})\}_K$  with support  $\{\neg \text{ab}_2(\text{Tim})\}$  suggests, nor is  $P(\text{Tim})$  a consequence of  $B(\text{Tim})$ .

Interestingly, we also have that, in context  $\{P(\text{Tim})\}_K$ , the assumption  $\neg \text{ab}_1(\text{Tim})$  is *dominated* by the assumption  $\neg \text{ab}_2(\text{Tim})$ . That is, we have both

$$P(\text{Tim}) \vdash_K \neg \text{ab}_2(\text{Tim})$$

and

$$P(\text{Tim}), \neg \text{ab}_2(\text{Tim}), K \vdash \text{ab}_1(\text{Tim}).$$

<sup>6</sup>Since, in such case, we can show  $\Gamma \vdash \text{ab}_k(\mathbf{b})$  by means of rules 1 and 3.

<sup>7</sup>See [Geffner and Pearl, 87] for proofs.

Rule 6, as we discussed above, can then be understood as asserting that the default instance  $\Delta_1(\text{Tim})$  can be ignored in order to evaluate whether it is legitimate to ‘jump’ from  $P(\text{Tim})$  to  $\neg F(\text{Tim})$  in the presence of new facts. Or, more graphically, that the link connecting **B** to **F**, in what **Tim** is concerned, can be ignored as long as  $\neg \text{ab}_2(\text{Tim})$  is believed. In particular then, we have that  $\text{CB}(\text{Tim})$  turns out to be irrelevant to  $\text{ab}_2(\text{Tim})$  in context  $\{P(\text{Tim})\}_K$  and, thus, we obtain

$$P(\text{Tim}), \text{CB}(\text{Tim}) \vdash_K \neg F(\text{Tim}).$$

**L** might be also regarded as legitimizing a weak form of contraposition. For instance, by virtue of Theorem 3 and the fact that we can derive both

$$B(\text{Tim}) \vdash_K F(\text{Tim})$$

and

$$B(\text{Tim}), P(\text{Tim}) \vdash_K \neg F(\text{Tim}),$$

we have that

$$B(\text{Tim}) \vdash_K \neg P(\text{Tim}),$$

also follows. That is, if we assume a bird to fly, though we know that penguin-birds do not fly, we are implicitly assuming that the bird is not a penguin. Stronger forms of contraposition, as deriving  $\neg B(\text{Tim})$  from  $\neg F(\text{Tim})$  however, are not sanctioned by **L**.

**Example 2.** Consider the background context  $K$  given by the defaults:

$$\forall x. P(x) \wedge \neg \text{ab}_1(x) \Rightarrow Q(x)$$

$$\forall x. Q(x) \wedge \neg \text{ab}_2(x) \Rightarrow R(x)$$

$$\forall x. S(x) \wedge \neg \text{ab}_3(x) \Rightarrow \neg R(x)$$

Clearly, for an individual **a**, we can derive  $P(\mathbf{a}) \vdash_K \neg \text{ab}_1(\mathbf{a})$  and, thus,  $P(\mathbf{a}) \vdash_K Q(\mathbf{a})$ . It turns out however, that the conclusion  $Q(\mathbf{a})$  results defeated if  $\neg R(\mathbf{a})$  is learned in such context. This is due to the fact that  $\neg R(\mathbf{a})$  does provide an argument for  $\text{ab}_1(\mathbf{a})$  supported by the assumption  $\neg \text{ab}_2(\mathbf{a})$ , and thus,  $I_K(\text{ab}_1(\mathbf{a}); \neg R(\mathbf{a}) | P(\mathbf{a}))$  does not hold.

What this indicates is that while **L** does not consider default contrapositives to be strong enough as to authorize deriving the negation of the antecedent from the negation of the consequent, **L** does consider default contrapositives to be strong enough to make the latter relevant to the former, and thus precluding certain inferences to take place. In terms of Nute [86], contrapositives are treated in **L** only as *defeaters*.

Indeed, not only does **L** preclude deriving  $Q(\mathbf{a})$  from  $P(\mathbf{a})$  when  $\neg R(\mathbf{a})$  is learned, but even when  $S(\mathbf{a})$  is. We find this latter type of behavior counterintuitive though.<sup>8</sup> In the next subsection we shall propose a refinement of the definition of the irrelevance predicate  $I(\cdot)$  given above which distinguishes between the two situations.

## 2.4 Contrapositives

The way **L** handles contraposition of defaults departs from other frameworks known to the author. Except for a weak form of contraposition, **L** does not permit to infer the negation of a default antecedent from the negation of its con-

<sup>8</sup>This type of behavior is also exhibited by circumscription and by Reiter’s default logic, when defaults are encoded as to allow contraposition (see [Morris, 87]).

sequent, though it makes the latter relevant to the former, thus precluding certain dubious derivations to take place.

Still, as we discussed above, contrapositives appear sometimes to interfere with derivations that appear to be intuitively valid. These situations usually arise from the conflict of two ‘expectation-evoking’ defaults with incompatible consequents. Here we propose a simple refinement of the definition of irrelevance given above, which draws on the ideas of [Pearl, 88a], and which leaves those derivations undisturbed.

Pearl essentially argues that causality should play a distinctive role in default reasoning, and that, in particular, reasoning chains involving ‘expectation-evoking’ defaults (e.g. “if it rained, the grass is wet”) followed by ‘explanation-evoking’ defaults (e.g. “if the grass is wet, the sprinkler was on”) should not be authorized.

In our case, due to the fact that we assume defaults to be ‘expectation-evoking,’ and their contrapositives to be ‘explanation-evoking,’<sup>9</sup> all we need to do, in order to enforce Pearl’s maxim, is to prevent such chains of reasoning when computing the irrelevance predicate  $I(\cdot)$ .

The definition of  $I(\cdot)$  above, amounts to sanction a set of sentences  $E'$  to be relevant to a sentence  $H$  in context  $E_K$ , whenever there is an argument for  $H$  in context  $\{E \cup E'\}_K$  with support  $AS$ , consistent with  $E_K$ . The extra requirement we add is simply that, whenever  $\Delta_j(\mathbf{a})$  and  $\Delta_k(\mathbf{b})$  represent two ‘expectation-evoking’ default instances with consequents inconsistent in  $K$ , then  $AS$  does not simultaneously include the assumptions  $\neg \text{ab}_j(\mathbf{a})$  and  $\neg \text{ab}_k(\mathbf{b})$ .

This simple proviso significantly improves the original account of irrelevance given above, and, in particular, correctly accounts for the type of counterintuitive behavior mentioned above.

From now on we will use  $I(\cdot)$  to stand for this improved definition and will refer to the pair of conflicting ‘expectation-evoking’ defaults as forming a *causal fork*. The new definition can then be understood simply as preventing relevance to ‘flow’ through causal forks. We will also refer to the assumptions that correspond to defaults forming a causal fork, as *conflicting assumptions*.

We illustrate next how such refinement endows **L** with the ability to properly handle the “Yale Shooting Problem.”<sup>10</sup>

**Example 3.** We consider next a version of the now famous “Yale Shooting Problem,” presented in [Hanks and McDermott, 86] as an example in which both Reiter’s logic and circumscription yield weaker conclusions than expected.

The puzzle says that people alive ( $A(t)$ ) typically remain alive ( $A(t+1)$ ) unless shot ( $S(t)$ ) with a loaded gun ( $L(t)$ ). Likewise, loaded guns ( $L(t)$ ) typically remain loaded ( $L(t+1)$ ).<sup>11</sup>

$$\forall t. L(t) \wedge \neg \text{ab}_1(t) \Rightarrow L(t+1)$$

$$\forall t. A(t) \wedge \neg \text{ab}_2(t) \Rightarrow A(t+1)$$

$$\forall t. S(t) \wedge L(t) \wedge \neg \text{ab}_3(t) \Rightarrow \neg A(t+1)$$

<sup>9</sup>Poole [87] makes a similar assumption.

<sup>10</sup>See also [Pearl, 88a].

<sup>11</sup>For clarity, we do not follow Hanks’ and McDermott’s use of a reified situation calculus. The formulation we use appears more comprehensible to us, while still serves to illustrate the difficulty detected by Hanks and McDermott in both circumscription and Reiter’s default logic.

$$\forall t. S(t) \wedge L(t) \Rightarrow ab_2(t)$$

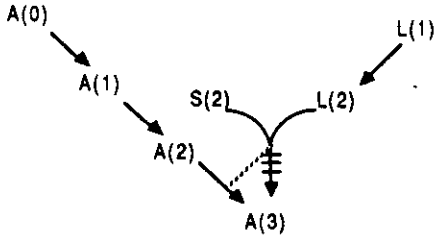


Figure 2: A version of the "Yale Shooting Problem"

We want to show that the person in question, called Fred, will most likely stop being alive if he is shot at time  $t = 2$ , with a gun loaded at  $t = 1$ , even if he was alive at time  $t = 0$ , i.e. we want to prove  $A(0), L(1), S(2) \vdash_K \neg A(3)$ . First notice that by virtue of rule 5, we have

$$S(2), L(2) \vdash_K \neg ab_3(2),$$

from which we can further infer, by means of rule 6

$$A(0), L(1), S(2), L(2) \vdash_K \neg A(3).$$

This in turn follows from the fact that

$$I_K(ab_3(2); A(0), L(1) | S(2), L(2))$$

holds, as a result of the assumption  $\neg ab_2(2)$  being dominated by the assumption  $\neg ab_3(2)$  in context  $\{S(2), L(2)\}_K$ .

Similar results would be obtained by circumscription and default logic. The behavior of  $L$  departs from these formalisms however, in which  $L$  is capable of further establishing<sup>12</sup>  $L(2)$  from  $L(1)$ ,  $S(2)$  and  $A(0)$ , and, thus, by rule 3, the expected conclusion

$$A(0), L(1), S(2) \vdash_K \neg A(3).$$

Notice first that, by means of rule 5, we have that

$$L(1) \vdash_K \neg ab_1(1).$$

In order to evaluate whether  $L(2)$  can be concluded upon learning both  $S(2)$  and  $A(0)$ , we need to test whether  $I_K(ab_1(1); A(0), S(2) | L(1))$  holds. In particular, we need to verify whether there is an argument for  $ab_1(1)$  in the resulting context whose support does not include conflicting assumptions.

It turns out that the only argument for  $ab_1(1)$  in context  $\{L(1), A(0), S(2)\}_K$ , does appeal to the conflicting assumptions  $\neg ab_2(2)$  and  $\neg ab_3(2)$  in its support, and, therefore, does not render  $E' = \{A(0), S(2)\}$  relevant to  $ab_1(1)$  in  $\{L(1)\}_K$ . It follows then by rule 6 that

$$A(0), L(1), S(2) \vdash_K L(2)$$

which, together with the previous result, leads by means of rule 3, to the expected conclusion

$$A(0), L(1), S(2) \vdash_K \neg A(3).$$

Let us finally remark that the derivation presented does not rest on a preference for 'reasoning forward' in time as opposed to 'reasoning backwards'. In particular, had we learned in addition that Fred survived the shot ( $A(3)$ ),  $L$  would correctly have failed to authorize the conclusion that the gun was loaded at the time of the shooting ( $L(2)$ ).

<sup>12</sup>Thanks to the improved definition of  $I(\cdot)$ . Otherwise  $L$  would have exhibited the same limitation.

### 3 Discussion

We have presented a system of defeasible inference motivated on probabilistic grounds and notions of relevance which provides an alternative basis for default reasoning. We have illustrated through examples how such an approach appears to overcome some of the weaknesses exhibited by other non-monotonic logics proposed in AI. In this section we want to propose some refinements and discuss some of the open issues.

**1. Supported Propositions** Circumscription and default logic appeal to either minimal models or fixed-point constructions in order to characterize the set of defeasible conclusions authorized in a given context. In particular, formulas that hold in a minimal model or extension of a given default theory, represents propositions which enjoy certain degree of support, while formulas which hold in none, stand for propositions with no support at all.

The classical example, is the "Nixon diamond:" we know quakers to be pacifists, republicans to be non-pacifists and Nixon to be both a quaker and a republican. Neither circumscription nor default logic expresses in this case any preference for believing either that Nixon is a pacifist or that he is a non-pacifist. Still, both formalisms distinguish between "Nixon is a pacifist," and, say, "Nixon is a soccer fan." The first proposition fails to be sanctioned because of conflicting evidence; the second, due to lack of support.

$L$  does not appeal to either minimal models or fixed-point constructions and, therefore, does not account for such a distinction: neither proposition is derivable by its rules.<sup>13</sup> Still, a simple account for such a distinction can be constructed on top of  $L$ . Let us say that a proposition  $H$  is supported in context  $E_K$ , if there is a candidate assumption set  $AS$  not ruled out by the evidence, i.e.  $E \not\vdash_K \neg AS$ , such that  $E, AS \vdash_K H$ .

From such a definition it is possible to show that "Nixon is a pacifist" is supported, while "Nixon is a soccer fan" is not. More interestingly, it can be shown by means of Theorem 2, that if  $H$  is derivable from  $E_K$ , then no proposition inconsistent with  $H$  in such a context will be regarded as supported.

**2. Background and Evidence.**  $L$  naturally handles implicit exceptions. We have seen in the example 1 how subclasses override conflicting superclasses properties, without having explicated the corresponding 'abnormalities.' This results from the probabilistic interpretation of defaults embedded in the rules of  $L$ , together with the distinction between the formulas taking part of the background  $K$  from the formulas taking part of the evidential set  $E$ .

The latter distinction is specially important; "penguins," for instance, would not have overridden "birds" with respect to "flying" if we had stated the fact that "penguins are birds" in  $E$  rather than in  $K$  (see [Geffner and Pearl, 87]). As we pointed out in section 2,  $K$  should contain those sentences whose truth does not depend on context, and "penguins are birds" is one such sentence.<sup>14</sup>

<sup>13</sup>This point was raised by D. Etherington in relation to [Geffner and Pearl, 87].

<sup>14</sup>There are other frameworks for default reasoning that have appealed to distinctions of this sort. Two such examples are Poole's [85] scheme for comparing conflicting default theories

We might also regard  $K$  as defining the vocabulary which is used in  $E$  to characterize a particular context. As such,  $K$  encodes information about classes with no commitment at all about what their members are.

This distinction for instance, in the framework of inheritance networks, amounts to include in  $K$  the expressions that correspond to links among classes, leaving in  $E$  those which correspond to links connecting individuals to classes.

The question that remains to be addressed is whether such criterions for distinguishing  $K$  from  $E$  are sufficient for validating rule 5. While a number of examples here and in [Geffner and Pearl 87] suggest so, we expect a more general answer to evolve.

**3. Soundness.** Rules 1–4 represent the core of  $L$ . They share the inferential power of a probabilistic sound and complete system of rules developed by Adams [66] for capturing what he called the probabilistic consequences of a set of indicative conditionals.<sup>15</sup> The addition of rules 5 and 6 amounts to augmenting the probabilistic interpretation of defaults embedded in rules 1–4 with a set of conditional independence assumptions, drawn on the basis of the syntactic structure of the knowledge base.

Other syntactic and non-syntactic means of determining reasonable conditional independence assumptions must be possible. We have illustrated for instance how a refinement of the definition of  $I(\cdot)$  originally provided, which takes into account the nature of the defaults involved, permitted certain reasonable inferences, otherwise precluded, to take place. Further refinements might be needed in order to capture other subtle aspects associated with causal defaults.

Another aspect that remains open, is a characterization of the provable consistent theories in the light of  $L$ . Though we expect such characterization to comprise most of the default ‘benchmarks’ reported in the literature, we are specially interested in those theories which can be mapped to graphs, and in which, reasoning, even in the presence of inconsistency, can be done ‘meaningfully’ and efficiently.

## Acknowledgments

Many of the intuitions that led to this work originated in conversations with Judea Pearl. I also want to thank him, M. Fuenmayer and M. Goldszmidt for comments on earlier drafts of this paper.

This work was partially supported by the National Science Foundation grant IRI 86-10155.

## References

- [Adams, 66] Adams E., “Probability and the Logic of Conditionals”, in *Aspects of Inductive Logic*, J. Hintikka and P. Suppes (Eds), North Holland Publishing Company, Amsterdam, 1966.
- [Delgrande, 87] Delgrande J., “An Approach to Default Reasoning Based on a First-Order Conditional Logic”. *Proceedings AAAI-87*, Seattle, 1987, pp 340-345.
- [Geffner and Pearl, 87] Geffner H. and Pearl J., “A Framework for Reasoning with Defaults”, TR-94b, October 1987, Cognitive Systems Lab., UCLA.
- [Ginsberg, 87] Ginsberg M., editor. *Readings in Non-Monotonic Logics*, Morgan Kaufmann, Palo Alto, 1987.
- [Hanks and McDermott, 86] Hanks S. and McDermott D., “Default Reasoning, Non-Monotonic Logics, and the Frame Problem”. *Proceedings AAAI-86*, Philadelphia, 1986, pp 328-333.
- [McCarthy 80] McCarthy J., “Circumscription—A Form of Non-Monotonic Reasoning”, *Artificial Intelligence 13*, 1980, pp 27-39.
- [McCarthy 86] McCarthy J., “Applications of Circumscription to Formalizing Commonsense Knowledge”, *Artificial Intelligence 28*, 1986, pp 89-116.
- [McDermott and Doyle, 80] McDermott D. and Doyle J., “Non-Monotonic Logic I”, *Artificial Intelligence 13*, 1980, pp 41-72.
- [Morris 87] Morris P., “Curing Anomalous Extensions”, *Proceedings of the AAAI-87*, Seattle, 1987, pp 437-442.
- [Nute 86] Nute D., “LDR: a Logic for Defeasible Reasoning”, ACMC Research Report 01-0013, University of Georgia, Athens, 1986.
- [Pearl and Verma, 87] Pearl J. and Verma T. “The Logic of Representing Dependencies by Directed Graphs” *Proceedings AAAI-87*, Seattle, 1987, pp 374-379. Also in [Pearl, 88b].
- [Pearl, 88a] Pearl J., “Embracing Causality in Default Reasoning”, *Artificial Intelligence 35*, 1988, pp 259-271. Also in [Pearl, 88b].
- [Pearl, 88b] Pearl J., *Probabilistic Reasoning in Intelligent Systems*, Morgan Kaufmann, Los Altos, 1988.
- [Pearl and Geffner, 88] Pearl J. and Geffner H., “Probabilistic Semantics for a Subset of Default Reasoning”, TR-93-III, March 1988, Cognitive Systems Lab., UCLA. Also in [Pearl, 88b].
- [Poole, 85] Poole D., “On the Comparison of Theories: Preferring the Most Specific Explanation”, *Proceedings of the IJCAI-85*, Los Angeles, 1985, pp 144-147.
- [Poole, 87] Poole D., “Defaults and Conjectures: Hypothetical Reasoning for Explanation and Prediction”, Report CS-87-4, October 1987, University Waterloo.
- [Reiter, 80] Reiter R., “A Logic for Default Reasoning” *Artificial Intelligence 13*, 1980, pp 81-132.
- [Reiter and Criscuolo, 81] Reiter R. and Criscuolo G., “On Interacting Defaults”, *Proceedings IJCAI-81*, Vancouver, 1981, pp 270-276.

and Delgrande’s [87] extended conditional logic, in which a distinction is made between sentences expressing necessary truths from those expressing contingent truths.

<sup>15</sup>Indicative conditionals of the form  $a - b$  are interpreted by Adams as asserting that the probability of  $b$  given  $a$  is infinitesimally close to one. See also [Pearl and Geffner, 88].