

Basic Beliefs and Argument-Based Beliefs in Awareness Epistemic Logic with Structured Arguments

Alfredo BURRIEZA ^a, Antonio YUSTE-GINEL ^{a,1}

^a*Department of Philosophy, University of Malaga, Spain*

Abstract. There are two intuitive principles governing belief formation and argument evaluation that can potentially clash. After arguing that adopting them unrestrictedly leads to an infinite regress, we propose a formal framework in which qualified versions of both principles can be subscribed without falling into such a regress. The proposal integrates tools from two different traditions: structured argumentation and awareness epistemic logic. We show that our formalism satisfies certain rationality postulates and argue that the rest of them can be seen as too ideal when modelling resource-bounded agents.

Keywords. epistemic logic, structured argumentation, awareness logic, beliefs

1. Introduction

There exists certain tension between the formation of some epistemic attitudes of an agent and the way she assesses her available arguments. For the sake of simplicity, we will restrict our attention to the case of beliefs in what follows. The mentioned tension arises when one tries to embrace two principles that, when taken separately, seem to be intuitively acceptable:

- P1 The beliefs of an agent should be partially determined by the evaluation she performs of her available arguments. To be more precise, if an agent is considering her doxastic attitude towards a sentence φ , she should first assess her available arguments about φ and then form her belief consequently (for instance, by believing φ if she owns an accepted argument in favour of φ).² In short: belief formation is conditioned by argument evaluation.
- P2 When an agent assesses her available arguments, she should take into account her beliefs with respect to the premises. In this sense, arguments with believed premises should be taken to be stronger by the agent than arguments whose premises are not believed. In short: argument evaluation is conditioned by belief formation.

¹Corresponding Author: Office 522, Department of Philosophy, Faculty of Humanities, University of Málaga, 29010, Spain; E-mail: antonioyusteginel@gmail.com.

²The term *accepted* is extremely vague at this point, but it will be discussed and clarified later on.

Adopting P1 and P2 unrestrictedly leads to an infinite regress. To see this, let us examine the following fictional dialogue with an agent embracing P1 and P2. We start the conversation by asking: “why do you believe φ ?”. By applying P1, she would reply something like: “because I own an accepted argument α that concludes φ ”. We could ask her, in turn: “why do you accept argument α ?”. The agent might reply, applying P2: “because I believe that its premises $\text{Prem}(\alpha) = \varphi_1, \dots, \varphi_n$ are true”. Then we would ask: “why do you believe so?” and she would invoke P1 again to say that she owns accepted arguments $\alpha_1, \dots, \alpha_n$ concluding $\varphi_1, \dots, \varphi_n$. It is easy to see that this conversation could go on indefinitely.

It is worth saying that an analogous form of regression is found in the epistemological literature about the foundation of epistemic justification. Concretely, it is used as a classical argument for foundationalist theories of epistemic justification [1], in which we found inspiration for the present work. Besides, it is interesting to note that different works from the fields of formal argumentation and epistemic logic have separately subscribed different versions of P1 or P2. Let us just mention and briefly comment some of them.

Regarding P1 within formal argumentation, the idea of founding the beliefs (or knowledge) of agents on the evaluation they perform of their available arguments is already present in the seminal work of Dung [2]. This idea is recovered and further developed by frameworks of structured argumentation (e.g. [3,4]), where the sentences believed by the agent can be explicitly stated. Concurrently, epistemic logic has recently focused on the problem of including the –heretofore ignored– justification component into its formal models of knowledge and belief. This has been done in multiple manners, among which we can distinguish between syntactic and semantic approaches –where the adjectives *syntactic* and *semantic* refer to the choices for modelling justification. As for the first group of approaches, it is customary to employ justification logic (e.g. [5,6,7]). As for the second one, they have focused on how to ground the beliefs and knowledge of an agent in (possibly conflicting) pieces of evidence [8,9]. Additionally, some works (among others [10,11]) have mixed tools from formal argumentation and epistemic logic in order to develop their particular view of P1.

Regarding P2, we could say that its explicit acceptance is less spread throughout the literature. Nevertheless, in formal argumentation the idea of ordering sets of premises according to their reliability (see Section 1.2 of [12] and the references given there) can be understood as a version of P2. Besides, some works in justification logic [6,7] define the acceptance of a complex piece of evidence as the agent having a (modal) belief of its premises being true.

The main aim of this paper is to present a simple formalism (Section 2) that allows embracing explicitly qualified versions of P1 and P2 without falling into the mentioned regress (Section 3). We do so by integrating tools from awareness epistemic logic and formal argumentation. Moreover, and locating our work in the field of epistemic logic, we are interested in resource-bounded agents. This implies overcoming at least two problems: i) the classical problem of logical omniscience and ii) certain idealizations that underlie structured argumentation formalisms and that have recently been examined critically [13]. In particular, we drop the extended assumption that agents generate *all* well-shaped arguments from a given knowledge base and analyse the (negative) effects of this choice on the satisfiability of [3]’s rationality postulates (Section 4).

2. An Awareness Logic for Belief and Argumentation

The main ingredients of our logic for belief and argumentation are: (i) epistemic logic [14,15,16], a well-known tool for modelling qualitatively beliefs and knowledge of several agents; (ii) its extension with awareness operators [17] to model explicit beliefs, which allows overcoming the problem of logical omniscience (see e.g. Section 9 of [15]) and (iii) ideas taken from ASPIC⁺ [4] to model structured arguments.³ Among the most relevant features of ASPIC⁺, we highlight the following ones: a) it deals with both deductive and non-deductive (defeasible) arguments; capturing also different kinds of attacks among arguments (attacking the premises, the conclusion or the inference link) and b) it has been shown to be comprehensive, in the sense that many other proposals in structured argumentation and non-monotonic logic can be seen as special cases of it (see [4]).

Definition 1 (Language). *Let \mathbb{P} be a fixed and denumerable set of atoms; the language \mathcal{L}_{BA} is defined as the the pair $(\mathcal{F}, \mathcal{A})$ of formulas and arguments which are respectively generated by the following grammars:*

$$\begin{aligned} \varphi ::= & p \mid \neg\varphi \mid (\varphi \wedge \varphi) \mid \Box\varphi \mid \text{aware}(\alpha) \mid \text{conc}(\alpha) = \varphi \mid \\ & \text{strict}(\alpha) \mid \text{undercuts}(\alpha, \alpha) \mid \text{wellshap}(\alpha) \quad p \in \mathbb{P}, \alpha \in \mathcal{A} \\ \alpha ::= & \langle \varphi \rangle \mid \langle \alpha_1, \dots, \alpha_n \rightarrow \varphi \rangle \mid \langle \alpha_1, \dots, \alpha_n \Rightarrow \varphi \rangle \quad \varphi \in \mathcal{L}_{BA} \end{aligned}$$

Elements of \mathbb{P} represent *factual atomic sentences*, i.e. sentences about states of affairs whose truth value is agent-independent. The rest of boolean connectives are defined and read as usual. Let us adopt the following intuitive reading for the remaining formulas and arguments: $\langle \varphi \rangle$ is an atomic argument, whose only premise and conclusion is φ . $\langle \alpha_1, \dots, \alpha_n \rightarrow \varphi \rangle$ represents an argument whose last inference link strictly (deductively) concludes φ . $\langle \alpha_1, \dots, \alpha_n \Rightarrow \varphi \rangle$ represents an argument whose last inference link defeasibly concludes φ . $\Box\varphi$ means that the agent has a basic-implicit belief that φ . Basic-implicit beliefs accept different intuitive readings, both positive (reasonable assumptions, sound observations, etc) and negative (prejudices, biases, etc). The adjective *basic* underlines the idea that their source is not inferential, while *implicit* points out that they are closed under logical consequence. $\text{aware}(\alpha)$ means that the agent is aware of argument α . As usual in awareness logic [17], the operator aware admits several informal readings. For the special case of atomic arguments ($\text{aware}(\langle \varphi \rangle)$), we propose to read them as follows: “the agent recognizes her doxastic attitude toward φ through non-inferential methods”. $\text{wellshap}(\alpha)$ means that argument α is well-shaped, i.e. it has been constructed properly for the sentence it says it argues for. In more detail, every subargument of α using a strict inference link has been produced by the application of a valid deductive rule and every subargument of α using a defeasible inference link has been produced using an accepted defeasible rule. $\text{conc}(\alpha) = \varphi$ means that φ is the conclusion of α . $\text{undercuts}(\alpha, \beta)$ means that α undercuts β (i.e. α attacks β 's inference link). Finally, $\text{strict}(\alpha)$ means that α does not make use of any defeasible rule, i.e. α only contains atomic arguments and arguments formed with \rightarrow .

³We remark that the formalism below does not intend to be an alternative to ASPIC⁺, but rather an application of it to solve the conceptual problem presented in the introduction.

Definition 2 (Argument structure [4]). *Let us define the following meta-syntactic functions for analysing an argument's structure:*

$\text{Prem}(\alpha)$ returns the *premises* of α and it is defined as follows: $\text{Prem}(\langle\varphi\rangle) = \{\varphi\}$, $\text{Prem}(\langle\alpha_1, \dots, \alpha_n \hookrightarrow \varphi\rangle) = \text{Prem}(\alpha_1) \cup \dots \cup \text{Prem}(\alpha_n)$ where $\hookrightarrow \in \{\rightarrow, \Rightarrow\}$. **Example:** $\text{Prem}(\langle\langle\langle p \rangle, \langle q \rangle \Rightarrow r \rangle \Rightarrow s \rangle \rightarrow s \vee t) = \{p, q\}$.

$\text{Conc}(\alpha)$ returns the *conclusion* of α and it is defined as follows $\text{Conc}(\langle\varphi\rangle) = \{\varphi\}$ and $\text{Conc}(\langle\alpha_1, \dots, \alpha_n \hookrightarrow \varphi\rangle) = \{\varphi\}$ where $\hookrightarrow \in \{\rightarrow, \Rightarrow\}$. Note that arguments of ASPIC⁺ have unique conclusions (differently to what happens, for instance, in justification logic [6] where the + operator allows for arguments with multiple conclusions). **Example:** $\text{Conc}(\langle\langle\langle p \rangle, \langle q \rangle \Rightarrow r \rangle \Rightarrow s \rangle \rightarrow s \vee t) = s \vee t$.

$\text{sub}_A(\alpha)$ returns the *subarguments* of α and it is defined as follows: $\text{sub}_A(\langle\varphi\rangle) = \{\langle\varphi\rangle\}$ and $\text{sub}_A(\langle\alpha_1, \dots, \alpha_n \hookrightarrow \varphi\rangle) = \{\langle\alpha_1, \dots, \alpha_n \hookrightarrow \varphi\rangle\} \cup \text{sub}_A(\alpha_1) \cup \dots \cup \text{sub}_A(\alpha_n)$ where $\hookrightarrow \in \{\rightarrow, \Rightarrow\}$. **Example:** $\text{sub}_A(\langle\langle\langle p \rangle, \langle q \rangle \Rightarrow r \rangle \Rightarrow s \rangle \rightarrow s \vee t) = \{\langle\langle\langle p \rangle, \langle q \rangle \Rightarrow r \rangle \Rightarrow s \rangle \rightarrow s \vee t, \langle\langle p \rangle, \langle q \rangle \Rightarrow r \rangle \Rightarrow s, \langle\langle p \rangle, \langle q \rangle \Rightarrow r \rangle, \langle p \rangle, \langle q \rangle\}$.

$\text{TopRule}(\alpha)$ returns the *top rule* of α , i.e. the last one applied in the formation of α . It is defined as follows: $\text{TopRule}(\langle\varphi\rangle)$ is left undefined, $\text{TopRule}(\langle\alpha_1, \dots, \alpha_n \rightarrow \varphi\rangle) = \text{TopRule}(\langle\alpha_1, \dots, \alpha_n \Rightarrow \varphi\rangle) = ((\text{Conc}(\alpha_1), \dots, \text{Conc}(\alpha_n)), \varphi)$. **Example:** $\text{TopRule}(\langle\langle\langle p \rangle, \langle q \rangle \Rightarrow r \rangle \Rightarrow s \rangle \rightarrow s \vee t) = (s, s \vee t)$.

$\text{DefRule}(\alpha)$ returns the set of *defeasible rules* of α and it is defined as $\text{DefRule}(\langle\varphi\rangle) = \emptyset$, $\text{DefRule}(\langle\alpha_1, \dots, \alpha_n \rightarrow \varphi\rangle) = \text{DefRule}(\alpha_1) \cup \dots \cup \text{DefRule}(\alpha_n)$ and $\text{DefRule}(\langle\alpha_1, \dots, \alpha_n \Rightarrow \varphi\rangle) = \{((\text{Conc}(\alpha_1), \dots, \text{Conc}(\alpha_n)), \varphi)\} \cup \text{DefRule}(\alpha_1) \cup \dots \cup \text{DefRule}(\alpha_n)$. **Example:** $\text{DefRule}(\langle\langle\langle p \rangle, \langle q \rangle \Rightarrow r \rangle \Rightarrow s \rangle \rightarrow s \vee t) = \{((p, q), r), ((r), s)\}$.

Let us also define *single negations*, for any $\varphi \in \mathcal{L}_{BA}$: $\sim \varphi := \psi$ if φ is of the form $\neg\psi$; else $\sim \varphi := \neg\varphi$.

Definition 3 (Model). *A model for \mathcal{L}_{BA} is a tuple $M = (W, \mathcal{B}, \mathcal{O}, \mathcal{D}, n, \|\cdot\|)$ where:*

- $W \neq \emptyset$ is a set of possible worlds
- $\mathcal{B} \subseteq W$ and $\mathcal{B} \neq \emptyset$ is the set of doxastically indistinguishable worlds
- $\mathcal{O} \subseteq \mathcal{A}$ is the (finite) set of available arguments or the awareness set of the agent
- $\mathcal{D} \subseteq \mathcal{L}_{BA}^n \times \mathcal{L}_{BA}$ (with $n \in \mathbb{N}$) is a finite set of accepted defeasible rules s.t. if $((\varphi_1, \dots, \varphi_n), \varphi) \in \mathcal{D}$, then $\{\varphi_1, \dots, \varphi_n, \varphi\} \not\vdash_0 \perp$; where \vdash_0 is the consequence relation of classical propositional logic
- $n: \mathcal{D} \rightarrow \mathbb{P}$ is a (possibly partial) naming function for defeasible rules, where $n(R)$ informally means “the defeasible rule R is applicable”
- $\|\cdot\|: \mathbb{P} \rightarrow \wp(W)$ is an atomic valuation

Definition 4 (Truth). *Formulas of \mathcal{L}_{BA} are interpreted in pointed models (M, w) where $w \in W$. $M, w \models \varphi$ means that φ is true in (M, w) . \models is defined for every kind of formulas as follows (we omit the clauses for propositional variables and boolean connectives):*

- $M, w \models \Box\varphi$ iff for all $w' \in W$: $w' \in \mathcal{B}$ implies $M, w' \models \varphi$
- $M, w \models \text{aware}(\alpha)$ iff $\alpha \in \mathcal{O}$
- $M, w \models \text{conc}(\alpha) = \varphi$ iff $\text{Conc}(\alpha) = \varphi$
- $M, w \models \text{strict}(\alpha)$ iff $\text{DefRule}(\alpha) = \emptyset$

- $M, w \models \text{undercuts}(\alpha, \beta)$ iff $\text{Conc}(\alpha) = \sim n(\text{TopRule}(\beta))$ ⁴
- $M, w \models \text{wellshap}(\langle \varphi \rangle)$
- $M, w \models \text{wellshap}(\langle \alpha_1, \dots, \alpha_n \rightarrow \varphi \rangle)$ iff $M, w \models \text{wellshap}(\alpha_i)$ for every $1 \leq i \leq n$ and $\{\text{Conc}(\alpha_1), \dots, \text{Conc}(\alpha_n)\} \vdash_0 \varphi$
- $M, w \models \text{wellshap}(\langle \alpha_1, \dots, \alpha_n \Rightarrow \varphi \rangle)$ iff $M, w \models \text{wellshap}(\alpha_i)$ for every $1 \leq i \leq n$ and $((\text{Conc}(\alpha_1), \dots, \text{Conc}(\alpha_n)), \varphi) \in \mathcal{D}$

Validity ($\models \varphi$) and local logical consequence ($\Gamma \models \varphi$) are defined as usual [18]. Note that our way of representing basic-implicit beliefs is equivalent (in the single-agent case) to have a Kripke model where the accessibility relation is serial, transitive and euclidean; therefore \Box satisfies *KD45* axioms (see [10,16]). Regarding the truth clauses for $\text{conc}(\alpha) = \varphi$ and $\text{strict}(\alpha)$; it is easy to show that these kinds of formulas are model independent (since they are based on argument structure, see Definition 2). This implies that, for these kinds of formulas they are true in a pointed model iff they are valid. Furthermore, note that the clause for $\Box\varphi$, $\text{undercuts}(\alpha, \beta)$, $\text{aware}(\alpha)$ and $\text{wellshap}(\alpha)$ makes the satisfiability of these kinds of formulas world-independent, i.e. they are true in a pointed model if they are globally true in the model. Consequently, we have that $\star \rightarrow \Box\star$ and $\neg\star \rightarrow \Box\neg\star$ are valid schemata, where $\star \in \{\text{aware}(\alpha), \text{conc}(\alpha) = \varphi, \text{undercuts}(\alpha, \beta), \text{wellshap}(\alpha), \text{strict}(\alpha)\}$. Informally, this amounts to assume that: i) awareness of arguments is fully introspective w.r.t. basic-implicit beliefs and ii) the agent is logically competent w.r.t. the arguments she is aware of. However, and unlike what is usual in structured argumentation [4]; our agent will not work with the whole set of all well-shaped arguments (which is by definition infinite), but rather with the (finite) set of arguments that she is aware of.

3. Basic Beliefs and AB-Beliefs in \mathcal{L}_{BA}

In order to solve the tension between P1 and P2, we distinguish between basic-explicit beliefs (Definition 5) and argument-based beliefs (AB-Beliefs, for short; Definition 9). While the notion of basic belief (both its implicit and explicit versions) only needs some informal clarification (Section 3.1); AB-beliefs force us to import some concepts from formal argumentation (sections 3.2, 3.3 and 3.4), especially from ASPIC⁺. Most of the central concepts used in ASPIC⁺ (or our adaptations) are definable in \mathcal{L}_{BA} .

3.1. Basic Beliefs

Recall that basic-implicit beliefs are represented through the primitive, normal modal operator \Box , hence they suffer from logical omniscience: $(M, w \models \Box\psi$ for all $\psi \in \Gamma$ and $\Gamma \models \varphi$) implies $M, w \models \Box\varphi$. This property has been extensively discussed in the epistemic logic literature, and it has been argued to be problematic when dealing with resource-bounded agents (see e.g. [17,15, Chapter 9]). The pitfall can be overcome by distinguishing between basic-implicit beliefs ($\Box\varphi$) and basic-explicit beliefs ($\Box^e\varphi$) following the awareness approach [17]:

Definition 5 (Basic-explicit beliefs). $\Box^e\varphi := \Box\varphi \wedge \text{aware}(\langle \varphi \rangle)$

⁴ Note that we do not need to consider undercuts as a primitive operator, since it could be defined through a (simpler) operator that captures the meaning of n . We make this choice for the sake of succinctness.

Informally, basic-explicit beliefs can be generally understood as actual beliefs of the agent whose justification is not inferential (it comes from other epistemic phenomena, such as observations or reliable communications). In fact, we can think of many beliefs that a (reasonable) epistemic agent may have that need no arguments to be justified. Imagine, for instance, that you walk into your classroom and you see three students in there. Consequently, you form the belief that “there are three student in the classroom”. Do you need any complex argument to justify such a belief? Our claim is that, in principle, you do not. Indeed, you can form arguments supporting the proposition if someone would question your belief. But, for the agent herself (you, in this case), mere observation is a good enough reason to believe that there are three students in the classroom.

3.2. Doxastic Preference

Premises are usually understood as a source of argument strength [12], regarding the *support dimension* (see [12] for the distinction between the three dimensions or tiers of argument strength). In structured argumentation [4,12], this is often modelled by stratifying a given set of formulas into different preference classes. Such a hierarchy is usually assumed to be primitive and its nature is abstracted away from the modelling process. Let us now show how basic beliefs induce a meaningful hierarchy of this kind. Let (M, w) be a pointed model and let $\alpha \in \mathcal{A}$, we can distinguish between three types of premises of α : $\text{Prem}(\alpha) = \text{Prem}^+(\alpha) \cup \text{Prem}^?(\alpha) \cup \text{Prem}^-(\alpha)$ where each component is defined as follows $\text{Prem}^+(\alpha) := \{\varphi \in \text{Prem}(\alpha) \mid M, w \models \Box\varphi\}$ (the set of trusted or believed premises); $\text{Prem}^?(\alpha) := \{\varphi \in \text{Prem}(\alpha) \mid \neg\Box\varphi \wedge \neg\Box\neg\varphi\}$ (the set of premises considered contingent by the agent) and $\text{Prem}^-(\alpha) := \{\varphi \in \text{Prem}(\alpha) \mid M, w \models \Box\neg\varphi\}$ (the set of disbelieved premises). The three kind of premises are pairwise disjoint (due to the consistency of basic beliefs) and possibly empty. Furthermore, this distinction induces another one within the set of all arguments $\mathcal{A} = \mathcal{A}^+ \cup \mathcal{A}^? \cup \mathcal{A}^-$ where each component is defined as follows: $\mathcal{A}^+ := \{\alpha \in \mathcal{A} \mid \text{Prem}(\alpha) = \text{Prem}^+(\alpha)\}$; $\mathcal{A}^? := \{\alpha \in \mathcal{A} \mid \text{Prem}(\alpha) = \text{Prem}^+(\alpha) \cup \text{Prem}^?(\alpha), \text{Prem}^?(\alpha) \neq \emptyset\}$ and $\mathcal{A}^- = \{\alpha \in \mathcal{A} \mid \text{Prem}^-(\alpha) \neq \emptyset\}$.⁵ It seems natural to assume the following preference ordering between the three classes of arguments $\mathcal{A}^+ \supset_p \mathcal{A}^? \supset_p \mathcal{A}^-$, that can be lowered to arguments straightforwardly: $\alpha >_p \beta$ iff $\alpha \in \mathcal{A}^+$, $\beta \in \mathcal{A}^-$ and $\mathcal{A}^+ \supset_p \mathcal{A}^-$ with $’, ’’ \in \{+, ?, -\}$. The relation $>_p$ is precisely our qualified version of P2: *argument evaluation is conditioned by basic belief formation*. Interestingly enough, this relation can be captured in \mathcal{L}_{BA} , as shown in [19], using the following shorthands: $\text{accept}(\alpha) := \bigwedge_{\varphi \in \text{Prem}(\alpha)} \Box\varphi$ ⁶ (*basic acceptance*); $\text{reject}(\alpha) := \bigvee_{\varphi \in \text{Prem}(\alpha)} \Box\neg\varphi$ (*basic rejection*); $\text{prem}^>(\alpha, \beta) := (\text{accept}(\alpha) \wedge \neg\text{accept}(\beta)) \vee (\neg\text{reject}(\alpha) \wedge \text{reject}(\beta))$; $\text{prem}^{\approx}(\alpha, \beta) := \neg\text{prem}^>(\alpha, \beta) \wedge \neg\text{prem}^>(\beta, \alpha)$:

Proposition 1. *Let (M, w) be a pointed model, we have that $M, w \models \text{prem}^>(\alpha, \beta)$ iff $\alpha >_p \beta$.*

⁵The *lifting principle* applied in order to go from preferences between premises to preference between arguments is the so-called *min-min* principle [12]. Note that basic beliefs permit more fine-grained distinctions regarding the relative strength of arguments. For instance, we could distinguish within $\mathcal{A}^?$ between arguments whose premises are jointly considered a doxastic possibility $\mathcal{A}^{?+} := \{\alpha \in \mathcal{A} \mid \diamond \bigwedge_{\varphi \in \alpha} \varphi\}$ and arguments that do not enjoy this property $\mathcal{A}^{?-} := \mathcal{A}^? / \mathcal{A}^{?+}$. Nonetheless, we adopt the current division for simplicity.

⁶This definition is inspired by [6].

Premises are not the only source of argument strength regarding the *support dimension*. The other main source are inference links. In order to keep things simple, we adopt a minimal (yet intuitively acceptable) principle to assess inference links: *ceteris paribus*, strict arguments should be preferred to defeasible arguments. This principle can be captured as follows: let $\mathcal{A}^{st} := \{\alpha \in \mathcal{A} \mid \text{DefRule}(\alpha) = \emptyset\}$ and let $\mathcal{A}^{df} := \mathcal{A} / \mathcal{A}^{st}$, we can define new preference classes by intersecting separately both sets with the previous hierarchy. Furthermore, we assume the following natural preference ordering:

$$\mathcal{A}^+ \cap \mathcal{A}^{st} \sqsupset_{il} \mathcal{A}^+ \cap \mathcal{A}^{df} \sqsupset_{il} \mathcal{A}^? \cap \mathcal{A}^{st} \sqsupset_{il} \mathcal{A}^? \cap \mathcal{A}^{df} \sqsupset_{il} \mathcal{A}^- \cap \mathcal{A}^{st} \sqsupset_{il} \mathcal{A}^- \cap \mathcal{A}^{df}$$

The new preference ordering can be lowered to arguments as follows: $\alpha >_{il} \beta$ iff $\alpha \in \mathcal{A}', \beta \in \mathcal{A}''$ and $\mathcal{A}' \sqsupset_{il} \mathcal{A}''$ with $\mathcal{A}', \mathcal{A}'' \in \{\mathcal{A}^+ \cap \mathcal{A}^{st}, \mathcal{A}^+ \cap \mathcal{A}^{df}, \mathcal{A}^? \cap \mathcal{A}^{st}, \mathcal{A}^? \cap \mathcal{A}^{df}, \mathcal{A}^- \cap \mathcal{A}^{st}, \mathcal{A}^- \cap \mathcal{A}^{df}\}$. Note that the relation satisfies $>_p \subset >_{il}$. Besides, it can be captured in \mathcal{L}_{BA} through the following schemes: $\text{strict}^>(\alpha, \beta) := \text{strict}(\alpha) \wedge \neg \text{strict}(\beta)$; $\alpha > \beta := \text{prem}^>(\alpha, \beta) \vee (\text{prem}^\approx(\alpha, \beta) \wedge \text{strict}^>(\alpha, \beta))$; $\alpha \geq \beta := \neg(\beta > \alpha)$; $\alpha \approx \beta := \alpha \geq \beta \wedge \beta \geq \alpha$.

Proposition 2. *Let (M, w) be a pointed model, we have that $M, w \models \alpha \geq \beta$ iff $\alpha \geq_{il} \beta$.*

Let us stress two points regarding the preference ordering \geq_{il} which are important for the study of [3]’s rationality postulates. First, it is *reasonable* in the sense of [4]. Second, \geq_{il} is a total preorder on \mathcal{A} . This fact, expressed in the object language has the form of the following valid schemas, for every $\alpha, \beta, \delta \in \mathcal{A}$: $\models (\alpha \geq \beta \wedge \beta \geq \delta) \rightarrow \alpha \geq \delta$ (transitivity) and $\models \alpha \geq \beta \vee \beta \geq \alpha$ (connectedness).

3.3. Attack and Defeat

Agents do not assess arguments in isolation, or merely pairwise, checking if certain features of the involved premises and inference links are good enough to support the conclusion. Another important dimension of argument strength is called the *dialectical tier* which, following [12], is “mainly represented by relations of argumentative attack and defeat between arguments”. \mathcal{L}_{BA} is rich enough to capture the three customary kinds of attacks discussed in structured argumentation:

Definition 6 (Argument attack). *Given a pointed model (M, w) and $\alpha, \beta \in \mathcal{A}$: we say that α undermines β iff $M, w \models \text{undermines}(\alpha, \beta)$, where $\text{undermines}(\alpha, \beta) := \bigvee_{\varphi \in \text{Prem}(\beta)} \text{conc}(\alpha) = \sim \varphi$; α rebuts β iff $M, w \models \text{rebuts}(\alpha, \beta)$, where $\text{rebuts}(\alpha, \beta) := \bigvee_{\langle \beta_1, \dots, \beta_n \mapsto \varphi \rangle \in \text{sub}_A(\beta)} \text{conc}(\alpha) = \sim \varphi$ where $\mapsto \in \{\rightarrow, \Rightarrow\}$; and α undercuts β iff $M, w \models \text{undercuts}^*(\alpha, \beta)$, where $\text{undercuts}^*(\alpha, \beta) := \bigvee_{\beta' \in \text{sub}_A(\beta)} \text{undercuts}(\alpha, \beta')$.*

Our definition of attack integrates a notion of *unrestricted rebuttal*, in the sense that rebuttals are permitted on any kind of complex argument. This is indeed polemic. While the creators of ASPIC⁺, amongst others, only allow rebuttals on the application of defeasible rules; others have argued that the unrestricted notion seems natural in dialectical contexts [20,21]. Moreover, [21] requires the rebutted argument to be defeasible (non-strict) while [20] does not require it but, in turn, this feature is implied by their setting. We permit *completely unrestricted rebuttals* for a simple reason: since awareness sets do not exhibit any closure property, in absence of completely unrestricted rebuttal direct consistency fails (see Section 4 for more details).

From an agent perspective, some attacks must be disregarded. Imagine, for instance, that an agent is aware of $\langle\langle p \rangle \rightarrow p \vee q \rangle$ and that she accepts it in a doxastic sense ($\Box p$). She then receives an undermining argument $\langle\langle r \rangle \Rightarrow \neg p \rangle$ but she does not accept it (she does not believe that r). It seems that such an attack must not be considered a threat for the agent. Consequently, the notion of *defeat* should take into account the preference relation defined above. We import the definition of defeat from ASPIC⁺ to our object language, introducing two essential differences. First, preferences do play a role when determining the success of undercutting attacks (the reason for doing so is offered below). Second, the agent only considers defeats among the well-shaped arguments that she is aware of, capturing that although her resources are bounded (w.r.t. argument generation) they are locally well applied. We proceed in two steps: defining a successful counterpart for each type of attack and adding the awareness/well-shapedness requirement.

Definition 7 (Successful attack, defeat). *Given a pointed model (M, w) and two arguments $\alpha, \beta \in \mathcal{A}$ we say that: α successfully undermines β iff $M, w \models \text{SuUndermines}(\alpha, \beta)$, where $\text{SuUndermines}(\alpha, \beta) := \bigvee_{\varphi \in \text{Prem}(\beta)} (\text{conc}(\alpha) = \sim\varphi \wedge \alpha \geq \langle\varphi\rangle)$; α successfully rebuts β iff $M, w \models \text{SuRebuts}(\alpha, \beta)$ where $\text{SuRebuts}(\alpha, \beta) := \bigvee_{\langle\beta_1, \dots, \beta_n \leftrightarrow \varphi\rangle \in \text{sub}_A(\beta)} (\text{conc}(\alpha) = \sim\varphi \wedge \alpha \geq \langle\beta_1, \dots, \beta_n \leftrightarrow \varphi\rangle)$; α successfully undercuts β iff $M, w \models \text{SuUndercuts}(\alpha, \beta)$, where $\text{SuUndercuts}(\alpha, \beta) := \bigvee_{\beta' \in \text{sub}_A(\beta)} (\text{undercuts}(\alpha, \beta') \wedge \alpha \geq \beta')$ and, finally, we say that α defeats β iff $M, w \models \text{defeat}(\alpha, \beta)$, where $\text{defeat}(\alpha, \beta) := (\text{SuUndermines}(\alpha, \beta) \vee \text{SuRebuts}(\alpha, \beta) \vee \text{SuUndercuts}(\alpha, \beta)) \wedge \text{aware}(\alpha) \wedge \text{aware}(\beta) \wedge \text{wellshap}(\alpha) \wedge \text{wellshap}(\beta)$.*

As mentioned above, it has been argued that undercutting attacks always succeed (independently from what the preferences are) [4]. This may lead to counter-intuitive cases in the current setting. Taking the same example that [4], due to Pollock, suppose that an agent considers that an object is red because she sees that it is red (she is aware of an argument $\langle\langle \text{SeeRed} \rangle \Rightarrow \text{IsRed} \rangle$). Suppose that someone suggests her to consider the undercutting “there might be a red shining, therefore the inference rule you are applying does not hold”. This can be modelled by putting into her awareness set an argument $\langle\langle \text{RedLight} \rangle \Rightarrow \neg D \rangle$ where D is an atomic proposition saying that the defeasible inference rule $(\langle\langle \text{SeeRed} \rangle, \text{IsRed} \rangle)$ is applicable. Suppose however that she believes that there is no such light in the room, $M, w \models \Box \neg \text{RedLight}$. It looks that, under this assumption, $\langle\langle \text{RedLight} \rangle \Rightarrow \neg D \rangle$ is not a good reason to prevent the agent from drawing her initial conclusion that IsRed holds.

3.4. AB-Beliefs

Given a set of well-shaped and owned arguments B , the agent is already able to determine the defeat relation among them. Nevertheless, the question of how to decide which subset(s) of B should be considered *justified* remains still open. This question has been called the *evaluation tier* of argument strength in [12] and it is notoriously solved by applying different semantics to an *argumentation framework* (first introduced by Dung in [2]). Note that each pointed model (M, w) naturally induces a Dung-style argumentation framework [2] (AF, for short), which will be the main construct to define AB-beliefs.

Definition 8 (Associated argumentation framework). *Let (M, w) be a pointed model where $M = (W, \mathcal{B}, \mathcal{C}, \mathcal{D}, \mathbf{n}, || \cdot ||)$. The argumentation framework associated to (M, w) ,*

denoted by AF^M is the pair $(\mathcal{O}^{ws}, \rightsquigarrow)$ where $\mathcal{O}^{ws} := \{\alpha \in \mathcal{O} \mid M, w \models \text{wellshap}(\alpha)\}$ and $\rightsquigarrow \subseteq \mathcal{O}^{ws} \times \mathcal{O}^{ws}$ is defined as $\alpha \rightsquigarrow \beta$ iff $M, w \models \text{defeat}(\alpha, \beta)$.⁷

The semantics of an AF is usually given in terms of *extensions*, i.e. subsets of \mathcal{O}^{ws} satisfying certain intuitive constraints to be an acceptable set [2]. Given a set of arguments $B \subseteq \mathcal{O}^{ws}$, typical minimal requirements are *conflict-freeness* (there are no $\alpha, \beta \in B$ s.t. $\alpha \rightsquigarrow \beta$) and *self-defence* (every defeater of members of B is in turn defeated by some member of B). A set of arguments B is a *complete extension* iff it contains precisely the arguments it defends. Finally, the *grounded extension* of AF^M , denoted by $GE(AF^M)$ is the minimal (w.r.t. set inclusion) complete extension. For a more precise definition of these notions and an extensive discussion about the existing semantics, the reader is referred to [22]. Our choice of grounded semantics for defining AB-beliefs is rooted on the arguments presented in [23] for such a decision regarding epistemic reasoning.

Definition 9 (AB-Beliefs). *Let (M, w) be a pointed model for \mathcal{L}_{BA} , and let $\varphi \in \mathcal{L}_{BA}$, we say that φ is AB-believed in (M, w) , denoted by $M, w \models B^{AB} \varphi$, iff $\exists \alpha \in GE(AF^M)$: $\text{Conc}(\alpha) = \varphi$.*

This definition captures our qualified version of P1: *AB-belief formation is conditioned by argument evaluation*. Moreover, note that the following schema is valid $\models \Box^e \varphi \rightarrow B^{AB} \varphi$ (i.e. basic-explicit beliefs are a special case of AB-beliefs). AB-beliefs cannot be captured in \mathcal{L}_{BA} . The reason for this is that its definition quantifies over arguments (and sets of arguments, since the grounded extension requires subset-minimality). This inconvenience could be circumvented in several ways that are out of the scope of this paper. Instead, let us just increase \mathcal{L}_{BA} with a new clause $B^{AB} \varphi$, where $\varphi \in \mathcal{L}_{BA}$, and adopt the truth clause of Definition 9 for the new kind of formulas. In the following example, we illustrate the difference between both kinds of beliefs and how our qualified versions of P1 and P2 work.

Example 1 (Assessing a survey). *A researcher in charge of a survey (in what follows, the agent) is assessing the last report of her team. In particular, the agent is wondering whether a Claim follows from some Data gathered by her team, as suggested in the report, i.e. she is determining the acceptability of $\langle \langle \text{Data} \rangle \Rightarrow \text{Claim} \rangle$. Model M , depicted in the top part of Figure 1 shows her implicit doxastic attitudes towards the involved propositions. The bottom-part of the same figure shows the associated AF, AF^M , where black arrows represent defeats and dashed arrows represent unsuccessful attacks. Some elements of the model are omitted in the representation (\mathcal{O} , \mathcal{D} and n), but they can be completed by observing the associated AF.*

The head of the laboratory has told the agent to consider the undercutting attack $\langle \langle \neg \text{Honest} \rangle \Rightarrow \neg \text{Reliab} \rangle$, according to which if her team is not behaving honestly, the defeasible rule $((\text{Data}), \text{Claim})$ should be considered suspicious (we fix $n((\text{Data}), \text{Claim})) = \text{Reliab}$). Nevertheless, the agent holds a basic-explicit belief that her team is behaving honestly, $M, w \models \Box^e \text{Honest}$; so she disregards the mentioned undercutting, $M, w \models \neg \text{SuUndercuts}(\langle \langle \neg \text{Honest} \rangle \Rightarrow \neg \text{Reliab} \rangle, \langle \langle \text{Data} \rangle \Rightarrow \text{Claim} \rangle)$. Moreover, she also considers the strict argument $\langle \langle \text{Defective} \rightarrow \neg \text{Data} \rangle, \langle \text{Defective} \rangle \rightarrow \neg \text{Data} \rangle$ ac-

⁷ Given the simplification of the modal semantics we have assumed, it can be shown that for every model M , with domain W , it holds that $AF^{M,w} = AF^{M,w'}$ for every $w, w' \in W$. This remark permits us to refer to $AF^{M,w}$ just as AF^M .

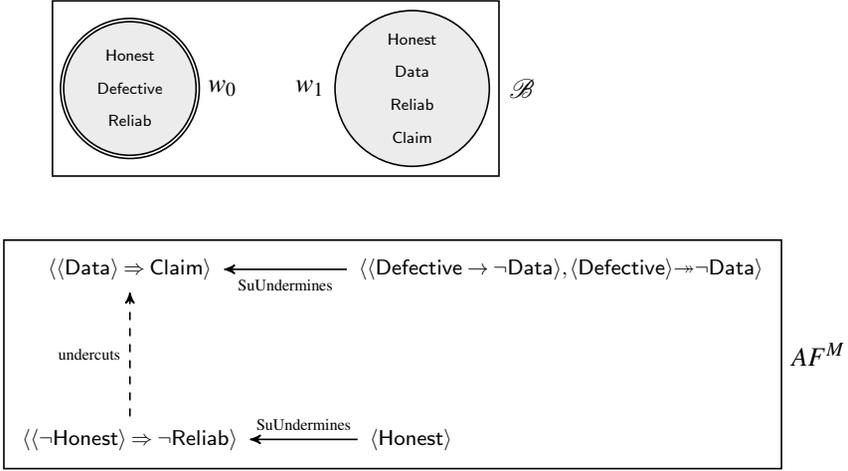


Figure 1. Pointed model (M, w_0) (top part) and its associated AF, AF^M (bottom part).

ording to which if one of the measure devices used in the study is defective, then the gathered data is not true. Note that this argument undermines $\langle\langle \text{Data} \rangle \Rightarrow \text{Claim} \rangle$. Due to previous problems with the mentioned device, she considers as doxastically possible a situation where it does not work properly (w_0), hence the undermining succeeds. Consequently, she keeps sceptic about the value of Claim: $M, w \models \neg B^{AB} \text{Claim} \wedge \neg B^{AB} \neg \text{Claim}$.

4. Rationality Postulates

In [3], Caminada and Amgoud provide a list of rationality postulates that a good argumentation formalism must satisfy. In [4], Modgil and Prakken discuss these postulates in relation to ASPIC⁺. In this section, we offer sufficient conditions for two of them to be satisfied and argue that the other two are too idealistic in an epistemic logic for resource-bounded agents. First of all, let us formulate the postulates in the current setting. Let AF^M be an associated AF, we say that AF^M satisfies:

- RP_{SUB} (sub-argument closure) iff for any $\alpha \in GE(AF^M)$, $\text{sub}_A(\alpha) \subseteq GE(AF^M)$
- RP_{DC} (direct consistency) iff $\nexists \varphi \in \mathcal{L}_{BA}$: $\varphi, \sim\varphi \in \text{Conc}(GE(AF^M))$ ⁸
- RP_{CL} (closure under strict rules) iff for all $\varphi \in \mathcal{L}_{BA}$ s.t. $\text{Conc}(GE(AF^M)) \vdash_0 \varphi$ it holds that $\varphi \in \text{Conc}(GE(AF^M))$
- RP_{IC} (indirect consistency) iff $\text{Conc}(GE(AF^M)) \not\vdash_0 \perp$

The following propositions establish sufficient conditions for RP_{SUB} (resp. RP_{DC}) to be satisfied by an associated AF:

Proposition 3. *Let (M, w) be a pointed model, where $M = (W, \mathcal{B}, \mathcal{O}, \mathcal{D}, n, \|\cdot\|)$. If \mathcal{O} is closed under subarguments (i.e. $\alpha \in \mathcal{O}$ implies $\text{sub}_A(\alpha) \subseteq \mathcal{O}$), then $GE(AF^M)$ is closed under subarguments.*

⁸We lift the domain of the function Conc from arguments to sets of arguments as follows: $\text{Conc}(\mathcal{S}) := \{\text{Conc}(\alpha) \mid \alpha \in \mathcal{S}\}$ for any $\mathcal{S} \subseteq \mathcal{A}$.

Proposition 4. *Let (M, w) be a pointed model, then AF^M satisfies direct consistency.*

Remark. *In the current setting, it is crucial for Proposition 4 to hold that we allow completely unrestricted rebuttals (see Definition 6 and the subsequent discussion).*

RP_{CL} and RP_{IC} are violated by the current framework. Let us show why this happens and why it is not an unavoidable inconvenience for our purposes. First of all, note that RP_{CL} cannot be satisfied by *any* associated AF. Note that $\text{Conc}(GE(AF^M)) \vdash_0 \{\varphi \in \mathcal{L}_{BA} \mid \vdash_0 \varphi\}$ for any model M . Therefore, for RP_{CL} to be true, it should hold that $\{\varphi \in \mathcal{L}_{BA} \mid \vdash_0 \varphi\} \subseteq \text{Conc}(GE(AF^M))$. But this is impossible since $\{\varphi \in \mathcal{L}_{BA} \mid \vdash_0 \varphi\}$ is infinite and $\text{Conc}(GE(AF^M))$ is finite by assumption (because awareness sets are finite by assumption). Nevertheless, RP_{CL} is just a special case of logical omniscience (propositional logical omniscience); so its satisfiability should not be pursued when modelling resource-bounded agents. As pointed out in [3], this problem can be avoided using query-based implementations for computing the grounded extension. This strategy does not seem appropriate in the current context, since it still would require to generate the whole set of well-shaped arguments.

As for RP_{IC} , its failure is more threatening. Moreover, our agent fails to have the following forms of consistency (that fall between direct and indirect consistency): (i) there is no $\varphi \in \text{Conc}(GE(AF^M))$ such that $\{\varphi\} \vdash_0 \perp$ and (ii) there are no $\varphi, \psi \in \text{Conc}(GE(AF^M))$ such that $\{\varphi, \psi\} \vdash_0 \perp$. These facts reveal the minimal character of our formalism. Note however, that the first case can be avoided by closing \mathcal{O} under conclusions and single negations. The second case can in turn be overcome by defining $\sim\varphi := \{\psi \mid \{\varphi, \psi\} \vdash_0 \perp\}$. Be as it may, failure of different forms of consistency are understood as pitfalls in many different contexts. However, at the same time, it seems plausible to claim that reasonable (yet not fully rational) agents can have indirectly inconsistent AB-beliefs; as far as they keep their AB-beliefs being directly consistent (see e.g. [24, §2] for a defence of this kind of inconsistencies). Note that although AB-beliefs might be indirectly inconsistent, they are not trivial (agents never end up believing *everything*). Moreover, if one wants to strengthen the reasoning skills of the modelled agent, two interesting questions arise. First, is there any set of sufficient conditions that guarantees the satisfaction of RP_{IC} in \mathcal{L}_{BA} while keeping awareness sets finite? A positive answer might not be trivial, since the satisfaction of RP_{IC} is usually proved as a corollary of RP_{DC} and RP_{CL} [3,4]. Second, given an indirectly inconsistent associated AF, is there an action (or sequence of actions) such that indirect inconsistency is recovered?

5. Future Work

Besides the open problems mentioned in the last section, there are several questions that require further study. We highlight the following ones. First, examining \mathcal{L}_{BA} on the view of additional postulates (see [25]). Second, it would also be interesting to study whether it is possible to characterize axiomatically the behaviour of B^{AB} , when treated as a primitive operator.

Acknowledgements We thank the anonymous reviewers for valuable and insightful comments, some of which had to be left out due to lack of space unfortunately. The research activity of A. Yuste-Ginel is supported by the predoctoral grant MEC-D-FPU 2016/04113.

References

- [1] Hasan A, Fumerton R. Foundationalist Theories of Epistemic Justification. In: Zalta EN, editor. *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University; 2018.
- [2] Dung PM. On the acceptability of arguments and its fundamental role in nonmonotonic reasoning, logic programming and n-person games. *Artificial Intelligence*. 1995;77(2):321–357.
- [3] Caminada M, Amgoud L. On the evaluation of argumentation formalisms. *Artificial Intelligence*. 2007;171(5-6):286–310.
- [4] Modgil S, Prakken H. A general account of argumentation with preferences. *Artificial Intelligence*. 2013;195:361–397.
- [5] Artemov S. The logic of justification. *The Review of Symbolic Logic*. 2008;1(4):477–513.
- [6] Baltag A, Renne B, Smets S. The Logic of Justified Belief Change, Soft Evidence and Defeasible Knowledge. In: Ong L, de Queiroz R, editors. *Logic, Language, Information and Computation*. WoLLIC 2012. LNCS. vol. 7456. Springer; 2012. p. 168–190.
- [7] Baltag A, Renne B, Smets S. The logic of justified belief, explicit knowledge, and conclusive evidence. *Annals of Pure and Applied Logic*. 2014;165(1):49–81.
- [8] van Benthem J, Pacuit E. Dynamic logics of evidence-based beliefs. *Studia Logica*. 2011;99(1-3):61.
- [9] Baltag A, Bezhanishvili N, Özgün A, Smets S. Justified Belief and the Topology of Evidence. In: Väänänen J, Hirvonen Å, de Queiroz R, editors. *Logic, Language, Information, and Computation*. Springer; 2016. p. 83–103.
- [10] Grossi D, van der Hoek W. Justified Beliefs by Justified Arguments. In: Baral C, Giacomo GD, Eiter T, editors. *Principles of Knowledge Representation and Reasoning: Proceedings of the Fourteenth International Conference*. AAAI Press; 2014.
- [11] Shi C, Smets S, Velázquez-Quesada FR. Beliefs supported by binary arguments. *Journal of Applied Non-Classical Logics*. 2018;28:2-3:165–188.
- [12] Beirlaen M, Heyninx J, Pardo P, Straßer C. Argument strength in formal argumentation. *IfCoLog Journal of Logics and their Applications*. 2018;5(3):629–675.
- [13] D’Agostino M, Modgil S. A Rational Account of Classical Logic Argumentation for Real-World Agents. In: Kaminka G, et al., editors. *Proceedings of the Twenty-Second European Conference on Artificial Intelligence*. ECAI’16. IOS Press; 2016. p. 141–149.
- [14] Hintikka J. *Knowledge and belief: an introduction to the logic of the two notions*. Cornell University Press; 1962.
- [15] Fagin R, Halpern JY, Moses Y, Vardi M. *Reasoning about knowledge*. MIT press; 2004.
- [16] Meyer JJC, van der Hoek W. *Epistemic logic for AI and computer science*. vol. 41. Cambridge University Press; 1995.
- [17] Fagin R, Halpern JY. Belief, awareness, and limited reasoning. *Artificial intelligence*. 1987;34(1):39–76.
- [18] Blackburn P, De Rijke M, Venema Y. *Modal Logic*. Cambridge University Press; 2002.
- [19] Burrieza A, Yuste-Ginel A. Argument evaluation in multi-agent justification logics. *Logic Journal of the IGPL*. 2019;DOI:10.1093/jigpal/jzz046.
- [20] Caminada MWA, Modgil S, Oren N. Preferences and Unrestricted Rebut. In: Parsons S, Oren N, Reed C, Cerutti F, editors. *Computational Models of Argument: Proceedings of COMMA 2014*. IOS Press; 2014. p. 209–220.
- [21] Heyninx J, Straßer C. Revisiting Unrestricted Rebut and Preferences in Structured Argumentation. In: Sierra C, editor. *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI-17*; 2017. p. 1088–1092.
- [22] Baroni P, Caminada M, Giacomin M. Abstract argumentation frameworks and their semantics. In: Baroni P, Gabbay DM, Giacomin M, editors. *Handbook of formal argumentation*. London: College Publications; 2018. p. 159–236.
- [23] Caminada M. On the issue of reinstatement in argumentation. In: Fisher M, van der Hoek W, Konev B, Lisitsa A, editors. *Logics in Artificial Intelligence, JELIA 2006*. vol. 4160 of LNCS. Springer; 2006. p. 111–123.
- [24] Parikh R. Sentences, belief and logical omniscience, or what does deduction tell us? *The Review of Symbolic Logic*. 2008;1(4):459–476.
- [25] Caminada M. Rationality Postulates: applying argumentation theory for non-monotonic reasoning. *Journal of Applied Logics*. 2017;4(8):2707–2734.