

# Belief Revision: A Critique\*

Nir Friedman

Institute of Computer Science

Hebrew University

Jerusalem, 91904, Israel

[nir@cs.huji.ac.il](mailto:nir@cs.huji.ac.il)

<http://www.cs.huji.ac.il/~nir>

Joseph Y. Halpern

Computer Science Department

Cornell University

Ithaca NY 14853.

[halpern@cs.cornell.edu](mailto:halpern@cs.cornell.edu)

<http://www.cs.cornell.edu/home/halpern>

May 21, 2000

## Abstract

We examine carefully the rationale underlying the approaches to belief change taken in the literature, and highlight what we view as methodological problems. We argue that to study belief change carefully, we must be quite explicit about the “ontology” or scenario underlying the belief change process. This is something that has been missing in previous work, with its focus on postulates. Our analysis shows that we must pay particular attention to two issues that have often been taken for granted: The first is how we model the agent’s epistemic state. (Do we use a set of beliefs, or a richer structure, such as an ordering on worlds? And if we use a set of beliefs, in what language are these beliefs expressed?) We show that even postulates that have been called “beyond controversy” are unreasonable when the agent’s beliefs include beliefs about her own epistemic state as well as the external world. The second is the status of observations. (Are observations known to be true, or just believed? In the latter case, how firm is the belief?) Issues regarding the status of observations arise particularly when we consider *iterated* belief revision, and we must confront the possibility of revising by  $\varphi$  and then by  $\neg\varphi$ .

**Keywords:** Belief revision, AGM postulates, Iterated revision.

---

\*Some of this work was done while both authors were at the IBM Almaden Research Center. The first author was also at Stanford while much of the work was done. IBM and Stanford’s support are gratefully acknowledged. This work was also supported in part by NSF under grants IRI-95-03109 and IRI-96-25901, by the Air Force Office of Scientific Research under grant F49620-96-1-0323, and by an IBM Graduate Fellowship to the first author. A preliminary version of this paper appeared in L. C. Aiello, J. Doyle, and S. C. Shapiro (Eds.) *Principles of knowledge representation and reasoning : Proc. Fifth International Conference (KR ’96)*, pp. 421–431, 1996. This version will appear in the *Journal of Logic, Language, and Information*.

# 1 Introduction

The problem of *belief change*—how an agent should revise her beliefs upon learning new information—has been an active area of research in both philosophy and artificial intelligence. The problem is a fascinating one in part because it is clearly no unique answer. Nevertheless, there is a strong intuition that one wants to make *minimal* changes, and all the approaches to belief change in the literature, such as [Alchourrón, Gärdenfors, and Makinson 1985; Gärdenfors 1988; Katsuno and Mendelzon 1991a], try to incorporate this principle. However, approaches differ on what constitutes a minimal change. This issue has come to the fore with the spate of recent work on *iterated* belief revision (see, for example, [Boutilier 1996; Boutilier and Goldszmidt 1993; Darwiche and Pearl 1994; Freund and Lehmann 1994; Lehmann 1995; Levi 1988; Williams 1994]).

The approaches to belief change typically start with a collection of postulates, argue that they are reasonable, and prove some consequences of these postulates. Often a semantic model for the postulates is provided and a representation theorem is proved (of the form that every semantic model corresponds to some belief revision process, and that every belief revision process can be captured by some semantic model). Our goal here is not to introduce yet another model of belief change, but to examine carefully the rationale underlying the approaches in the literature. The main message of the paper is that describing postulates and proving a representation theorem is not enough. While it may have been reasonable when research on belief change started in the early 1980s to just consider the implications of a number of seemingly reasonable postulates, it is our view that it should no longer be acceptable just to write down postulates and give short English justifications for them. While postulates do provide insight and guidance, it is also important to describe what, for want of a better word, we call the underlying *ontology* or scenario for the belief change process. Roughly speaking, this means describing carefully what it means for something to be believed by an agent and what the status is of new information that is received by the agent. This point will hopefully become clearer as we present our critique. We remark that even though the issue of ontology is tacitly acknowledged in a number of papers (for example, in the last paragraph of [Lehmann 1995]), it rarely enters into the discussion in a significant way.<sup>1</sup> We hope to show that ontology must play a central role in all discussions of belief revision.

Our focus is on approaches that take as their starting point the postulates for belief revision proposed by Alchourrón, Gärdenfors, and Makinson (AGM from now on) [1985], but our critique certainly applies to other approaches as well, in particular, Katsuno and Mendelzon’s *belief update* [1991b]; see Section 5. The AGM approach assumes that an agent’s epistemic state is represented by a *belief set*, that is, a set  $K$  of formulas in a logical language  $\mathcal{L}$ .<sup>2</sup> What the agent learns is assumed to be characterized by some

---

<sup>1</sup>A recent manuscript by Hansson [1998a] does raise some issues of ontology in the context of justification of beliefs, but no particular ontology for belief revision is discussed.

<sup>2</sup>For example, Gärdenfors [1988, p. 21] says “A simple way of modeling the epistemic state of an individual is to represent it by a *set* of sentences.”

formula  $\varphi$ , also in  $\mathcal{L}$ ;  $K * \varphi$  describes the belief set of an agent that starts with belief set  $K$  and learns  $\varphi$ .<sup>3</sup>

There are two assumptions implicit in this notation:

- The functional form of  $*$  suggests that all that matters regarding how an agent revises her beliefs is the belief set and what is learnt.
- The notation suggests that the second argument of  $*$  can be an arbitrary formula in  $\mathcal{L}$ . But what does it mean to revise by *false*? In what sense can *false* be learnt? More generally, is it reasonable to assume that an arbitrary formula can be learnt in a given epistemic state?

The first assumption is particularly problematic when we consider the postulates that AGM require  $*$  to satisfy. These essentially state that the agent is consistent in her choices, in the sense that she acts as though she has an ordering on the strength of her beliefs [Gärdenfors and Makinson 1988; Grove 1988], or an ordering on possible worlds [Boutilier 1994; Grove 1988; Katsuno and Mendelzon 1991b], or some other predetermined manner of choosing among competing beliefs [Alchourrón, Gärdenfors, and Makinson 1985]. However, the fact that an agent's epistemic state is characterized by a collection of formulas means that the epistemic state cannot include information about relative strength of beliefs (as required for the approach of, say, [Gärdenfors and Makinson 1988]), unless this information is expressible in the language. Note that if  $\mathcal{L}$  is propositional logic or first-order logic, such information cannot be expressed. On the other hand, if  $\mathcal{L}$  contains *conditional* formulas of the form  $p > q$ , interpreted as “if  $p$  is learnt, then  $q$  will be believed”, then constraints on the relative strength of beliefs can be expressed (indirectly, by describing which beliefs will be retained after a revision).

Problems arise when the language is not rich enough to express relative degrees of strength in beliefs. Consider, for example, a situation where  $K = Cl(p \wedge q)$  (the logical closure of  $p \wedge q$ ; that is, the agent's beliefs are characterized by the formula  $p \wedge q$  and its logical consequences), and then the agent learns  $\varphi = \neg p \vee \neg q$ . We can imagine that an agent whose belief in  $p$  is stronger than her belief in  $q$  would have  $K * \varphi = Cl(p)$ . That is, the agent gives up her belief in  $q$ , but retains a belief in  $p$ . On the other hand, if the agent's belief in  $q$  is stronger than her belief in  $p$ , it seems reasonable to expect that  $K * \varphi = Cl(q)$ . This suggests that it is unreasonable to take  $*$  to be a function if the representation language is not rich enough to express what may be significant details of an agent's epistemic state.

We could, of course, assume that information about the relative strength of beliefs in various propositions is implicit in the choice of the revision operator  $*$ , even if it is

---

<sup>3</sup>One of the reviewers considered the usage of the terms *learns* and *observes* as introducing a bias regarding the nature of the belief revision process. We feel that this comment highlights the need for a more concrete ontology for belief revision. We continue to use these terms, since they are quite standard in the literature, but we admit that they are indeed biased. There may well be ontologies for which they are inappropriate.

not contained in the language. This is perfectly reasonable, and also makes it more reasonable that  $*$  be a function. However, note that we can then no longer assume that we use the same  $*$  when doing iterated revision, since there is no reason to believe that the relative strength of beliefs is maintained after we learn a formula. In fact, in a number of recent papers [Boutilier 1996; Boutilier and Goldszmidt 1993; Friedman and Halpern 1998; Nayak 1994; Spohn 1988; Williams 1994],  $*$  is defined as a function from (epistemic states  $\times$  formulas) to epistemic states, but the epistemic states are no longer just belief sets; they include information regarding relative strengths of beliefs. The revision function on epistemic states induces a mapping from (belief sets  $\times$  formulas) to belief sets, but at the level of belief sets, the mapping may not be functional; for a belief set  $K$  and formula  $\varphi$ , the belief set  $K * \varphi$  may depend on what epistemic state induced  $K$ . Thus, the effect of  $*$  on belief sets may change over time.<sup>4</sup>

There is certainly no agreement on what postulates belief change should satisfy. However, the following two postulates are almost universal:

- $\varphi \in K * \varphi$
- if  $K$  is consistent and  $\varphi \in K$ , then  $K * \varphi = K$ .

These postulates have been characterized by Rott [1989] as being “beyond controversy”. Nevertheless, we argue that they are not as innocent as they may at first appear.

The first postulate says that the agent believes the last thing she learns. Making sense of this requires some discussion of the underlying ontology. It certainly makes sense if we take (as Gärdenfors [1988] does) the belief set  $K$  to consist of all formulas that are accepted by the agent (where “accepted” means “treated as true”), and the agent revises by  $\varphi$  only if  $\varphi$  has somehow come to be accepted. However, note that deciding when a formula has come to be accepted is nontrivial. In particular, just observing  $\varphi$  will not in general be enough for accepting  $\varphi$ . Acceptance has a complex interaction with what is already believed. For example, imagine a scientist who believes that heavy objects drop faster than light ones, climbs the tower of Pisa, drops a 5 kilogram textbook and a 500 milligram novel, and observes they hit the ground at the same time. This scientist will probably not accept that the time for an object to fall to the ground is independent of its weight, on the basis of this one experiment (although perhaps repeated experiments may lead her to accept it).

While the acceptance point of view is certainly not unreasonable, the fact that just observing  $\varphi$  is not necessarily enough for acceptance often seems forgotten. It also seems hard to believe that *false* would ever be accepted. More generally, it is far from obvious that in a given epistemic state  $K$  we should allow arbitrary consistent formulas to be accepted. Intuitively, this does not allow for the possibility that some beliefs are held

---

<sup>4</sup>Freund and Lehmann [1994] have called the viewpoint that  $*$  may change over time the *dynamic* point of view. However, this seems somewhat of a misnomer when applied to papers such as [Boutilier 1996; Boutilier and Goldszmidt 1993; Friedman and Halpern 1998; Williams 1994], since there  $*$  in fact is static, when viewed as a function on epistemic states and formulas.

so firmly that their negations could never be accepted. (Later we describe an ontology where observations are taken to be known in which in fact some consistent formulas will not be accepted in some epistemic states.)

While the second postulate is perhaps plausible if we cannot talk about epistemic importance or strength of belief in the language, it is less so once we can talk about such things (or if either epistemic importance or strength of belief is encoded in the epistemic state some other way). For suppose that  $\varphi \in K$ . Why should  $K * \varphi = K$ ? It could well be that being informed of  $\varphi$  raises the importance of  $\varphi$  in the epistemic ordering, or the agent's strength of belief in  $\varphi$ . If strength of belief can be talked about in the language, then a notion of minimal change should still allow strengths of belief to change, even when something expected is observed. Even if we cannot talk about strength of belief in the language, this observation has an impact on iterated revisions. For example, one assumption made by Lehmann [1995] (his postulate I4) is that if  $p$  is believed after revising by  $\varphi$ , then revising by  $[\varphi \cdot p \cdot \psi]$ —that is, revising by  $\varphi$  then  $p$  then  $\psi$ —is equivalent to revising by  $[\varphi \cdot \psi]$ . But consider a situation where after revising by  $\varphi$ , the agent believes both  $p$  and  $q$ , but her belief in  $q$  is stronger than her belief in  $p$ . We can well imagine that after learning  $\neg p \vee \neg q$  in this situation, she would believe  $\neg p$  and  $q$ . However, if she first learned  $p$  and then  $\neg p \vee \neg q$ , she would believe  $p$  and  $\neg q$ , because, as a result of learning  $p$ , she would give  $p$  higher epistemic importance than  $q$ . In this case, we would not have  $[\varphi \cdot p \cdot (\neg p \vee \neg q)] = [\varphi \cdot (\neg p \vee \neg q)]$ . In light of this discussion, it is not surprising that the combination of the second postulate with a language that can talk about epistemic ordering leads to technical problems such as Gärdenfors' *triviality result* [1988].

To give a sense of our concerns here, we discuss two basic ontologies. The first ontology that seems (to us) reasonable assumes that the agent has some knowledge as well as beliefs. We can think of the formulas that the agent knows as having the highest state of epistemic importance. In keeping with the standard interpretation of knowledge, we also assume that the formulas that the agent knows are true in the world. Since agents typically do not have certain knowledge of very many facts, we assume that the knowledge is augmented by beliefs (which can be thought of as defeasible guides to action). Thus, the set of formulas that are known form a subset of the belief set. We assume that the agent observes the world using reliable sensors; thus, if the agent observes  $\varphi$ , then the agent is assumed to know  $\varphi$ .<sup>5</sup> After observing  $\varphi$ , the agent adds  $\varphi$  to his stock of knowledge, and may revise his belief set. Since the agent's observations are taken to be knowledge, the agent will believe  $\varphi$  after observing  $\varphi$ . However, the agent's epistemic state may change even if she observes a formula that she previously believed to be true. In particular, if the formula observed was believed to be true but not known to be true,

---

<sup>5</sup>In fact, the phrase “reliable sensors” is somewhat too strong. We can deal with unreliable sensors in our framework by explicitly modeling the difference between the sensor reading the agent observes, which we take to be a reliable observation, and the actual state of the external world that the sensor is, unreliably, measuring or reporting. See [Boutilier, Friedman, and Halpern 1998] for a more detailed discussion of this point.

after the observation it is known. Note that, in this ontology, the agent never observes *false*, since *false* is not true of the world. In fact, the agent never observes anything that contradicts her knowledge. Thus,  $K * \varphi$  is defined only for formulas  $\varphi$  that are compatible with the agent’s knowledge. Moving to iterated revision, this means we cannot have a revision by  $\varphi$  followed by a revision by  $\neg\varphi$ . This ontology underlies some of our earlier work [1997, 1997, 1998]. As we show here, by taking a variant of Darwiche and Pearl’s approach [1994], we can also capture this ontology. This ontology captures the essence of Bayesian updating in probabilistic reasoning. By taking observations to be known, we are essentially giving observed events probability 1 and conditioning on them.

We can consider a second ontology that has a different flavor. In this ontology, if we observe something, we believe it to be true and perhaps even assign it a strength of belief. But this assignment does not represent the strength of belief of the observation in the resulting epistemic state. Rather, the belief in the observation must “compete” against current beliefs if it is inconsistent with these beliefs. In this ontology, it is not necessarily the case that  $\varphi \in K * \varphi$ , just as it is not the case that a scientist will necessarily adopt the consequences of his most recent observation into his stock of beliefs (at least, not without doing some additional experimentation). Of course, to flesh out this ontology, we need to describe how to combine a given strength of belief in the observation with the strengths of the beliefs in the original epistemic state. Perhaps the closest parallel in the uncertainty literature is something like the Dempster-Shafer rule of combination [Shafer 1976], which gives a rule for combining two separate bodies of belief. We believe that this type of ontology deserves further study. We sketch one particular approach to modeling this ontology in Section 5 and refer the reader to [Boutilier, Friedman, and Halpern 1998] for a more thorough treatment of it. A quite different treatment of this problem, more in the spirit of the AGM approach, has been studied by Hansson [1991], Makinson [1997], and others; see [Hansson 1998b] for an overview of this work.

The rest of the paper is organized as follows. In Section 2, we review the AGM framework, and point out some problems with it. In Section 3, we consider proposals for belief change and iterated belief change from the literature due to Boutilier [1996], Darwiche and Pearl [1994], Freund and Lehmann [1994], and Lehmann [1995], and try to understand the ontology implicit in the proposal (to the extent that one can be discerned). In Section 4, we consider the first ontology discussed above in more detail. We conclude with some discussion in Section 5.

## 2 AGM Belief Revision

In this section we review the AGM approach to belief revision. As we said earlier, this approach assumes that beliefs and observations are expressed in some language  $\mathcal{L}$ . It is assumed that  $\mathcal{L}$  is closed under negation and conjunction, and comes equipped with a compact consequence relation  $\vdash_{\mathcal{L}}$  that contains the propositional calculus and satisfies the deduction theorem. The agent’s epistemic state is represented by a belief set, that

is, a set of formulas in  $\mathcal{L}$  closed under deduction. There is also assumed to be a revision operator  $*$  that takes a belief set  $K$  and a formula  $\varphi$  and returns a new belief set  $K * \varphi$ , intuitively, the result of revising  $K$  by  $\varphi$ . The following AGM postulates are an attempt to characterize the intuition of “minimal change”:

- R1.**  $K * \varphi$  is a belief set
- R2.**  $\varphi \in K * \varphi$
- R3.**  $K * \varphi \subseteq Cl(K \cup \{\varphi\})$
- R4.** If  $\neg\varphi \notin K$  then  $Cl(K \cup \{\varphi\}) \subseteq K * \varphi$
- R5.**  $K * \varphi = Cl(false)$  if and only if  $\vdash_{\mathcal{L}} \neg\varphi$
- R6.** If  $\vdash_{\mathcal{L}} \varphi \Leftrightarrow \psi$  then  $K * \varphi = K * \psi$
- R7.**  $K * (\varphi \wedge \psi) \subseteq Cl(K * \varphi \cup \{\psi\})$
- R8.** If  $\neg\psi \notin K * \varphi$  then  $Cl(K * \varphi \cup \{\psi\}) \subseteq K * (\varphi \wedge \psi)$

The essence of these postulates is the following. Revising  $K$  by  $\varphi$  gives a belief set (Postulate R1) that includes  $\varphi$  (R2). If  $\varphi$  is consistent with  $K$ , then  $K * \varphi$  consist precisely of those beliefs implied by the combination of the old beliefs with the new belief (R3 and R4). Note that it follows from R1–R4 that if  $\varphi \in K$ , then  $K = K * \varphi$ . The next two postulates discuss the coherence of beliefs. R5 states that as long as  $\varphi$  is consistent, then so is  $K * \varphi$ , and R6 states that the syntactic form of the new belief does not affect the revision process. The last two postulates enforce a certain coherency on the outcome of revisions by related beliefs. Basically they state that if  $\psi$  is consistent with  $K * \varphi$  then  $K * (\varphi \wedge \psi)$  is the result of adding  $\psi$  to  $K * \varphi$ .

The intuitions described by AGM is based on one-step (noniterated) revision. Nevertheless, the AGM postulates do impose some restrictions on iterated revisions. For example, suppose that  $q$  is consistent with  $K * p$ . Then, according to R2 and R3,  $(K * p) * q = Cl(K * p \cup \{q\})$ . Using R7 and R8 we can conclude that  $(K * p) * q = K * (p \wedge q)$ .

There are several representation theorems for AGM belief revision; perhaps the easiest to understand is due to Grove [1988]. We discuss a slight modification, due to Boutilier [1994] and Katsuno and Mendelzon [1991b]: Let an  $\mathcal{L}$ -world be a complete and consistent truth assignment to the formulas in  $\mathcal{L}$ . Let  $\mathcal{W}$  consist of all the  $\mathcal{L}$ -worlds, and let  $\preceq$  be a *ranking*, that is, a total preorder, on the worlds in  $\mathcal{W}$ . Let  $\text{min}_{\preceq}$  consist of all the minimal worlds with respect to  $\preceq$ , that is, all the worlds  $w$  such that there is no  $w'$  with  $w' \prec w$ . With  $\preceq$  we can associate a belief set  $\text{Bel}(\preceq)$ , consisting of all formulas  $\varphi$  that are true in all the worlds in  $\text{min}_{\preceq}$ . Moreover, we can define a revision operator  $*$ , by taking  $\text{Bel}(\preceq) * \varphi$  to consist of all formulas  $\psi$  that are true in all the minimal  $\varphi$ -worlds according to  $\preceq$ . It can be shown that  $*$  satisfies the AGM postulates (when its first argument is  $\text{Bel}(\preceq)$ ). Thus, we can define a revision operator by taking a collection of orderings  $\preceq_K$ , one for each belief set  $K$ . To define  $K * \varphi$  for a belief set  $K$ , we apply

the procedure above, starting with the ranking  $\preceq_K$  corresponding to  $K$ .<sup>6</sup> Furthermore, Grove [1988], Katsuno and Mendelzon [1991b], and Boutilier [1994] show that every belief revision operator satisfying the AGM axioms can be characterized in this way.

This elegant representation theorem also brings out some of the problems with the AGM postulates. First, note that a given revision operator  $*$  is represented by a family of rankings, one for each belief set. There is no necessary connection between the rankings corresponding to different belief sets. It might seem more reasonable to have a more global setting (perhaps one global ranking) from which each element in the family of rankings arises.

A second important point is that the epistemic state here is a function not only of the belief set, but also of the ranking. The latter, however, is represented only implicitly. Each ranking  $\preceq$  is associated with a belief set  $\text{Bel}(\preceq)$ , but it is the ranking that gives the information required to describe how revision is carried out. The belief set does not suffice to determine the revision; there are many rankings  $\preceq$  for which the associated belief set  $\text{Bel}(\preceq)$  is  $K$ . Since the revision process only gives us the revised belief set, not the revised ranking, the representation does not support iterated revision.

This suggests that we should consider, not how to revise belief sets, but how to revise rankings. More generally, whatever we take to be our representation of the epistemic state, it seems appropriate to consider how these representations should be revised. We can define an analogue of the AGM postulates for epistemic states in a straightforward way (cf. [Friedman and Halpern 1994]): Taking  $E$  to range over epistemic states and  $\text{Bel}(E)$  to represent the belief set associated with epistemic state  $E$ , we have

**R1'.**  $E * \varphi$  is an epistemic state

**R2'.**  $\varphi \in \text{Bel}(E * \varphi)$

**R3'.**  $\text{Bel}(E * \varphi) \subseteq \text{Cl}(\text{Bel}(E) \cup \{\varphi\})$

and so on, with the obvious syntactic transformation. In fact, as we shall see in the next section, a number of processes for revising epistemic states have been considered in the literature, and in fact they all do satisfy these modified postulates.

Finally, even if we restrict attention to belief sets, we can consider what happens if the underlying language  $\mathcal{L}$  is rich enough to talk about how revision should be carried out. For example, suppose  $\mathcal{L}$  includes conditional formulas, and we want to find some ranking  $\preceq$  for which the corresponding belief set is  $K$ . Not just any ranking  $\preceq$  such that  $\text{Bel}(\preceq) = K$  will do here. The beliefs in  $K$  put some constraints on the ranking. For example, if  $p > q$  is in  $K$  and  $p \notin K$ , then the minimal  $\preceq$ -worlds satisfying  $p$  must all satisfy  $q$ , since after  $p$  is learnt,  $q$  is believed. Once we restrict to rankings  $\preceq_K$  that are consistent with  $K$  then the AGM postulates are no longer sound. This point has essentially been made before [Boutilier 1992; Rott 1989]. However, it is worth stressing

---

<sup>6</sup>In this construction, for each belief set  $K$  other than the inconsistent belief set, we have  $\text{Bel}(\preceq_K) = K$ . The inconsistent belief set gets special treatment here.

the sensitivity of the AGM postulates to the underlying language and, more generally, to the choice of epistemic state.

### 3 Proposals for Iterated Revision

We now briefly review some of the previous proposals for iterated belief change, and point out how the impact of the observations we have been making on the approaches. Most of these approaches start with the AGM postulates, and augment them to get seemingly appropriate restrictions on iterated revision. This is not an exhaustive review of the literature on iterated belief revision by any stretch of the imagination. Rather, we have chosen a few representative approaches that allow us to bring out our methodological concerns.

#### 3.1 Boutilier's Approach

As we said in the previous section, Boutilier takes the agent's epistemic state to consist of a ranking of possible worlds. Boutilier [1996] describes a particular revision operator  $*_B$  on epistemic states. This revision operator maps a ranking  $\preceq$  of possible worlds and an observation  $\varphi$  to a revised ranking  $\preceq *_B \varphi$  such that (a)  $\preceq *_B \varphi$  satisfies the conditions of the representation theorem described above, that is, the minimal worlds in  $\preceq *_B \varphi$  are precisely the minimal  $\varphi$ -worlds in  $\preceq$ , and (b) in a precise sense,  $\preceq *_B \varphi$  is the result of making the minimal number of changes to  $\preceq$  required to guarantee that all the minimal worlds in  $\preceq *_B \varphi$  satisfy  $\varphi$ . Given a ranking  $\preceq$  and a formula  $\varphi$ , the ranking  $\preceq *_B \varphi$  is identical to  $\preceq$  except that the minimal  $\varphi$ -worlds according to  $\preceq$  have the minimal rank in the revised ranking, while the relative ranks of all other worlds remains unchanged.

Boutilier characterizes the properties of his approach as follows. Suppose that, starting in some epistemic state, we revise by  $\varphi_1, \dots, \varphi_n$ . Further suppose  $\varphi_{i+1}$  is consistent with the beliefs after revising by  $\varphi_1, \dots, \varphi_i$ . Then the beliefs after revising by  $\varphi_1, \dots, \varphi_n$  are precisely the beliefs after observing  $\varphi_1 \wedge \dots \wedge \varphi_n$ . (More precisely, given any ranking  $\preceq$ , the belief set associated with the ranking  $\preceq *_B \varphi_1 *_B \dots *_B \varphi_n$  is the same as that associated with the ranking  $\preceq *_B (\varphi_1 \wedge \dots \wedge \varphi_n)$ . Note, however, that  $\preceq *_B \varphi_1 *_B \dots *_B \varphi_n \neq \preceq *_B (\varphi_1 \wedge \dots \wedge \varphi_n)$  in general.) Thus, as long as the agent's new observations are not surprising, the agent's beliefs are exactly the ones she would have had had she observed the conjunction of all the observations. This is an immediate consequence of the AGM postulates, and thus holds for any approach that attempts to extend the AGM postulates to iterated revision.

What happens when the agent observes a formula  $\varphi_{n+1}$  that is inconsistent with her current beliefs? Boutilier shows that in this case the new observation nullifies the impact of the all the observations starting with the most recent one that is inconsistent with  $\varphi_{n+1}$ . More precisely, suppose  $\varphi_{i+1}$  is consistent with the belief after observing  $\varphi_1, \dots, \varphi_i$  for all  $i \leq n$ , but  $\varphi_{n+1}$  is inconsistent with the beliefs after observing

$\varphi_1, \dots, \varphi_n$ . Let  $k$  be the maximal index such that  $\varphi_{n+1}$  is consistent with the beliefs after learning  $\varphi_1, \dots, \varphi_k$ . The agent's beliefs after observing  $\varphi_{n+1}$  are the same as her beliefs after observing  $\varphi_1, \dots, \varphi_k, \varphi_{n+1}$ . Thus, the agent acts as though she did not observe  $\varphi_{k+1}, \dots, \varphi_n$ .

Boutilier does not provide any argument for the reasonableness of this ontology. In fact, Boutilier's presentation (like almost all others in the literature) is not in terms of an ontology at all; he presents his approach as an attempt to minimize changes to the ranking. While the intuition of minimizing changes to the ranking seems reasonable at first, it becomes less reasonable when we realize its ontological implications. The following example, due to Darwiche and Pearl [1994], emphasizes this point. Suppose we encounter a strange new animal and it appears to be a bird, so we believe it is a bird. On closer inspection, we see that it is red, so we believe that it is a red bird. However, an expert then informs us that it is not a bird, but a mammal. Applying Boutilier's revision operator, we would no longer believe that the animal is red. This does not seem so reasonable.

One more point is worth observing: As described by Boutilier [1996], his approach does not allow revision by *false*. While we could, of course, modify the definition to handle *false*, it is more natural simply to disallow it. This suggests that, whatever ontology is used to justify Boutilier's approach, in that ontology, revising by *false* should not make sense.

### 3.2 Freund and Lehmann's Approach

Freund and Lehmann [1994] stick close to the original AGM approach. They work with belief sets, not more general epistemic states. However, they are interested in iterated revision. They consider the effect of adding just one more postulate to the basic AGM postulates, namely

**FL.** If  $\neg\varphi \in K$ , then  $K * \varphi = K_\perp * \varphi$ ,

where  $K_\perp$  is the inconsistent belief set, which consists of all formulas.

Suppose  $*$  satisfies R1–R8 and FL. Just as with Boutilier's approach, if  $\varphi_{i+1}$  is consistent with the beliefs after learning  $\varphi_1, \dots, \varphi_i$  for  $i \leq n - 1$ , then  $K * \varphi_1 * \dots * \varphi_n = K * (\varphi_1 \wedge \dots \wedge \varphi_n)$ . However, if we then observe  $\varphi_{n+1}$ , and it is inconsistent with  $K * \varphi_1 \wedge \dots \wedge \varphi_n$ , then  $K * \varphi_1 * \dots * \varphi_{n+1} = K_\perp * \varphi_{n+1}$ . That is, observing something inconsistent causes us to retain none of our previous beliefs, but to start over from scratch. While the ontology here is quite simple to explain, as Freund and Lehmann themselves admit, it is a rather severe form of belief revision. Darwiche and Pearl's red bird example applies to this approach as well.

### 3.3 Darwiche and Pearl's Approach

Darwiche and Pearl [1994] suggest a set of postulates extending the AGM postulates, and claim to provide a semantics that satisfies them. Their intuition is that the revision operator should retain as much as possible certain parts of the ordering among worlds in the ranking. In particular, if  $w$  and  $w'$  both satisfy  $\varphi$ , then a revision by  $\varphi$  should not change the relative rank of  $w$  and  $w'$ . Similarly, if both  $w$  and  $w'$  satisfy  $\neg\varphi$ , then a revision should not change their relative rank. They describe four postulates that are meant to embody these intuitions:

- C1. If  $\varphi \vdash \psi$ , then  $(K * \psi) * \varphi = K * \varphi$
- C2. If  $\varphi \vdash \neg\psi$ , then  $(K * \psi) * \varphi = K * \varphi$
- C3. If  $\psi \in K * \varphi$ , then  $\psi \in (K * \psi) * \varphi$
- C4. If  $\neg\psi \notin K * \varphi$ , then  $\neg\psi \notin (K * \psi) * \varphi$

Freund and Lehmann [1994] point out that C2 is inconsistent with the AGM postulates. This observation seems inconsistent with the fact that Darwiche and Pearl claim to provide an example of a revision method that is consistent with their postulates. What is going on here? It turns out that the issues raised earlier help clarify the situation.

The semantics that Darwiche and Pearl use as an example is based on a special case of Spohn's *ordinal conditional functions* (OCFs) [1988] called  $\kappa$ -rankings [Goldszmidt and Pearl 1992]. A  $\kappa$ -ranking associates with each world either a natural number  $n$  or  $\infty$ , with the requirement that for at least one world  $w_0$ , we have  $\kappa(w_0) = 0$ . We can think of  $\kappa(w)$  as the rank of  $w$ , or as denoting how surprising it would be to discover that  $w$  is the actual world. If  $\kappa(w) = 0$ , then world  $w$  is unsurprising; if  $\kappa(w) = 1$ , then  $w$  is somewhat surprising; if  $\kappa(w) = 2$ , then  $w$  is more surprising, and so on. If  $\kappa(w) = \infty$ , then  $w$  is impossible.<sup>7</sup> OCFs provide a way of ranking worlds that is closely related to, but has a little more structure than the orderings considered by Boutilier (as well as Grove and Katsuno and Mendelzon). The extra structure makes it easier to define a notion of conditioning.

Given a formula  $\varphi$ , let  $\kappa(\varphi) = \min\{\kappa(w) : w \models \varphi\}$ ; we define  $\kappa(\text{false}) = \infty$ . We say that  $\varphi$  is *believed with firmness*  $\alpha \geq 0$  in OCF  $\kappa$  if  $\kappa(\varphi) = 0$  and  $\kappa(\neg\varphi) = \alpha$ . Thus,  $\varphi$  is believed with firmness  $\alpha$  if  $\varphi$  is unsurprising and the least surprising world satisfying  $\neg\varphi$  has rank  $\alpha$ . We define  $\text{Bel}(\kappa)$  to consist of all formulas that are believed with firmness at least 1.

Spohn defined a notion of conditioning on OCFs. Given an OCF  $\kappa$ , a formula  $\varphi$  such that  $\kappa(\varphi) < \infty$ , and  $\alpha \geq 0$ ,  $\kappa_{\varphi,\alpha}$  is the unique OCF satisfying the property desired by Darwiche and Pearl—namely, if  $w$  and  $w'$  both satisfy  $\varphi$  or both satisfy  $\neg\varphi$ , then revision

---

<sup>7</sup>Spohn allowed ranks to be arbitrary ordinals, not just natural numbers, and did not allow a rank of  $\infty$ , since, for philosophical reasons, he did not want to allow a world to be considered impossible. As we shall see, there are technical advantages to introducing a rank of  $\infty$ .

by  $\varphi$  should not change the relative rank of  $w$  and  $w'$ , that is,  $\kappa_{\varphi,\alpha}(w) - \kappa_{\varphi,\alpha}(w') = \kappa(w) - \kappa(w')$ —such that  $\varphi$  is believed with firmness  $\alpha$  in  $\kappa_{\varphi,\alpha}$ . It is defined as follows:

$$\kappa_{\varphi,\alpha}(w) = \begin{cases} \kappa(w) - \kappa(\varphi) & \text{if } w \text{ satisfies } \varphi \\ \kappa(w) - \kappa(\neg\varphi) + \alpha & \text{if } w \text{ satisfies } \neg\varphi. \end{cases}$$

Notice that  $\kappa_{\varphi,\alpha}$  is defined only if  $\kappa(\varphi) < \infty$ , that is, if  $\varphi$  is considered possible.

Darwiche and Pearl defined the following revision function on OCFs:

$$\kappa *_{DP} \varphi = \begin{cases} \kappa & \text{if } \kappa(\neg\varphi) \geq 1 \\ \kappa_{\varphi,1} & \text{otherwise.} \end{cases}$$

Thus, if  $\varphi$  is already believed with firmness at least 1 in  $\kappa$ , then  $\kappa$  is unaffected by a revision by  $\varphi$ ; otherwise, the effect of revision is to modify  $\kappa$  by conditioning so that  $\varphi$  ends up being believed with degree of firmness 1. Intuitively, this means that if  $\varphi$  is not believed in  $\kappa$ , in  $\kappa * \varphi$  it is believed, but with the minimal degree of firmness.

It is not hard to show that if we take an agent's epistemic state to be represented by an OCF, then Darwiche and Pearl's semantics satisfies all the AGM postulates modified to apply to epistemic states (that is, R1'–R8' in Section 2), except that revising by *false* is disallowed, just as in Boutilier's approach, so that R5' holds vacuously; in addition, this semantics satisfies Darwiche and Pearl's C1–C4, modified to apply to epistemic states. For example, C2 becomes

**C2'.** If  $\varphi \vdash \neg\psi$ , then  $\text{Bel}((E * \psi) * \varphi) = \text{Bel}(E * \varphi)$ .

Indeed, as Darwiche and Pearl observe, Boutilier's revision operator also satisfies C1'–C4'; however, it has properties that they view as undesirable. Thus, Darwiche and Pearl's claim that their postulates are consistent with AGM is correct, if we think at the level of general epistemic states. On the other hand, Freund and Lehmann are quite right that R1–R8 and C1–C4 are incompatible; indeed, as they point out, R1–R4 and C2 are incompatible. The importance of making clear exactly whether we are considering the postulates with respect to the OCF  $\kappa$  or the belief set  $\text{Bel}(\kappa)$  is particularly apparent here.<sup>8</sup>

The fact that Boutilier's revision operator also satisfies C1'–C4' clearly shows that these postulates do not capture all of Darwiche and Pearl's intuitions. Their semantics embodies further assumptions. Some of them seem *ad hoc*. Why is it reasonable to believe  $\varphi$  with a *minimal* degree of firmness after revising by  $\varphi$ ? Rather than trying to come up with an improved collection of postulates (which Darwiche and Pearl themselves suggest might be a difficult task), it seems to us that a more promising approach is to find an appropriate ontology.

---

<sup>8</sup>We note that in a recent version of their paper, Darwiche and Pearl [1997] use a similar technique to deal with the inconsistency of C2.

### 3.4 Lehmann’s Revised Approach

Finally, we consider Lehmann’s “revised” approach to belief revision [1995]. With each sequence  $\sigma$  of observations, Lehmann associates a belief set that we denote  $\text{Bel}(\sigma)$ . Intuitively, we can think of  $\text{Bel}(\sigma)$  as describing the agent’s beliefs after making the sequence  $\sigma$  of observations, starting from her initial epistemic state. Lehmann allows all possible sequences of consistent formulas. Thus, he assumes that the agent does not observe *false*. We view Lehmann’s approach essentially as taking the agent’s epistemic state to be the sequence of observations made, with the obvious revision operator that concatenate a new observation to the current epistemic state. The properties of belief change depend on the function  $\text{Bel}$ . Lehmann requires  $\text{Bel}$  to satisfy the following postulates (where  $\sigma$  and  $\rho$  denote sequences of formulas, and  $\cdot$  is the concatenation operator):

- I1.  $\text{Bel}(\sigma)$  is a consistent belief set
- I2.  $\varphi \in \text{Bel}(\sigma \cdot \varphi)$
- I3. If  $\psi \in \text{Bel}(\sigma \cdot \varphi)$ , then  $\varphi \Rightarrow \psi \in \text{Bel}(\sigma)$
- I4. If  $\varphi \in \text{Bel}(\sigma)$ , then  $\text{Bel}(\sigma \cdot \varphi \cdot \rho) = \text{Bel}(\sigma \cdot \rho)$
- I5. If  $\psi \vdash \varphi$ , then  $\text{Bel}(\sigma \cdot \varphi \cdot \psi \cdot \rho) = \text{Bel}(\sigma \cdot \psi \cdot \rho)$
- I6. If  $\neg\psi \notin \text{Bel}(\sigma \cdot \varphi)$ , then  $\text{Bel}(\sigma \cdot \varphi \cdot \psi \cdot \rho) = \text{Bel}(\sigma \cdot \varphi \cdot \varphi \wedge \psi \cdot \rho)$
- I7.  $\text{Bel}(\sigma \cdot \neg\varphi \cdot \varphi) \subseteq \text{Cl}(\text{Bel}(\sigma) \cup \{\varphi\})$

We refer the interested reader to [Lehmann 1995] for the motivation for these postulates. As Lehmann argues, the spirit of the original AGM postulates is captured by these postulates. Lehmann views I5 and I7 as two main additions to the basic AGM postulates. He states that “Since postulates I5 and I7 seem secure, i.e., difficult to reject, the postulates I1–I7 may probably be considered as a reasonable formalization of the intuitions of AGM” [Lehmann 1995, Section 5]. Our view is that it is impossible to decide whether to accept or reject postulates such as I5 or I7 (or, for that matter, any of the other postulates) without an explicit ontology. There may be ontologies for which I5 and I7 are reasonable, and others for which they are not. “Reasonableness” is not an independently defined notion; it depends on the ontology. The ontology of the next section emphasizes this point.

## 4 Taking Observations to be Knowledge

We now consider an ontology where observations are taken to be knowledge.<sup>9</sup>

---

<sup>9</sup>We remark that Rott [1991, Section 6] outlines an approach where observations are treated as knowledge, but does not provide an underlying ontology. The high-level intuition behind his approach is similar to the one we present here, although the details differ.

As we said in the introduction, in this ontology, the agent has some (closed) set of formulas that he *knows* to be true, which is included in a larger set of formulas that he *believes* to be true. The belief set can be viewed as the result of applying some nonmonotonic reasoning system grounded in the observations. We can think of there being an ordering on the strength of his beliefs, with the formulas known to be true—the observations and their consequences—having the greatest strength of belief. Because observations are taken to be knowledge, any formula observed is added to the stock of knowledge (and must be consistent with what was previously known). In this ontology, it is impossible to observe *false*. In fact, it is impossible to make any inconsistent sequence of observations. That is, if  $\varphi_1, \dots, \varphi_n$  is observed, then  $\varphi_1 \wedge \dots \wedge \varphi_n$  must be consistent (although it may not be consistent with the agent’s original beliefs).

In earlier work [Friedman 1997; Friedman and Halpern 1998], we presented one way of formalizing this ontology, based on the framework of Halpern and Fagin [1989] for modeling multi-agent systems (see [Fagin, Halpern, Moses, and Vardi 1995] for more details). For modeling belief revision, we use this framework restricted to a single agent.<sup>10</sup> The key assumption in this framework is that we can characterize the system by describing it in terms of a *state* that changes over time. Formally, we assume that at each point in time, the agent is in some *local state*. Intuitively, this local state encodes the information the agent has observed thus far. There is also an *environment*, whose state encodes relevant aspects of the system that are not part of the agent’s local state. A *global state* is a tuple  $(s_e, s_a)$  consisting of the environment state  $s_e$  and the local state  $s_a$  of the agent. A *run* of the system is a function from time (which, for ease of exposition, we assume ranges over the natural numbers) to global states. Thus, if  $r$  is a run, then  $r(0), r(1), \dots$  is a sequence of global states that, roughly speaking, is a complete description of what happens over time in one possible execution of the system. We take a *system* to consist of a set of runs. Intuitively, these runs describe all the possible behaviors of the system, that is, all the possible sequences of events that could occur in the system over time.

Given a system  $\mathcal{R}$ , we refer to a pair  $(r, m)$  consisting of a run  $r \in \mathcal{R}$  and a time  $m$  as a *point*. If  $r(m) = (s_e, s_a)$ , we define  $r_a(m) = s_a$  and  $r_e(m) = s_e$ . We say two points  $(r, m)$  and  $(r', m')$  are *indistinguishable* to the agent, and write  $(r, m) \sim_a (r', m')$ , if  $r_a(m) = r'_a(m')$ , i.e., if the agent has the same local state at both points. Finally, an *interpreted system* is a tuple  $(\mathcal{R}, \pi)$ , consisting of a system  $\mathcal{R}$  together with a mapping  $\pi$  that associates with each point a truth assignment to the primitive propositions.

To capture the AGM framework, we consider a special class of interpreted systems: We fix a propositional language  $\mathcal{L}$ . We assume that the agent makes observations, which are characterized by formulas in  $\mathcal{L}$ , and that her local state consists of the sequence of observations that she has made. We assume that the environment’s local state describes which formulas are actually true in the world, so that it is a truth assignment to the

---

<sup>10</sup>Although the multi-agent aspects of this framework does not play a role in our analysis here, we note that it allows for a natural extension of our results to multi-agent belief revision. This, however, is beyond the scope of this paper.

formulas in  $\mathcal{L}$ . As observed by Katsuno and Mendelzon [1991a], the AGM postulates assume that the world is *static*; to capture this, we assume that the environment state does not change over time. Formally, we are interested in the unique interpreted system  $(\mathcal{R}^{AGM}, \pi)$  that consists of all runs satisfying the following two assumptions for every point  $(r, m)$ :

- The environment's state  $r_e(m)$  is a truth assignment to the formulas in  $\mathcal{L}$  that agrees with  $\pi$  at  $(r, m)$  (that is,  $\pi(r, m) = r_e(m)$ ), and  $r_e(m) = r_e(0)$ .
- The agent's state  $r_a(m)$  is a sequence of the form  $\langle \varphi_1, \dots, \varphi_m \rangle$ , such that  $\varphi_1 \wedge \dots \wedge \varphi_m$  is true according to the truth assignment  $r_e(m)$  and  $r_a(m-1) = \langle \varphi_1, \dots, \varphi_{m-1} \rangle$ .

Notice that the form of the agent's state makes explicit an important implicit assumption: that the agent remembers all her previous observations.

In an interpreted system, we can talk about an agent's knowledge: the agent knows  $\varphi$  at a point  $(r, m)$  if  $\varphi$  holds in all points  $(r', m')$  such that  $(r, m) \sim_a (r', m')$ . It is easy to see that, according to this definition, if  $r_a(m) = \langle \varphi_1, \dots, \varphi_m \rangle$ , then the agent knows  $\varphi_1 \wedge \dots \wedge \varphi_m$  at the point  $(r, m)$ : the agent's observations are known to be true in this approach. We are interested in talking about the agent's beliefs as well as her knowledge. To allow this, we added a notion of *plausibility* to interpreted systems in [Friedman and Halpern 1997]. We consider a variant of this approach here, using OCFs, since it makes it easier to relate our observations to Darwiche and Pearl's framework.

We assume that we start with an OCF  $\kappa$  on runs such that  $\kappa(r) \neq \infty$  for any run  $r$ . Intuitively,  $\kappa$  represents our prior ranking on runs. Initially, no runs is viewed as impossible. We then associate, with each point  $(r, m)$ , an OCF  $\kappa^{(r,m)}$  on the runs. We define  $\kappa^{(r,m)}$  by induction on  $m$ . We take  $\kappa^{(r,0)} = \kappa$ , and we take  $\kappa^{(r,m+1)} = \kappa_{\varphi_{m+1}, \infty}^{(r,m)}$ , where  $r_a(m+1) = \langle \varphi_1, \dots, \varphi_{m+1} \rangle$ . Thus,  $\kappa^{(r,m+1)}$  is the result of conditioning  $\kappa^{(r,m)}$  on the last observation the agent made, giving it degree of firmness  $\infty$ . Thus, the agent is treating the observations as knowledge in a manner compatible with the semantics for knowledge in interpreted systems. Moreover, since observations are known, they are also believed.

Note that we could have described the same belief change process by considering an OCF on formulas, rather than runs. However, we feel that using runs and systems provides a better model of what is going on, and gives a more explicit ontology. We can then *derive* an OCF on formulas from the OCF on runs that we consider.

As we show in [Friedman 1997; Friedman and Halpern 1998], this framework satisfies the AGM postulates R1'–R8', interpreted on epistemic states. (Here we take the agent's epistemic state at the point  $(r, m)$  to consist of  $r_a(m)$  together with  $\kappa^{(r,m)}$ .) Moreover, we show this framework satisfies an additional postulate, which we call R9':

**(R9')** If  $\not\models_{\mathcal{L}} \neg(\varphi \wedge \psi)$  then  $\text{Bel}(E * \varphi * \psi) = \text{Bel}(E * \varphi \wedge \psi)$ .

This postulate captures the intuition that observations are taken to be knowledge, and thus observing  $\varphi$  and then  $\psi$  is equivalent to observing  $\varphi \wedge \psi$ . In fact, we show there that R1'–R9' characterize, in a precise sense, revision in this framework.

It is easy to verify that the framework also satisfies Darwiche and Pearl’s postulates (appropriately modified to apply to epistemic states), except that the contentious C2 is now vacuous, since it is illegal to revise by  $\psi$  and then by  $\varphi$  if  $\varphi \vdash \neg\psi$ .

How does this framework compare to Lehmann’s? Like Lehmann’s, there is an explicit attempt to associate beliefs with a sequence of revisions. However, we have restricted the sequence of revisions, since we are treating observations as knowledge. It is easy to see that I1–I3 and I5–I7 hold in our framework. However, since we have restricted the sequence of observations allowed, some of these postulates are much weaker in our framework than in Lehmann’s. In particular, I7 is satisfied vacuously, since we do not allow a sequence of the form  $\sigma \cdot \neg\varphi \cdot \varphi$ . On the other hand, I4 is not satisfied in our framework. Our discussion in the introduction suggests a counterexample. Suppose that initially,  $\kappa(p \wedge q) = 0$ ,  $\kappa(\neg p \wedge q) = 1$ ,  $\kappa(p \wedge \neg q) = 2$ , and  $\kappa(\neg p \wedge \neg q) = 3$ . Thus, initially the agent believes both  $p$  and  $q$ , but believes  $p$  with firmness 1 and  $q$  with firmness 2. If the agent then observes  $\neg p \vee \neg q$ , he will then believe  $q$  but not  $p$ . On the other hand, suppose the agent first observes  $p$ . He still believes both  $p$  and  $q$ , of course, but now  $p$  is believed with firmness  $\infty$ . That means if he then observes  $\neg p \vee \neg q$ , he will believe  $p$ , but not  $q$ , violating I4. However, a weaker variant of I4 does hold in our system: if the agent *knows*  $\varphi$ , then observing  $\varphi$  will not change her future beliefs.

Did we really need all the machinery of runs and systems here? While we could no doubt get away without it, we believe that modeling time and the agent’s state explicitly adds a great deal. In particular, by including time, we capture the belief revision process within the model. By including the agent’s state and the environment, we can capture various assumptions about the agent’s observational capabilities (this is especially relevant once we allow inaccurate observations) and generalize to allow multiple agents. See [Boutilier, Friedman, and Halpern 1998; Friedman and Halpern 1998; Friedman and Halpern 1997] for further discussion and demonstration of the advantages of this framework.

## 5 Discussion

The goal of this paper was to highlight what we see as some methodological problems in much of the literature on belief revision. There has been (in our opinion) too much attention paid to postulates, and not enough to the underlying ontology. An ontology must make clear what the agent’s epistemic state is, what types of observations the agent can make, the status of observations, and how the agent goes about revising the epistemic state.

We have (deliberately) not been very precise about what counts as an ontology, and clearly there are different levels of detail that one can provide. Although some papers on belief revision have attempted to provide something in the way of an ontology, the ontology has typically been insufficiently detailed to verify the reasonableness of the postulates. For example, Gärdenfors [1988]—whose ontology is far better developed

than most—takes belief sets to consist of formulas that are accepted, and revision to be by formulas that are accepted. As we have seen, unless belief sets or the revision operator contain additional information (such as epistemic importance or strengths of beliefs) this ontology will violate R1. On the other hand, unless we are careful about how we add such information, we may well violate some other axioms (such as R3 or R5). In any case, it is the job of the ontology to make completely clear such issues as whether observations are believed to be true or known to be true, and if they are believed, what the strength of belief is. This issue is particularly important if we have epistemic states like rankings that are richer than belief sets. If observations are believed, but not necessarily known, to be true, then it is not clear how to go about revising such a richer epistemic state. With what degree of firmness should the new belief be held? No particular answer seems to us that well motivated. It may be appropriate for the user to attach degrees of firmness to observations, as was done in [Goldszmidt 1992; Williams 1994; Wobcke 1995] (following the lead of Spohn [1988]); we can even generalize to allowing uncertain observations [Dubois and Prade 1992].

We are currently interested in finding alternative ontologies where observations are not taken to be knowledge. One such ontology, which is motivated by ideas from stochastic processes, is described in [Boutilier, Friedman, and Halpern 1998]. This ontology explicitly models how “noisy” observations come about. Roughly speaking, we assign a prior plausibility to observing a formula  $\varphi$  in a world  $w$ . When we observe  $\varphi$  we update our plausibility measure on worlds by considering the plausibility of observing  $\varphi$  in each of these worlds. Of course, we can then consider various assumptions on this prior plausibility. For example, we might say that we believe observations to be true. That is, our prior plausibility of observing  $\varphi$  in worlds where it is true is higher than the plausibility of observing  $\varphi$  in worlds where it false. We note that this assumption does *not* imply that  $\varphi \in \text{Bel}(K * \varphi)$ , since we have to combine the plausibility of observing  $\varphi$  with the plausibility of  $\varphi$  before the observation was made. If  $\varphi$  is considered implausible in  $K$ , then the single observation of  $\varphi$  might not suffice to make  $\varphi$  more plausible than  $\neg\varphi$ . (For example, the medieval scientist probably would not change her beliefs about the speed of falling objects after a single experiment.)

In this paper we have focused on belief revision. However, the need to clarify the underlying ontology goes far beyond belief revision. Much the same comments can be made for all the work on nonmonotonic logic as well (this point is essentially made in [Halpern 1993]). Not surprisingly, our critique applies to other approaches belief change as well, and in particular to Katsuno and Mendelzon’s *belief update* [1991a]. Although the motivation described by Katsuno and Mendelzon is different than that of revision, the discussion of update is stated in terms of postulates about belief sets. Thus, for example, the distinction between the agent’s belief set and epistemic state arises in update as well: update is defined as function from (belief sets  $\times$  formulas) to belief sets. However, the agent’s belief set does not uniquely determine the outcome of update. This is demonstrated, for example, by Katsuno and Mendelzon’s semantic characterization that requires the agent to have a ternary relation on possible worlds such that  $w_1 <_w w_2$

holds if the agent considers  $w_1$  to be “closer” to  $w$  than  $w_2$ . Other issues we raise here also apply to update for similar reasons.

The ontology we propose in Section 4, where observations are treated as knowledge, can be applied to update as well. Moreover, the assumption that observations are true is less problematic for iterated update: since update does not assume that propositions are static, we can consider runs where  $\varphi$  is true at one time point, and false at the next one. Thus, assuming that observations are known to be true does not rule out sequences of observations of the form  $\varphi, \neg\varphi, \dots$ . In fact, all sequences of consistent observations are allowed. We refer the reader to [Friedman 1997; Friedman and Halpern 1998] for more details.

It seems to us that many of the intuitions that researchers in the area have are motivated by thinking in terms of observations as known, even if this is not always reflected in the postulates considered. We have examined carefully one particular instantiation of this ontology, that of treating observations as knowledge. We have shown that, in this ontology, some postulates that seem reasonable, such as Lehmann’s I4, do not hold. We do not mean to suggest that I4 is “wrong” (whatever that might mean in this context). Rather, it shows that we cannot blithely accept postulates without making the underlying ontology clear. We would encourage the investigation of other ontologies for belief change.

### Acknowledgments

The authors are grateful to Craig Boutilier, Adnan Darwiche, Adam Grove, Daniel Lehmann, and an anonymous reviewer for comments on the paper and useful discussions relating to this work.

## References

- Alchourrón, C. E., P. Gärdenfors, and D. Makinson (1985). On the logic of theory change: partial meet functions for contraction and revision. *Journal of Symbolic Logic* 50, 510–530.
- Boutilier, C. (1992). Normative, subjective and autoepistemic defaults: adopting the Ramsey test. In *Principles of Knowledge Representation and Reasoning: Proc. Third International Conference (KR ’92)*, pp. 685–696. San Francisco, Calif.: Morgan Kaufmann.
- Boutilier, C. (1994). Unifying default reasoning and belief revision in a modal framework. *Artificial Intelligence* 68, 33–85.
- Boutilier, C. (1996). Iterated revision and minimal change of conditional beliefs. *Journal of Philosophical Logic* 25, 262–305.

- Boutilier, C., N. Friedman, and J. Y. Halpern (1998). Belief revision with unreliable observations. In *Proceedings, Fifteenth National Conference on Artificial Intelligence (AAAI '96)*, pp. 127–134.
- Boutilier, C. and M. Goldszmidt (1993). Revising by conditional beliefs. In *Proceedings, Eleventh National Conference on Artificial Intelligence (AAAI '93)*, pp. 648–654.
- Darwiche, A. and J. Pearl (1994). On the logic of iterated belief revision. In *Principles of Knowledge Representation and Reasoning: Proc. Fourth International Conference (KR '94)*, pp. 5–23. San Francisco, Calif.: Morgan Kaufmann.
- Darwiche, A. and J. Pearl (1997). On the logic of iterated belief revision. *Artificial Intelligence* 89, 1–29.
- Dubois, D. and H. Prade (1992). Belief change and possibility theory. In P. Gärdenfors (Ed.), *Belief Revision*. Cambridge, U.K.: Cambridge University Press.
- Fagin, R., J. Y. Halpern, Y. Moses, and M. Y. Vardi (1995). *Reasoning about Knowledge*. Cambridge, Mass.: MIT Press.
- Freund, M. and D. Lehmann (1994). Belief revision and rational inference. Technical Report TR 94-16, Hebrew University.
- Friedman, N. (1997). *Modeling Beliefs in Dynamic Systems*. Ph. D. thesis, Stanford.
- Friedman, N. and J. Y. Halpern (1994). A knowledge-based framework for belief change. Part II: revision and update. In J. Doyle, E. Sandewall, and P. Torasso (Eds.), *Principles of Knowledge Representation and Reasoning: Proc. Fourth International Conference (KR '94)*, pp. 190–201. San Francisco, Calif.: Morgan Kaufmann.
- Friedman, N. and J. Y. Halpern (1997). Modeling belief in dynamic systems. part I: foundations. *Artificial Intelligence* 95(2), 257–316.
- Friedman, N. and J. Y. Halpern (To appear, 1998). Modeling belief in dynamic systems. Part II: revision and update. *Journal of A.I. Research*.
- Gärdenfors, P. (1988). *Knowledge in Flux*. Cambridge, Mass.: MIT Press.
- Gärdenfors, P. and D. Makinson (1988). Revisions of knowledge systems using epistemic entrenchment. In *Proc. Second Conference on Theoretical Aspects of Reasoning about Knowledge*, pp. 83–95. San Francisco, Calif.: Morgan Kaufmann.
- Goldszmidt, M. (1992). *Qualitative probabilities: a normative framework for common-sense reasoning*. Ph. D. thesis, University of California Los Angeles.
- Goldszmidt, M. and J. Pearl (1992). Rank-based systems: A simple approach to belief revision, belief update and reasoning about evidence and actions. In *Principles of Knowledge Representation and Reasoning: Proc. Third International Conference (KR '92)*, pp. 661–672. San Francisco, Calif.: Morgan Kaufmann.
- Grove, A. (1988). Two modelings for theory change. *Journal of Philosophical Logic* 17, 157–170.

- Halpern, J. Y. (1993). A critical reexamination of default logic, autoepistemic logic, and only knowing. In *Proceedings, 3rd Kurt Gödel Colloquium*, pp. 43–60. Springer-Verlag.
- Halpern, J. Y. and R. Fagin (1989). Modelling knowledge and action in distributed systems. *Distributed Computing* 3(4), 159–179. A preliminary version appeared in *Proc. 4th ACM Symposium on Principles of Distributed Computing*, 1985, with the title “A formal model of knowledge, action, and communication in distributed systems: preliminary report”.
- Hansson, S. O. (1991). *Belief Base Dynamics*. Ph. D. thesis, Uppsala University.
- Hansson, S. O. (1998a). Belief revision from an epistemological point of view. Unpublished manuscript.
- Hansson, S. O. (1998b). A survey of non-prioritized belief revision. Unpublished manuscript.
- Katsuno, H. and A. Mendelzon (1991a). On the difference between updating a knowledge base and revising it. In *Principles of Knowledge Representation and Reasoning: Proc. Second International Conference (KR '91)*, pp. 387–394. San Francisco, Calif.: Morgan Kaufmann.
- Katsuno, H. and A. Mendelzon (1991b). Propositional knowledge base revision and minimal change. *Artificial Intelligence* 52(3), 263–294.
- Lehmann, D. (1995). Belief revision, revised. In *Proc. Fourteenth International Joint Conference on Artificial Intelligence (IJCAI '95)*, pp. 1534–1540.
- Levi, I. (1988). Iteration of conditionals and the Ramsey test. *Synthese* 76, 49–81.
- Makinson, D. (1997). Screened revision.
- Nayak, A. C. (1994). Iterated belief change based on epistemic entrenchment. *Erkenntnis* 41, 353–390.
- Rott, H. (1989). Conditionals and theory change: revision, expansions, and additions. *Synthese* 81, 91–113.
- Rott, H. (1991). Two methods of constructing contractions and revisions of knowledge systems. *Journal of Philosophical Logic* 20, 149–173.
- Shafer, G. (1976). *A Mathematical Theory of Evidence*. Princeton, N.J.: Princeton University Press.
- Spohn, W. (1988). Ordinal conditional functions: a dynamic theory of epistemic states. In W. Harper and B. Skyrms (Eds.), *Causation in Decision, Belief Change, and Statistics*, Volume 2, pp. 105–134. Dordrecht, Netherlands: Reidel.
- Williams, M. (1994). Transmutations of knowledge systems. In *Principles of Knowledge Representation and Reasoning: Proc. Fourth International Conference (KR '94)*, pp. 619–629. San Francisco, Calif.: Morgan Kaufmann.

Wobcke, W. (1995). Belief revision, conditional logic, and nonmonotonic reasoning.  
*Notre Dame Journal of Formal Logic* 36(1), 55–102.