

# Editing Boolean Classifiers: A Belief Change Perspective

Nicolas Schwind<sup>1</sup>, Katsumi Inoue<sup>2,3</sup>, Pierre Marquis<sup>4,5</sup>

<sup>1</sup> National Institute of Advanced Industrial Science and Technology, Tokyo, Japan

<sup>2</sup> National Institute of Informatics, Tokyo, Japan

<sup>3</sup> The Graduate University for Advanced Studies, SOKENDAI, Tokyo, Japan

<sup>4</sup> Univ. Artois, CNRS, CRIL, F-62300 Lens, France

<sup>5</sup> Institut Universitaire de France

nicolas-schwind@aist.go.jp, inoue@nii.ac.jp, marquis@cril.fr

## Abstract

This paper is about *editing* Boolean classifiers, i.e., determining how a Boolean classifier should be modified when new pieces of evidence must be incorporated. Our main goal is to delineate what are the rational ways of making such edits. This goes through a number of rationality postulates inspired from those considered so far for belief revision. We give a representation theorem and present some families of edit operators satisfying the postulates.

## Introduction

Alice, a bank employee, receives Bob, a customer who wants to obtain a loan. Bob has a low income, but no debts. His record shows that he had already requested a loan in the past, and had fully reimbursed it. The bank management has recently provided Alice with an AI algorithm (a pre-trained predictor) to help her decide which issue to give to any loan application. Alice is asked to use this algorithm which recommends against granting Bob the requested loan due to the fact that he is not the owner of his principal residence. However, Alice is experienced and remembers of two customers Cindy and Dan with a profile similar to Bob's, who both had previously been granted a loan without any issue. Hence, Alice's expertise led her not to follow the recommendation of the AI algorithm and to grant Bob the loan requested. But Alice would like to do more to avoid that the problem encountered arises again with future clients having similar profiles. She wonders what could be done to this end.

The research question tackled in this paper is relevant to Alice's concern. We focus on Boolean classifiers  $\varphi$ : given an instance represented as a world, i.e., a truth assignment of all the variables of interest,  $\varphi$  classifies the instance as positive when it is a model of  $\varphi$ , and as negative when it is a counter-model of  $\varphi$ . The concept associated with  $\varphi$  is the set of all positive instances. Our very purpose is to determine how a Boolean classifier  $\varphi$  that has already been learned should be modified when new pieces of positive evidence / negative evidence  $\mu$  (that may conflict with predictions of the classifier) are considered. We call such change operations on Boolean classifiers *positive edit* / *negative edits* (respectively), and we note them  $\diamond^+$  and  $\diamond^-$ .

Copyright © 2023, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

We assume that both inputs  $\varphi$  and  $\mu$  are represented as propositional formulae. Doing so,  $\varphi$  identifies both the classifier under consideration and, through its set of models, the associated concept. On the other hand,  $\mu$ 's models represent new positive (resp. negative) pieces of evidence in the case of a positive (resp. negative) edit. This representation choice allows one to deal with a number of existing ML classifiers: in an eXplainable AI perspective, many works have shown recently how ML classifiers  $C$  of various types can be associated with Boolean circuits  $\varphi_C$ , exhibiting the same input-output behaviours (i.e., the predictions made using  $C$  are precisely the same ones as those made using  $\varphi_C$ ). The ML models that are concerned include not only decision trees (Izza, Ignatiev, and Marques-Silva 2020; Audemard et al. 2021) and decision lists (Ignatiev and Silva 2021), but also a number of ML models that are usually considered as less interpretable, like random forests (Audemard, Koriche, and Marquis 2020; Izza and Marques-Silva 2021), gradient boosted trees (Ignatiev 2020), some Bayes nets (Shih, Choi, and Darwiche 2018, 2019), and binary neural networks (Narodytska et al. 2018; Shi et al. 2020). Accordingly, classifiers  $C$  from those families can be taken into account in our framework, using  $\varphi_C$  as a representation of  $C$  since the two are prediction-equivalent.

Edit operations are connected to *incremental concept learners*, like Mitchell's Candidate Elimination Algorithm (Mitchell 1977), Schlimmer and Granger's STAGGER (Schlimmer and Fisher 1986), Fisher's COBWEB (Fisher 1987), and Gallant's Pocket Algorithm (Gallant 1988). Such systems, also referred to as *on-line learning systems*, are suited to learning scenarios when a whole training set is not available a priori but examples arrive over time. Borrowing the criteria used in (Malooof and Michalski 2000) to draw a typology of such systems, edit operations characterize on-line learning systems with (full) concept memory (the role played by the classifier  $\varphi$ ), temporal batch (the set of models of  $\mu$  is a new set of examples in the case of a positive edit, and a new set of counter-examples in the case of a negative edit), and no instance memory (the examples and counter-examples used to induce  $\varphi$  are not stored). However, previous work about incremental concept learners was typically centered on aspects that are not considered in this paper. These included the design of a number of on-line learners (based on specific concept representations, e.g., decision

trees or decision rules), the evaluation of their empirical accuracy but also of their run-time efficiency (this can be a critical aspect since items in a data stream can be received at a so high rate that real-time guarantees are required to handle all of them (Domingos and Hulten 2000)), and finally the choice of examples that must be kept at each learning step (Maloof and Michalski 2004).

Contrastingly, in our work, the focus is on *an axiomatic approach*. We do not consider any specific concept representation, and do not make any assumption about how the batch of new examples or counter-examples are represented. We nevertheless suppose that the new piece of evidence  $\mu$  that triggers the edit operation of  $\varphi$  is certain, i.e., not pervaded by any noise. Thus, stepping back to the loan scenario, Alice is sure that Bob should be granted the loan. Here, our main goal is to delineate the rational ways of making such edits. This goes through a number of rationality postulates.

To determine such postulates, we look back at the *core principles of belief revision* which aims to incorporate, in a rational way, a new piece of information into the belief set of an agent (Alchourrón, Gärdenfors, and Makinson 1985; Alchourrón and Makinson 1985; Gärdenfors 1988). The AGM postulates (for Alchourrón, Gärdenfors and Makinson 1985) aim to formalize a set of rationality conditions based on three main principles: *primacy of update* (the new information must be believed after the change), *consistency* (the resulting belief set must be kept consistent when the new information is consistent), and *minimal change* (if simply adding the new information to the belief set raises no conflict, then nothing else should be added or removed).

Adapting the postulates of belief revision into edit is not a trivial task: provided that the beliefs of an agent  $\varphi$  and the new information  $\mu$  are represented by propositional formulae, when the conjunction of  $\varphi$  and  $\mu$  is consistent, the revision of  $\varphi$  by  $\mu$  corresponds to that conjunction (Katsuno and Mendelzon 1991). However, by representing a Boolean classifier  $\varphi$  and a set of incoming examples  $\mu$  by two propositional formulae, one cannot reasonably require the edited classifier to be the conjunction of  $\varphi$  and  $\mu$  whenever consistent: this process would unconditionally remove positive instances not explicitly questioned by  $\mu$ , while also not incorporating the examples from  $\mu$  previously classified negatively by  $\varphi$ .

Edit differs from belief revision in that the objects under consideration (all of which being represented by propositional formulae) are nevertheless of different nature. Thus, an agent's beliefs (represented by  $\varphi$ ) correspond to a set of possible worlds to whom the one actual "true" world is believed to belong, while in edit, it makes perfect sense for several instances to be both members of the concept represented by a Boolean classifier. Likewise, every Boolean classifier  $\varphi$  is essentially "consistent": when  $\varphi$  has no model, it simply represents the empty concept. This explains also why the consistency principle is irrelevant to an edit operation.

Nevertheless, the primacy of update and minimality of change principles can be adapted to the edit context. For this purpose, after some formal preliminaries, we introduce the edit postulates in the context of *positive* edit first (incorporating a batch of positive instances into a Boolean classifier).

We also give a representation theorem and present some examples of positive edit operators. Then, we show how these postulates can be adapted to the case of a *negative* edit (i.e., when the arriving batch is interpreted as a set of negative instances) and make precise how a correspondence between the two operations can be formalized through a duality result. We then consider the case of a *full edit*, where both positive and negative instances can be considered in the same batch. Lastly, related work is discussed just before the conclusion. The proofs of propositions are available online.<sup>1</sup>

## Formal Preliminaries

We consider a propositional language  $\mathcal{L}_{PS}$  built from a finite set  $PS$  of variables and the standard connectives. A *world* is a truth assignment of all variables from  $PS$ . The set of all worlds is denoted by  $\Omega$ , and the set of models of a propositional formula  $\varphi \in \mathcal{L}_{PS}$  (i.e., the set of worlds that make  $\varphi$  true) is denoted by  $[\varphi]$ . Given two formulae  $\alpha, \beta$ , we write  $\alpha \models \beta$  whenever  $[\alpha] \subseteq [\beta]$  and  $\alpha \equiv \beta$  when  $[\alpha] = [\beta]$ .

Belief revision aims to incorporate into the beliefs of an agent (a formula  $\varphi$ ) a new piece of information (a formula  $\mu$ ). Thus a revision operator  $\circ$  associates formulae  $\varphi, \mu$  with a revised formula  $\varphi \circ \mu$ , and is expected to satisfy a set of rationality postulates:<sup>2</sup>

**Definition 1** (KM revision operator). *A revision operator  $\circ$  is said to be a KM revision operator if it satisfies the following postulates:*

- (R1)  $\varphi \circ \mu \models \mu$
- (R2) If  $[\varphi \wedge \mu] \neq \emptyset$ , then  $\varphi \circ \mu \equiv \varphi \wedge \mu$
- (R3) If  $[\mu] \neq \emptyset$ , then  $[\varphi \circ \mu] \neq \emptyset$
- (R4) If  $\varphi \equiv \varphi'$  and  $\mu \equiv \mu'$ , then  $\varphi \circ \mu \equiv \varphi' \circ \mu'$
- (R5)  $(\varphi \circ \mu) \wedge \mu' \models \varphi \circ (\mu \wedge \mu')$
- (R6) If  $[(\varphi \circ \mu) \wedge \mu'] \neq \emptyset$ , then  $\varphi \circ (\mu \wedge \mu') \models (\varphi \circ \mu) \wedge \mu'$

(R1) is the success postulate, it relates to the primacy of update principle: the new information must be believed after revision. (R3) is the consistency postulate. (R4) is the syntax-irrelevance postulate. And (R2), (R5) and (R6) express the minimality of change conditions. We refer the reader to (Alchourrón, Gärdenfors, and Makinson 1985; Katsuno and Mendelzon 1991) for a deeper discussion about the rationale of these postulates.

## Positive Edit

We now intend to define a change operation  $\diamond^+$  that consists in editing an (already learned) Boolean classifier  $\varphi$  according to a new information  $\mu$ . We assume that  $\varphi$  is represented by a propositional formula. In this context, each world represents an *instance*, and a world  $\omega$  is a model of  $\varphi$  if and only if it is a positive instance of the concept represented by  $\varphi$  (so each instance is either classified as positive or negative by  $\varphi$ ). The new information  $\mu$  is called a *positive dataset* and is

<sup>1</sup><https://nicolas-schwind.github.io/SIM-AAAI23-proofs.pdf>

<sup>2</sup>We give here the KM postulates (Katsuno and Mendelzon 1991), which are the translation of the AGM postulates in finite propositional logic.

also represented by a propositional formula. The set of models of  $\mu$  represents a batch of arriving positive instances, also called positive *examples* (i.e., when referring to the models of  $\mu$ ). We do not make any further assumption on the way  $\varphi$  and  $\mu$  are represented (e.g.,  $\varphi$  could be a decision tree and  $\mu$  a DNF formula, but it does not have to be the case).

**Example 1.** Let us formalize the scenario provided in the introduction. We set  $PS = \{p, q, r, s\}$  where  $p$  means that the applicant “has a high income”,  $q$  stands for “owns her principal residence”,  $r$  means “has no debts”, and  $s$  means “has reimbursed a previous loan”. We assume that  $\varphi = p \wedge q \wedge r$ , i.e., the predictor recommends granting a loan precisely to those residence owners having a high income and no debts. Then, let Bob have the profile  $\omega_1 = 0011$ , i.e., he is not owning his residence, has a low income, but has no debts and has reimbursed a previous loan; and let Cindy and Dan be identified with the same profile  $\omega_2 = 0101$ . The positive dataset  $\mu$  is then defined as any propositional formula such that  $[\mu] = \{\omega_1, \omega_2\}$ , e.g.,  $\mu = \neg p \wedge s \wedge (q \leftrightarrow \neg r)$ .

An edit operator  $\diamond^+$  associates every Boolean classifier  $\varphi$  and every positive dataset  $\mu$  with an edited Boolean classifier  $\varphi \diamond^+ \mu$ . Our key assumption is that the new piece of evidence  $\mu$  that triggers the edit operation of  $\varphi$  is provided by a domain expert: it is therefore certain, i.e., not pervaded by any noise. This can be ensured in a number of scenarios (thus, stepping back to the example given in the introduction, Alice is sure that applicants with the same profiles as Bob, Cindy and Dan should be granted the loan).

We are ready to introduce our postulates for positive edit:

**Definition 2** (Positive Edit operator). *An operator  $\diamond^+$  is said to be a positive edit operator (PE operator for short) if it satisfies the following postulates:*

- (P1)  $\mu \models \varphi \diamond^+ \mu$
- (P2) If  $\mu \models \varphi$ , then  $\varphi \diamond^+ \mu \equiv \varphi$
- (P3) If  $\varphi_1 \equiv \varphi_2$  and  $\mu_1 \equiv \mu_2$ , then  $\varphi_1 \diamond^+ \mu_1 \equiv \varphi_2 \diamond^+ \mu_2$
- (P4) If  $\psi \models \varphi \diamond^+ \mu$ , then  $\varphi \diamond^+ \mu \equiv \varphi \diamond^+ (\mu \vee \psi)$

(P1) relates to the primacy of update principle. Since the incoming positive dataset  $\mu$  is assumed to be certain, (P1) requires the edited classifier to “comply” with  $\mu$ , i.e., to correctly classify all examples from  $\mu$  as positive instances. This can be viewed as the counterpart of (R1) in belief revision. (P2) is a minimality of change postulate: if the initial classifier already complies with  $\mu$ , then there is no need to change it. It is reminiscent to (R2) in belief revision, but (P2) and (R2) differ in their premise. Indeed, (P2) does not say anything when  $\mu \not\models \varphi$  and  $[\varphi \wedge \mu] \neq \emptyset$ : if  $\varphi$  does not comply with  $\mu$  (i.e., some positive examples from  $\mu$  were previously classified as negative instances by  $\varphi$ ), then it makes perfect sense to question the concept membership of any instance  $\omega \notin [\mu]$ . Note that when  $\mu \models \varphi$ , the conclusion of (P2) can equivalently be stated as  $\varphi \diamond^+ \mu \equiv \varphi \vee \mu$ , from which the similarity with (R2) is clearer:  $\mu$  is simply “added” to  $\varphi$ , which results in not changing  $\varphi$  at all. (P3) is the syntax-independence postulate, which is the direct counterpart of (R4). (P4) is another minimality of change postulate. Its counterparts in belief revision are (R5) and (R6), which together express that if  $\varphi$  revised by a first piece of information  $\mu_1$  is consistent with another piece of information  $\mu_2$ ,

then revising  $\varphi$  by both pieces of information taken together (i.e., by  $\mu_1 \wedge \mu_2$ ) boils down to “adding”  $\mu_2$  to the revision of  $\varphi$  by  $\mu_1$ . Likewise, in our setting, (P4) says that if the edit of a classifier  $\varphi$  by a first positive dataset  $\mu$  complies with another positive dataset  $\psi$ , then its edit by the two batches taken together (i.e., by  $\mu \vee \psi$ ) boils down to “adding”  $\psi$  to the edit of  $\varphi$  by  $\mu$ : indeed, if  $\psi \models \varphi \diamond^+ \mu$ , then  $\varphi \diamond^+ \mu \equiv (\varphi \diamond^+ \mu) \vee \psi$  and thus the conclusion of (P4) can equivalently be written as  $(\varphi \diamond^+ \mu) \vee \psi \equiv \varphi \diamond^+ (\mu \vee \psi)$ . Please note that (R3), the consistency postulate in belief revision, is the only postulate with no counterpart in our setting, since every Boolean classifier is essentially “consistent”: if  $[\varphi] = \emptyset$ , then  $\varphi$  characterizes an empty concept.

Notably, the edit postulates (P1-P4) are also reminiscent to properties that can be sought for incremental concept learners in the absence of noise. Thus, (P1) states that once the edit operation has been performed, the resulting concept  $\varphi \diamond^+ \mu$  must be consistent (in the sense of (Mitchell 1982)) with the new examples given by  $\mu$ , which precisely means that those examples must be positive instances of  $\varphi \diamond^+ \mu$ . (P2) requires not to change the concept  $\varphi$  when it classifies correctly the new examples given by  $\mu$ . This condition is achieved, for instance, by the perceptron update rule (Rosenblatt 1958). (P3) requires syntax not to play any role in the on-line learning process, which makes sense if the specific representation  $\varphi$  of the concept at hand is irrelevant (this is one of our starting assumptions). Finally, provided that (P2) holds, (P4) can be viewed as a relaxation of an order-independence condition that is satisfied by some on-line learners. This last property roughly states that while the new examples arrives over time, once the whole input sequence has been processed, the classifier has been transformed in the same way as if all pieces of evidence were available as a whole. Such an order-independence condition is ensured by ID5 (Utgoff 1989), that has been shown to compute the same decision tree as the one that would be generated by ID3 (Quinlan 1986), provided that the whole set of examples would be available at start. Formally, in our setting, the order-independence condition can be stated as  $(\varphi \diamond^+ \mu) \diamond^+ \psi \equiv \varphi \diamond^+ (\mu \vee \psi)$ . Note that this condition is quite demanding and not satisfied by every on-line learner (e.g., the perceptron update rule may easily question the way an instance  $\mu$  has been classified when editing further the linear classifier by taking a new instance  $\psi$  into account). Accordingly, we focused on a weaker condition (it is easy to show that (P4) is a logical consequence of the order-independence condition when (P2) holds).

At this point, one can already identify a few simple operators from the class:

**Definition 3** (Some PE operators). *The trivial, basic, and drastic operators, respectively noted  $\diamond_T^+$ ,  $\diamond_B^+$ , and  $\diamond_D^+$ , are defined for each classifier  $\varphi$  and each positive dataset  $\mu$  as:*

- $\varphi \diamond_T^+ \mu = \varphi$  if  $\mu \models \varphi$ , otherwise  $\varphi \diamond_T^+ \mu = \top$
- $\varphi \diamond_B^+ \mu = \varphi \vee \mu$
- $\varphi \diamond_D^+ \mu = \varphi$  if  $\mu \models \varphi$ , otherwise  $\varphi \diamond_D^+ \mu = \mu$

The trivial operator  $\diamond_T^+$  just requires a classifier to classify all worlds as positive instances as soon as it does not

initially comply with the input positive dataset. The basic operator  $\diamond_B^+$  is simply defined as disjunction: it adds as positive instances all the examples provided by  $\mu$  which were not already classified as positive. The drastic operator  $\diamond_D^+$  leaves the classifier  $\varphi$  unchanged if already compliant with  $\mu$ . Otherwise, similarly to the trivial operator, it “forgets” everything, but classifies as positive instances precisely the ones explicitly provided as examples by  $\mu$ . This can be viewed as the PE counterpart of the drastic revision operator  $\circ_D$  defined as  $\varphi \circ_D \mu = \varphi \wedge \mu$  if  $[\varphi \wedge \mu] \neq \emptyset$ , otherwise  $\varphi \circ_D \mu = \mu$ .

It is quite easy to check that these operators satisfy **(P1-P4)** (the proof is direct):

**Proposition 1.**  $\diamond_T^+$ ,  $\diamond_B^+$  and  $\diamond_D^+$  are PE operators.

**A Representation Theorem.** Let us now show how PE operators can be characterized in terms of so-called *positive assignments*:

**Definition 4** (Positive assignment). A positive assignment is a mapping associating every classifier  $\varphi$  with a mapping  $f_\varphi : \mathcal{P}(\Omega) \mapsto \mathcal{P}(\Omega)$ , such that for all classifiers  $\varphi, \varphi'$  and all subsets of worlds  $W, W' \in \mathcal{P}(\Omega)$ , the following properties are satisfied:

1.  $W \subseteq f_\varphi(W)$
2. If  $W \subseteq [\varphi]$ , then  $f_\varphi(W) = [\varphi]$
3. If  $\varphi \equiv \varphi'$ , then  $f_\varphi = f_{\varphi'}$
4. If  $W \subseteq W'$  and  $W' \subseteq f_{\varphi'}(W)$ , then  $f_\varphi(W) = f_{\varphi'}(W)$

**Proposition 2.** An operator  $\diamond^+$  is a PE operator if and only if there is a positive assignment  $\varphi \mapsto f_\varphi$  such that for each classifier  $\varphi$  and each positive dataset  $\mu$ ,  $[\varphi \diamond \mu] = f_\varphi([\mu])$ .

This is a “strong” representation result, in the sense that different positive assignments define different PE operators.

Interestingly, every mapping  $f_\varphi$  satisfies the property of *idempotence*:

**Proposition 3.** For each positive assignment  $\varphi \mapsto f_\varphi$  and each  $W \subseteq \Omega$ , we have that  $f_\varphi(f_\varphi(W)) = f_\varphi(W)$ .

Accordingly, a consequence of the PE postulates is that  $(\varphi \diamond^+ \mu) \diamond^+ \mu \equiv \varphi \diamond^+ \mu$ . This idempotence property reflects a very simple form of minimal change and is standard in belief change: it is satisfied by belief revision operators but also by other forms of change operations, e.g., contraction (Caridroit, Konieczny, and Marquis 2017).

Noteworthy, condition 4 corresponds to the condition of Irrelevance of Rejected Contracts (IRC) in matching theory (Hatfield and Milgrom 2005). In that context, this property requires the removal of rejected contracts not to affect a choice set, and is a necessary condition to guarantee the existence of stable allocations (Aygün and Sönmez 2013), without implying rationalizability<sup>3</sup>(Yang 2020).

**Distance-based PE operators.** We now introduce two classes of operators, called *dilation* operators and *min-generalization* operators. These operators are parameterized by a distance between worlds, i.e., a mapping  $d : \Omega \times \Omega \mapsto \mathbb{N}$

<sup>3</sup>A choice function  $\sigma : \mathcal{P}(E) \mapsto \mathcal{P}(E)$ , i.e., a mapping such that  $\sigma(S) \subseteq S$ , is *rationalizable* if it can be characterized in terms of preference relation over elements from  $E$ .

such that  $d(\omega, \omega') = 0$  if and only if  $\omega = \omega'$ , and that satisfies the triangular inequality property, i.e.,  $d(\omega, \omega'') \leq d(\omega, \omega') + d(\omega', \omega'')$ , for all worlds  $\omega, \omega', \omega''$ .

Let us start with dilation operators, whose definition is inspired from the notion of formula dilation from (Bloch and Lang 2002; Dalal 1988). Given a classifier  $\varphi$  such that  $[\varphi] \neq \emptyset$ , and an integer  $k$ , the *k-dilation* of  $\varphi$  w.r.t.  $d$ , denoted by  $D_\varphi^d(k)$ , is defined by  $D_\varphi^d(k) = \{\omega \in \Omega \mid d(\omega, \varphi) \leq k\}$ , where  $d(\omega, \varphi) = \min\{d(\omega, \omega') \mid \omega' \models \varphi\}$ .

**Definition 5** (Dilation operator). The dilation operator  $\diamond_{\text{dil},d}^+$  induced by  $d$  is defined for each classifier  $\varphi$  and each positive dataset  $\mu$  by  $[\varphi \diamond_{\text{dil},d}^+ \mu] = [\mu]$  if  $[\varphi] = \emptyset$ , otherwise  $[\varphi \diamond_{\text{dil},d}^+ \mu] = \arg \min_k (\{D_\varphi^d(k) \mid [\mu] \subseteq D_\varphi^d(k)\})$ .

A number of dilation operators can be defined depending on the choice of  $d$ . For instance, consider the Hamming distance between worlds, denoted by  $d_H$ , defined for all worlds  $\omega, \omega' \in \Omega$  as  $d_H(\omega, \omega') = \{x \in PS \mid \omega(x) \neq \omega'(x)\}$  (Dalal 1988). Then the Hamming-based dilation operator  $\diamond_{\text{dil},d_H}^+$  consists in  $k$ -dilating  $\varphi$  w.r.t.  $d_H$  where  $k$  is the least integer for which the resulting set of models includes every model of  $\mu$  (see Example 1 below).

Let us now introduce the class of min-generalization operators. Given a distance  $d$  and a world  $\omega$ , let  $\leq_\omega^d$  be the total preorder over worlds induced by  $\omega$  and  $d$  and defined by  $\omega' \leq_\omega^d \omega''$  iff  $d(\omega', \omega) \leq_\omega^d d(\omega'', \omega)$ . Given a classifier  $\varphi$  such that  $[\varphi] \neq \emptyset$ , the set  $\min([\varphi], \leq_\omega^d)$  denotes the set of models of  $\varphi$  that have a distance to  $\omega$  which is minimal among all models of  $\varphi$ , i.e.,  $\min([\varphi], \leq_\omega^d) = \{\omega' \in [\varphi] \mid \forall \omega'' \in [\varphi], d(\omega', \omega) \leq d(\omega'', \omega)\}$ .

**Definition 6** (Min-generalization operator). The min-generalization operator  $\diamond_{\text{gen},d}^+$  induced by  $d$  is defined for each classifier  $\varphi$  and each positive dataset  $\mu$  by  $[\varphi \diamond_{\text{gen},d}^+ \mu] = [\mu]$  if  $[\varphi] = \emptyset$ , otherwise  $[\varphi \diamond_{\text{gen},d}^+ \mu] = \{\omega \in \Omega \mid \exists \omega', \omega'' \in \Omega, \omega' \models \mu, \omega'' \in \min([\varphi], \leq_{\omega'}^d), d(\omega, \omega') + d(\omega, \omega'') \leq d(\omega', \omega'')\}$ .

The min-generalization operator consists in considering as positive instances for  $\varphi \diamond_{\text{gen},d}^+ \mu$  every world  $\omega$  that is “in-between” (w.r.t.  $d$ ) a model  $\omega'$  of  $\mu$  and a model  $\omega''$  of  $\varphi$  that is among the closest ones (w.r.t.  $d$ ) to  $\omega'$ . When  $d = d_H$ , the min-generalization operator can be characterized using the *most specific generalization* (msg) of the worlds involved.<sup>4</sup> Let  $\text{msg}(\omega, \omega')$  be the term  $\bigwedge_{x \in PS \mid \omega(x) = \omega'(x) = 1} x \wedge \bigwedge_{x \in PS \mid \omega(x) = \omega'(x) = 0} \neg x$ . Then one can check that  $\varphi \diamond_{\text{gen},d_H}^+ \mu \equiv \bigvee \{\text{msg}(\omega, \omega') \mid \omega \in [\mu], \omega' \in \min([\varphi], \leq_\omega^d)\}$ .

All the operators from these classes satisfy **(P1-P4)**:

**Proposition 4.** For every distance  $d$ , the operators  $\diamond_{\text{dil},d}^+$  and  $\diamond_{\text{gen},d}^+$  are PE operators.

**Example 1** (continued). Let us go back to our loan scenario, and recall that  $PS = \{p, q, r, s\}$ ,  $\varphi = p \wedge q \wedge r$ , and  $\mu = \neg p \wedge s \wedge (q \leftrightarrow \neg r)$ . Figure 1 depicts through

<sup>4</sup>Most specific generalization is a key concept in machine learning, see e.g., (Plotkin 1970; Mitchell 1977).

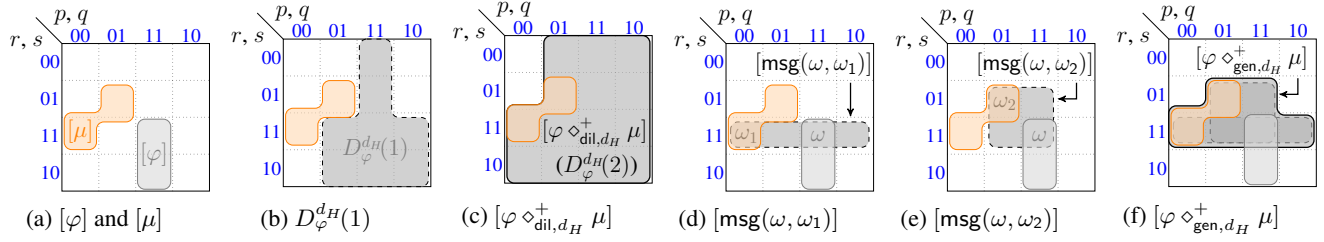


Figure 1: An example of Hamming-based dilation (Fig. 1c) and min-generalization (Fig. 1f) positive edits.

Karnaugh maps the models of  $\varphi$  and  $\mu$  (Fig. 1a), the 1-dilation of  $\varphi$  (Fig. 1b), the Hamming-based dilation edit of  $\varphi$  by  $\mu$ , which corresponds to the 2-dilation of  $\varphi$  (Fig. 1c), and the Hamming-based min-generalization edit of  $\varphi$  by  $\mu$  (Fig. 1f), which corresponds to the disjunction of the two msgs given in Fig. 1d and 1e. Accordingly, we get that  $\varphi \diamond_{\text{dil},d_H}^+ \mu \equiv p \vee q \vee r$  and  $\varphi \diamond_{\text{gen},d_H}^+ \mu \equiv s \wedge (q \vee r)$ .

As it can be verified on the example, dilation and min-generalization PE operators can easily add to  $\varphi$  models that are neither models of  $\varphi$  nor models of  $\mu$ , thus questioning negative instances (the counter-models of  $\varphi$ ). Such a generalization power (required by incremental learning) is not forbidden but not mandatory for PE operators (e.g., consider the basic and the drastic PE operators in Definition 3). Min-generalization PE operators may also question positive instances (again, this is expected when used for learning).

### Negative Edit

Let us now consider the incorporation of negative instances into a Boolean classifier, alias negative edits. This time, the models of the change formula  $\mu$  represent counter-examples of the target concept (a *negative dataset*).

**Definition 7** (Negative Edit operator). *An operator  $\diamond^-$  is said to be a negative edit operator (NE operator for short) if it satisfies the following postulates:*

- (N1)  $\mu \models \neg(\varphi \diamond^- \mu)$
- (N2) If  $\mu \models \neg\varphi$ , then  $\varphi \diamond^- \mu \equiv \varphi$
- (N3) If  $\varphi_1 \equiv \varphi_2$  and  $\mu_1 \equiv \mu_2$ , then  $\varphi_1 \diamond^- \mu_1 \equiv \varphi_2 \diamond^- \mu_2$
- (N4) If  $\psi \models \neg(\varphi \diamond^- \mu)$ , then  $\varphi \diamond^- \mu \equiv \varphi \diamond^- (\mu \vee \psi)$

Similarly to Harper and Levi's identities which show how a revision operator can be defined from a contraction operator and vice-versa (see e.g., (Caridroit, Konieczny, and Marquis 2017)), one can identify a correspondence between PE and NE operators. With an operator  $\diamond^* : \mathcal{L} \times \mathcal{L} \mapsto \mathcal{L}$ , let us associate an operator  $\sigma(\diamond^*) : \mathcal{L} \times \mathcal{L} \mapsto \mathcal{L}$  defined as:

$$\varphi \sigma(\diamond^*) \mu = \neg(\neg\varphi \diamond^* \mu),$$

for every classifier  $\varphi$  and every formula  $\mu$ . Then:

**Proposition 5.**  *$\sigma$  is an involution, that is, for each operator  $\diamond^* : \mathcal{L} \times \mathcal{L} \mapsto \mathcal{L}$ ,  $\sigma(\sigma(\diamond^*)) = \diamond^*$ . Moreover,  $\sigma(\diamond^*)$  is an NE operator if and only if  $\diamond^*$  is a PE operator.*

We say that the operator  $\sigma(\diamond^*)$  is the *dual* of  $\diamond^*$ . For instance, consider again the trivial, basic and drastic PE operators introduced in Definition 3. Then it is easy to see

that the dual of these operators, i.e., the trivial, basic, and drastic NE operators, respectively denoted by  $\diamond_T^- = \sigma(\diamond_T^+)$ ,  $\diamond_B^- = \sigma(\diamond_B^+)$ , and  $\diamond_D^- = \sigma(\diamond_D^+)$ , are defined for each classifier  $\varphi$  and each negative dataset  $\mu$  by:

- $\varphi \diamond_T^- \mu = \varphi$  if  $\mu \models \neg\varphi$ , otherwise  $\varphi \diamond_T^- \mu = \perp$
- $\varphi \diamond_B^- \mu = \varphi \wedge \neg\mu$
- $\varphi \diamond_D^- \mu = \varphi$  if  $\mu \models \neg\varphi$ , otherwise  $\varphi \diamond_D^- \mu = \neg\mu$

Dual operators of dilation operators and min-generalization operators can also be easily defined. Interestingly, the operators dual to dilation operators involve an operation of formula *erosion* (Bloch and Lang 2002), which is an operation on formulae dual to the one of dilation. For instance, the Hamming-based erosion operator, denoted by  $\diamond_{\text{ero},d_H}^-$ , is defined for each classifier  $\varphi$  and each negative dataset  $\mu$  as  $\varphi \diamond_{\text{ero},d_H}^- \mu = \neg(\neg\varphi \diamond_{\text{dil},d_H}^+ \mu)$ .

### Full Edit

Let us finally consider the more general case when the new piece of evidence consists of both positive and negative instances that must be incorporated into the classifier. We call such a piece of evidence a *dataset*, i.e., a pair  $(\mu^+, \mu^-)$  such that  $\mu^+$  is a positive dataset (a set of examples),  $\mu^-$  is a negative dataset (a set of counter-examples), and such that  $[\mu^+ \wedge \mu^-] = \emptyset$ . The set  $\mathcal{D}$  denotes the set of all datasets.

**Definition 8** (Full Edit operator). *An operator  $\diamond : \mathcal{L} \times \mathcal{D} \mapsto \mathcal{L}$  is said to be a full edit operator (FE operator for short) if for each a classifier  $\varphi$  and each dataset  $(\mu^+, \mu^-)$ , it satisfies the following postulates:*

- (F1)  $\mu^+ \models \varphi \diamond (\mu^+, \mu^-)$
- (F2)  $\mu^- \models \neg(\varphi \diamond (\mu^+, \mu^-))$
- (F3) If  $\mu^+ \models \varphi$  and  $\mu^- \models \neg\varphi$ , then  $\varphi \diamond (\mu^+, \mu^-) \equiv \varphi$
- (F4) If  $\varphi_1 \equiv \varphi_2$ ,  $\mu_1^+ \equiv \mu_2^+$  and  $\mu_1^- \equiv \mu_2^-$ , then  $\varphi_1 \diamond (\mu_1^+, \mu_1^-) \equiv \varphi_2 \diamond (\mu_2^+, \mu_2^-)$
- (F5) If  $\psi \models \varphi \diamond (\mu^+, \mu^-)$  and  $\alpha \models \neg(\varphi \diamond (\mu^+, \mu^-))$ , then  $\varphi \diamond (\mu^+, \mu^-) \equiv \varphi \diamond (\mu^+ \vee \psi, \mu^- \vee \alpha)$

The postulate (F1) (resp. (F2)) corresponds to (P1) (resp. (N1)), while (F3) (resp. (F4), (F5)) is a (weak) combination of (P2) and (N2) (resp. (P3) and (N3), (P4) and (N4)).

A number of FE operators can be defined by means of a PE operator or an NE operator. Given an operator  $\diamond^+ : \mathcal{L} \times \mathcal{L} \mapsto \mathcal{L}$ , let us define the operator  $\diamond_{\diamond^+} : \mathcal{L} \times \mathcal{D} \mapsto \mathcal{L}$  for each classifier  $\varphi$  and each dataset  $(\mu^+, \mu^-)$  as:

$$\varphi \diamond_{\diamond^+} (\mu^+, \mu^-) = (\varphi \diamond^+ \mu^+) \wedge \neg\mu^-.$$

We say that  $\diamond_{\circ+}$  is *positively induced* by  $\diamond^+$ . Then:

**Proposition 6.**  $\diamond_{\circ+}$  is an FE operator if and only if  $\diamond^+$  is a PE operator.

Proposition 6 gives us a constructive way to define an FE operator from a PE operator. Consider for instance the dilation operator  $\diamond_{\text{dil},d}^+$ , where  $d$  is any distance between worlds (cf. Definition 5). Then the operator  $\diamond_{\text{dil},d}^{\circ+}$  consists in first “dilating” an input classifier  $\varphi$  so as to include all positive examples from  $\mu^+$ , and then removing all instances introduced in the dilation step according to  $\mu^-$ . As a consequence of Proposition 6, this operator satisfies **(F1-F5)**.

Likewise, each NE operator also defines an FE operator. Given an operator  $\diamond^- : \mathcal{L} \times \mathcal{L} \mapsto \mathcal{L}$ ,  $\diamond$  is said to be *negatively induced* by  $\diamond^-$ , denoted by  $\diamond = \diamond_{\circ-}$ , if it is defined for each classifier  $\varphi$  and each dataset  $(\mu^+, \mu^-)$  by  $\varphi \diamond_{\circ-} (\mu^+, \mu^-) = (\varphi \diamond^- \mu^-) \vee \mu^+$ . Echoing Proposition 6, we get that:

**Proposition 7.**  $\diamond_{\circ-}$  is an FE operator if and only if  $\diamond^-$  is an NE operator.

Remark that inducing an operator  $\diamond$  by a PE operator and an NE operator, e.g., as  $\varphi \diamond (\mu^+, \mu^-) = (\varphi \diamond^+ \mu^+) \diamond^- \mu^-$ , does not always define an FE operator, even when  $\diamond^+$  and  $\diamond^-$  are dual. To give an example when this kind of construction does not work, let us consider our loan scenario again:

**Example 1** (continued). Assume now that Alice receives an additional applicant, Emir, with profile  $\omega_3 = 0100$ . Since Emir has a low income, debts, and has not yet reimbursed his previous loan, Alice is certain that Emir is not eligible for a new loan. We are then given both a dataset  $\mu = (\mu^+, \mu^-)$ , where  $\mu^+ = \neg p \wedge s \wedge (q \leftrightarrow \neg r)$  with  $[\mu^+] = \{\omega_1, \omega_2\}$  (Bob / Cindy and Dan are positive examples), and  $\mu^- = \neg p \wedge q \wedge \neg r \wedge \neg s$ , i.e.,  $[\mu^-] = \{\omega_3\}$  (Emir is a negative example). Let us consider the operator  $\diamond$  defined by  $\varphi \diamond \mu = (\varphi \diamond_{\text{dil},d_H}^+ \mu^+) \diamond_{\text{ero},d_H}^- \mu^-$ , i.e., the classifier is first edited according to  $\mu^+$  using the Hamming-based dilation edit  $\diamond_{\text{dil},d_H}^+$ , and is then edited again according to  $\mu^-$  using the Hamming-based erosion edit  $\diamond_{\text{ero},d_H}^-$ , that is, the negative edit operator dual to  $\diamond_{\text{dil},d_H}^+$ . Recall first that  $\varphi' = \varphi \diamond_{\text{dil},d_H}^+ \mu^+ = p \vee q \vee r$  (cf. Fig. 1c). Then we get that  $\varphi'' = \varphi \diamond \mu = \varphi' \diamond_{\text{ero},d_H}^- \mu^- \equiv (p \vee q) \wedge (r \vee (p \wedge q))$ , with  $[\varphi''] = D_{\varphi}^{d_H}(1)$  (cf. Fig. 1b). Yet  $\{\omega_1, \omega_2\} \cap [\varphi''] = \emptyset$ , i.e., Bob / Cindy and Dan are not classified as positive instances by the edited classifier  $\varphi''$ . Hence,  $\diamond$  does not satisfy **(F1)**, i.e.,  $\diamond$  is not an FE operator.

At that stage, a natural question is whether one can find an FE operator that is not induced by a PE operator or an NE operator. We provide below a positive answer to this question. In fact, we intend to introduce an operator which is not “decomposable” in any way by means of a combination of a PE operator and an NE operator. Formally, given an FE operator  $\diamond$ , a PE operator  $\diamond^+$  and an NE operator  $\diamond^-$ , we say that the pair  $(\diamond^+, \diamond^-)$  is *faithful to*  $\diamond$  if for each classifier  $\varphi$  and each dataset  $(\mu^+, \mu^-)$ ,  $\varphi \diamond (\mu^+, \mu^-) \equiv (\varphi \diamond^+ \mu^+) \diamond^- \mu^-$  or  $\varphi \diamond (\mu^+, \mu^-) \equiv (\varphi \diamond^- \mu^-) \diamond^+ \mu^+$ . An FE operator  $\diamond$  is then said to be *decomposable* if  $\diamond$  admits a faithful pair  $(\diamond^+, \diamond^-)$ . In particular, positively induced FE operators  $\diamond_{\circ+}$  are decomposable: it can be easily verified that for each PE

operator  $\diamond^+$ , the FE operator  $\diamond_{\circ+}$  admits the faithful pair  $(\diamond^+, \diamond_B^-)$ , where  $\diamond_B^-$  is the basic NE operator (recall that  $\varphi \diamond_B^- \mu = \varphi \wedge \neg \mu$ , for each classifier  $\varphi$  and each negative dataset  $\mu$ ). And similarly, negatively induced FE operators  $\diamond_{\circ-}$  are decomposable as well since they admit the faithful pair  $(\diamond_B^+, \diamond^-)$ .

Now, let us consider the operator  $\diamond_*$  defined for each classifier  $\varphi$  and each dataset  $(\mu^+, \mu^-)$  as:

$$\varphi \diamond_* (\mu^+, \mu^-) \begin{cases} \varphi, & \text{if } \mu^+ \models \varphi \text{ and } \mu^- \models \neg \varphi, \\ \top, & \text{if } \mu^+ \not\models \varphi \text{ and } [\mu^-] = \emptyset, \\ \mu^+, & \text{otherwise.} \end{cases}$$

This operator simply leaves unchanged the edited classifier in the case when it already complies with the input dataset (as required by **(F3)**). In the remaining cases, it behaves like the PE trivial operator  $\diamond_T^+$  if the negative batch is empty, otherwise it behaves like the PE drastic operator  $\diamond_D^+$  (cf. Definition 3). We can show that:

**Proposition 8.**  $\diamond_*$  is an FE operator that is not decomposable.

This leaves us the interesting open question of whether a characterization result for FE operators can be found. This is a perspective for further research.

## Related Work

*Theory revision* is a change operation studied by ML researchers in the nineties that is connected to the edit one. Theory revision is an important component for concept formation, and as such it has been investigated and implemented as part of knowledge acquisition and machine learning systems (see e.g., MOBAL (Morik et al. 1994) and EITHER (Ourston and Mooney 1994)). Typically, in theory revision, a theory  $\Sigma$  is a logical representation (most of the time, a FOL formula) linking together atoms, denoting features used for describing instances and targeted concepts. An instance  $x$  is classified by  $\Sigma$  as an element of a concept  $y$  (represented by an atom) whenever  $y$  can be deduced from  $\Sigma$  and  $x$ . When an instance together with its right concept (given by the change formula) is not classified by  $\Sigma$  as expected, a theory revision operator can be exploited to modify  $\Sigma$  so as to ensure that the instance is not classified incorrectly any longer in the revised theory. AGM contraction operators (Alchourrón, Gärdenfors, and Makinson 1985) can be used to this end (Wrobel 1993). Note that it can be the case that an instance  $x$  is not classified by  $\Sigma$  as an element of any concept. Accordingly, the representation to be changed  $\Sigma$  does not necessarily represent a “full” classifier as in the edit case. Furthermore, the basic operations that are used to derive the revised theory are usually not syntax-independent. This is typically the case, e.g., in the *learning from interpretations* setting (De Raedt 1997; De Raedt and Dehaspe 1997) where both the set of examples and the revised theory (called hypothesis) are full clausal theories, and other various formalizations of concept learning in logical settings, including inductive logic programming (Muggleton and De Raedt 1994; Flach 1997). This also makes them distinct from edit operations. Finally, works on theory revision are typically focused on defining specific approaches to

achieve a revision of the input theory (possibly using few basic operations (Goldsmith et al. 2004; Goldsmith and Sloan 2005)), but they do not adopt an axiomatic perspective for delineating all the rational theory revision operators.

More recently, Zhou (2019) emphasized again the importance of integrating learning and reasoning in modern learning systems. The idea is to improve the decisions made by an underlying ML system  $C$  (the classifier), taking advantage of a reasoning module  $M$ . Roughly speaking, whenever a prediction  $P$  is made by the classifier  $C$ , it is transmitted to a reasoning module that checks whether the prediction is correct or not. If it is correct, nothing should be changed; otherwise, the corrected prediction  $P'$  found by  $M$  is transmitted back to  $C$  that is trained again using  $P'$ . Our edit framework is similar in essence, where a classifier  $\varphi$  plays the role of  $C$  and the positive / negative dataset  $\mu$  is provided by an underlying expert module  $M$ . One of the strengths of Zhou’s approach is that it is model-agnostic: the ML system can be any black box function. This is reminiscent to our edit framework where no further assumption is made on the representation of the input classifier  $\varphi$  and batch  $\mu$ , besides being propositional formulae. However, in (Zhou 2019), the correction step is achieved by learning, which means that there is no guarantee that the repair is effective in the general case. In comparison, our framework, by its principled nature, guarantees the classifier to become fully compliant with the input batch after edit (cf. **(P1)**, **(N1)**, **(F1)** and **(F2)**).

Modifying a predictor so as to better take account for instances that are misclassified, as done with edit operations, is also at the core of *boosting*, a key principle in ML. In adaptative boosting for binary classification (Freund and Schapire 1997), the predictor has the form of an ensemble of weak learners, often decision trees reduced to decision stumps. The output of those weak learners is combined into a weighted sum that represents the final output of the boosted classifier. AdaBoost is an iterative learning algorithm: at each iteration, the algorithm samples the training set, taking account for the distribution given by the weights associated with the instances (at start the uniform distribution is considered), then it looks for a weak classifier which minimizes the total weighted error, uses this to calculate the error rate and the weight of the weak classifier that has been generated, and finally update the weights of the instances so as to favor at the next step the selection of instances that have been misclassified by the generated weak learner. After a preset number of iterations, the algorithm stops. It turns out that the boosted classifier generated after an iteration may still misclassify the instances that were already misclassified by the boosted classifier before the iteration. Accordingly, the update operation at work in AdaBoost for improving the current boosted tree at each iteration is not a positive edit operator: **(P1)** is not satisfied.

Lastly, closely related to our work is a recent paper about *classifier rectification* (Coste-Marquis and Marquis 2021). Unlike the present paper, more than two classes can be considered in (Coste-Marquis and Marquis 2021) (classes are explicitly represented). Thus, two subsets  $X$  and  $Y$  of  $PS$  are used to encode, on the one hand, instances (positive ones and negative ones) and on the other hand, classes. When

only two classes are targeted (the class of positive instances, a subset of  $\Omega_X$ , the worlds over  $X$ , and its complementary set in  $\Omega_X$  containing the negative instances), a singleton  $Y = \{y\}$  is enough. Coste-Marquis and Marquis [2021] point out rules to be obeyed by any rational change operation  $\star$  on Boolean classifiers  $\Sigma$ , when new pieces of evidence  $T$  must be taken into account. Boolean classifiers  $\Sigma$  are formulae from  $\mathcal{L}$  satisfying the so-called  $XY$ -classification property. When  $Y = \{y\}$ , this precisely means that  $\Sigma$  is equivalent to  $\varphi_X \Leftrightarrow y$  where  $\varphi_X$  is a formula over  $X$ . Thus,  $\Sigma$  classifies a given instance  $x \in \Omega_X$  as positive (resp. negative) whenever  $x \models \varphi_X$  (resp.  $x \models \neg\varphi_X$ ). Accordingly, every PE operation (resp. NE operation) of  $\varphi_X$  by a change formula  $\mu_X$  corresponds to a rectification operation of  $\Sigma = \varphi_X \Leftrightarrow y$  by  $T = \mu_X \Rightarrow y$  (resp.  $T = \mu_X \Rightarrow \neg y$ ). Postulates for the rectification operation have been provided in (Coste-Marquis and Marquis 2021). Though some connections between rectification postulates and PE / NE postulates exist, it is not the case that every PE (or NE) operator induces a rectification operator. Indeed, the rectification postulate **(RE2)** (see (Coste-Marquis and Marquis 2021) for details) makes formal a very demanding view of minimal change: when a change concerning an example (positive instance)  $\mu_X$  is triggered, the classifications achieved by the rectified classifier coincide with those achieved by the classifier  $\Sigma$  at start, except possibly for  $\mu_X$  (accordingly, rectification operators do not allow any generalization to take place and thus are not convenient for incremental learning).

## Conclusion

The paper was focused on the question of editing Boolean classifiers, i.e., determining how a Boolean classifier should be modified when new pieces of evidence must be incorporated, an issue at the crossroads of ML and KR. Though the performance of ML models in terms of accuracy is impressive most of the time (especially when classifiers are learned from a sufficient amount of data), the error risk cannot be totally removed (this is intrinsic to inductive generalization). Thus, it is important to design and study approaches to determine how a classifier should be modified whenever it does not label an instance in the right way. The reported work is a step in this direction, centered on the identification of first principles (postulates) for characterizing what a rational change could be when dealing with Boolean classifiers.

One of our key assumptions in this paper is that the input dataset is fully reliable, which is reflected by the success postulate **(P1)** (in the case of positive edit). However, a number of standard learning algorithms take noisy examples into account, e.g., the k-NN algorithm, the perceptron algorithm, and algorithms for generating decision trees with pruning; and those algorithms do not satisfy **(P1)**. To extend the edit setting to noisy data, we plan to investigate how **(P1)** could be relaxed, so that an example is incorporated only if the corresponding piece of evidence is considered “sufficiently often” by the learning algorithm. For capturing such a behavior, *improvement* appears as a promising candidate (Konieczny, Medina Grespan, and Pino Pérez 2010), and determining the extent to which edit and improvement could be combined looks as a valuable perspective for further work.

## Acknowledgments

This work has benefited from the support of the JSPS KAKENHI Grant Number JP21H04905 and JST CREST Grant Number JPMJCR22D3, Japan. Pierre Marquis has benefited from the support of the AI Chair EXPEKCTATION (ANR-19-CHIA-0005-01) of the French National Research Agency. He was also partially supported by TAILOR, a project funded by EU Horizon 2020 research and innovation programme under GA No952215.

## References

- Alchourrón, C. E.; Gärdenfors, P.; and Makinson, D. 1985. On the logic of theory change: partial meet contraction and revision functions. *Journal of Symbolic Logic*, 50: 510–530.
- Alchourrón, C. E.; and Makinson, D. 1985. On the logic of theory change: safe contraction. *Studia Logica*, 44(4): 405–422.
- Audemard, G.; Bellart, S.; Bounia, L.; Koriche, F.; Lagniez, J.-M.; and Marquis, P. 2021. On the Computational Intelligibility of Boolean Classifiers. In *Proceedings of the 18th International Conference on Principles of Knowledge Representation and Reasoning (KR'21)*, 74–86.
- Audemard, G.; Koriche, F.; and Marquis, P. 2020. On Tractable XAI Queries based on Compiled Representations. In *Proceedings of the 17th International Conference on Principles of Knowledge Representation and Reasoning (KR'20)*, 838–849.
- Ayguin, O.; and Sönmez, T. 2013. Matching with Contracts: Comment. *American Economic Review*, 103(5): 2050–2051.
- Bloch, I.; and Lang, J. 2002. *Towards mathematical morpho-logics*, 367–380. Heidelberg: Physica-Verlag HD.
- Caridroit, T.; Konieczny, S.; and Marquis, P. 2017. Contraction in propositional logic. *International Journal of Approximate Reasoning*, 80: 428–442.
- Coste-Marquis, S.; and Marquis, P. 2021. On Belief Change for Multi-Label Classifier Encodings. In *Proceedings of the 30th International Joint Conference on Artificial Intelligence (IJCAI'21)*, 1829–1836.
- Dalal, M. 1988. Investigations into a theory of knowledge base revision: preliminary report. In *Proceedings of the 7th National Conference on Artificial Intelligence (AAAI'88)*, 475–479.
- De Raedt, L. 1997. Logical Settings for Concept-Learning. *Artificial Intelligence*, 95(1): 187–201.
- De Raedt, L.; and Dehaspe, L. 1997. Clausal Discovery. *Machine Learning*, 26(2-3): 99–146.
- Domingos, P. M.; and Hulten, G. 2000. Mining high-speed data streams. In *Proceedings of the 6th International Conference on Knowledge Discovery and Data Mining (SIGKDD'00)*, 71–80.
- Fisher, D. H. 1987. Knowledge Acquisition via Incremental Conceptual Clustering. *Machine Learning*, 2(2): 139–172.
- Flach, P. A. 1997. Normal Forms for Inductive Logic Programming. In *Proceedings of the 7th International Workshop of Inductive Logic Programming (ILP'97)*, 149–156.
- Freund, Y.; and Schapire, R. E. 1997. A Decision-Theoretic Generalization of On-Line Learning and an Application to Boosting. *Journal of Computer and System Sciences*, 55(1): 119–139.
- Gallant, S. I. 1988. Connectionist Expert Systems. *Communications of the ACM*, 31(2): 152–169.
- Gärdenfors, P. 1988. *Knowledge in flux*. MIT Press.
- Goldsmith, J.; and Sloan, R. H. 2005. New Horn Revision Algorithms. *Journal of Machine Learning Research*, 6: 1919–1938.
- Goldsmith, J.; Sloan, R. H.; Szörényi, B.; and Turán, G. 2004. Theory revision with queries: Horn, read-once, and parity formulas. *Artificial Intelligence*, 156(2): 139–176.
- Hatfield, J. W.; and Milgrom, P. R. 2005. Matching with Contracts. *American Economic Review*, 95(4): 913–935.
- Ignatiev, A. 2020. Towards Trustable Explainable AI. In *Proceedings of the 29th International Joint Conference on Artificial Intelligence and the 17th Pacific Rim International Conference on Artificial Intelligence (IJCAI-PRICAI'20)*, 5154–5158.
- Ignatiev, A.; and Silva, J. M. 2021. SAT-Based Rigorous Explanations for Decision Lists. In *Proceedings of the 24th International Conference on Theory and Applications of Satisfiability Testing (SAT'21)*, 251–269.
- Izza, Y.; Ignatiev, A.; and Marques-Silva, J. 2020. On Explaining Decision Trees. *CoRR*, abs/2010.11034.
- Izza, Y.; and Marques-Silva, J. 2021. On Explaining Random Forests with SAT. In *Proceedings of the 30th International Joint Conference on Artificial Intelligence (IJCAI'21)*, 2584–2591.
- Katsuno, H.; and Mendelzon, A. O. 1991. Propositional knowledge base revision and minimal change. *Artificial Intelligence*, 52: 263–294.
- Konieczny, S.; Medina Grespan, M.; and Pino Pérez, R. 2010. Taxonomy of Improvement Operators and the Problem of Minimal Change. In *Proceedings of the 12th International Conference on Principles of Knowledge Representation and Reasoning (KR'10)*, 161–170.
- Maloof, M. A.; and Michalski, R. S. 2000. Selecting Examples for Partial Memory Learning. *Machine Learning*, 41(1): 27–52.
- Maloof, M. A.; and Michalski, R. S. 2004. Incremental learning with partial instance memory. *Artificial Intelligence*, 154(1-2): 95–126.
- Mitchell, T. M. 1977. Version Spaces: A Candidate Elimination Approach to Rule Learning. In *Proceedings of the 5th International Joint Conference on Artificial Intelligence (IJCAI'77)*, 305–310.
- Mitchell, T. M. 1982. Generalization as Search. *Artificial Intelligence*, 18(2): 203–226.
- Morik, K.; Potamias, G.; Moustakis, V.; and Charissis, G. 1994. Knowledgeable learning using MOBAL: A medical case study. *Applied Artificial Intelligence*, 8(4): 579–592.
- Muggleton, S.; and De Raedt, L. 1994. Inductive Logic Programming: Theory and Methods. *Journal of Logic Programming*, 19/20: 629–679.



- Narodytska, N.; Kasiviswanathan, S. P.; Ryzhyk, L.; Sagiv, M.; and Walsh, T. 2018. Verifying Properties of Binarized Deep Neural Networks. In *Proceedings of the 32nd AAAI Conference on Artificial Intelligence (AAAI'18)*, 6615–6624.
- Ourston, D.; and Mooney, R. 1994. Theory Refinement Combining Analytical and Empirical Methods. *Artificial Intelligence*, 66: 273–309.
- Plotkin, G. D. 1970. A Note on Inductive Generalization. *Machine Intelligence*, 5: 153–163.
- Quinlan, J. R. 1986. Induction of Decision Trees. *Machine Learning*, 1(1): 81–106.
- Rosenblatt, F. 1958. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological Review*, 65(6): 386–408.
- Schlimmer, J. C.; and Fisher, D. H. 1986. A Case Study of Incremental Concept Induction. In *Proceedings of the 5th National Conference on Artificial Intelligence (AAAI'86)*, 496–501.
- Shi, W.; Shih, A.; Darwiche, A.; and Choi, A. 2020. On Tractable Representations of Binary Neural Networks. In *Proceedings of the 17th International Conference on Principles of Knowledge Representation and Reasoning (KR'20)*, 882–892.
- Shih, A.; Choi, A.; and Darwiche, A. 2018. A Symbolic Approach to Explaining Bayesian Network Classifiers. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence (IJCAI'18)*, 5103–5111.
- Shih, A.; Choi, A.; and Darwiche, A. 2019. Compiling Bayesian Networks into Decision Graphs. In *Proceedings of the 33rd AAAI Conference on Artificial Intelligence (AAAI'19)*, 7966–7974.
- Utgoff, P. E. 1989. Incremental Induction of Decision Trees. *Machine Learning*, 4: 161–186.
- Wrobel, S. 1993. On the Proper Definition of Minimality in Specialization and Theory Revision. In *Proceedings on the European Conference on Machine Learning (ECML'93)*, 65–82.
- Yang, Y.-Y. 2020. Rationalizable choice functions. *Games and Economic Behavior*, 123: 120–126.
- Zhou, Z. 2019. Abductive learning: towards bridging machine learning and logical reasoning. *Science China Information Science*, 62(7): 76101:1–76101:3.