# Simple models of deliberated preference

Olivier Cailloux

6th November, 2017

## 1. Motivation

We define the deliberated preference of an individual. We are interested in a context of a choice of a subset of alternatives from a set of possible alternatives given a priori. We define a model of the DP and give conditions and a validation procedure that ensure that what the model claims is in the DP is indeed. We want in particular to allow for the model to be applicable to various individuals that start with different knowledge: the model may use different counter-arguments to convince different individuals.

More in detail. A full-blown proof may be so long that nobody would read it. (That's why mathematicians do not write their proofs in set-theoretic bare notation.) Hence, shortcuts are taken, which rely on the individual knowing stuff already. More generally, assume you want to convince someone of fact $F$ with a text (we call this an argument), in our case, we want to convince a person that such alternative from the set is a good choice for her. The argument may assume the reader knows already why such or such possible counter-arguments are invalid, and the argument will thus not delve into the details of some points. When shown in interactive mode (our context here), an individual may however question some points (which we call a counter-argument). The argumenter should then be able to provide an answer to the query (a counter-counter-argument).

This article aims at defining such a procedure and conditions of its validity.

## 2. Decision situation

Here is the definition of a decision situation.

- All alternatives $\mathscr{A}$.

- Topic $T^* = \{ t_a, a \in \mathscr{A} \} \cup \{ t_{\neg a}, a \in \mathscr{A} \}$. Denoted simply $a$ and $\neg a$. We define $\neg t$, with $t = t_a$, as equal to $t_{\neg a}$ and $\neg t$, with $t = t_{\neg a}$, as equal to $t_a$.

- All possible arguments: $S^*$. It contains at least $\varnothing$, the empty argument. It contains all the arguments and more. As an example, $S^*$ could be the set of all strings (assuming all arguments can be transmitted textually).

- $s' \not\rhd_\forall^t s$: read $s'$, supporting $\neg t$ sure, never attacks $s$, supporting $t$. We show arguments $s$ and $s'$ to $i$ and ask whether $t$ is considered as a possibility by $i$ (meaning $i$ would consider picking $a$, if $t = t_a$) or rather $\neg t$ is sure (meaning $i$ would not consider picking $a$ as a choice). $s' \not\rhd_\forall^t s$ denotes a stable preference for $t$, in the sense that $s'$ does not render $s$ invalid; it is sure that $s'$ has no impact on $s$, even assuming that $s'$ would in turn resist all counter-arguments to it. This is an hypothesized relation, unobservable, used to define the deliberated preferences of $i$. Suffices that the attack occurs at least once over the considered time frame and unstability factors (such as submitting $i$ to other counter-arguments) to negate $s' \not\rhd_\forall^t s$. We do not condition on $s'$ surviving: $s'$ is declared incorrect, with no necessity of pursuing the debate and no hope of reinstatement. Example: $s'$ has already been taken into account and countered in $s$; or $s'$ does not talk about $t$ at all; or is not understood by $i$. Note that when $s' \not\rhd_\forall^t s$, further attacks to $s'$ have no chance to change that fact (assuming some properties over the way $i$ reason). The negation of this should read: $s'$ may attack $s$ in at least some cases. For example, $\neg(s' \not\rhd_\forall^t s)$ ($s'$ may attack $s$) in case $i$ suspects that $s'$ is invalid (because of some counter-argument $s_2$ to $s'$ that $i$ has in mind) but wants to leave the door open to reinstatement of $s'$.

- Define $s' \rhd_\exists^t s$ iff $\neg(s' \not\rhd_\forall^t s)$. Note that $s$ or $s'$ (or both) may equal $\varnothing$, which is useful to query the attitude of $i$ towards an argument against her own knowledge or her default attitude when no argument is given.

- Define $s \rhd_\exists^{\neg t,\text{sure}} s'$ iff $s$ may attack $s'$, where $s$ has a weak claim ($t$ is possible) and $s'$ has a strong claim ($\neg t$ is sure). Similarly, $s \not\rhd_\forall^{\neg t,\text{sure}} s'$ is a constant absence of attack. [1]

- Require (axiom A1) that $s' \rhd_\exists^t s \vee s \rhd_\exists^{\neg t,\text{sure}} s'$. Equivalently: $\neg(s \not\rhd_\forall^{\neg t,\text{sure}} s' \wedge s' \not\rhd_\forall^t s)$; $s \not\rhd_\forall^{\neg t,\text{sure}} s' \Rightarrow s' \rhd_\exists^t s$.

---

[1] Previously: $\neg(s' \rhd_\forall^t s) \Leftrightarrow s' \not\rhd_\exists^t s \Leftrightarrow s \rhd_\exists^{\neg t,\text{sure}} s' \Leftrightarrow \neg(s \not\rhd_\forall^{\neg t,\text{sure}} s')$.

- Axiom A2: $s \not\vDash_{\forall}^{\neg t, \text{sure}} s' \Rightarrow s \not\vDash_{\forall}^{\neg t} s'$. Equivalently: $s \vartriangleright_{\exists}^{\neg t} s' \Rightarrow s \vartriangleright_{\exists}^{\neg t, \text{sure}} s'$. [2] [3]

Define $T_i \subseteq T^*$ as the set of propositions that are weakly accepted.

**Definition 1** (Weak acceptance). *Define a situation $(\mathcal{A}, S^*, \{ \vartriangleright_{\exists}^{t} \})$. A proposition $t \in T^*$ is weakly accepted, $t \in T_i$, iff $\exists s \in S^* \mid \forall s' \in S^* : s' \not\vDash_{\forall}^{t} s$.*

It follows that $\neg t \notin T_i$ iff $\forall s : \exists s' \mid s' \vartriangleright_{\exists}^{\neg t} s$. Equivalently, $\neg t \notin T_i \Leftrightarrow \vartriangleright_{\exists}^{\neg t}(S^*) = S^*$.

**Definition 2** (Sure acceptance). *$\neg t \in T_i^{\text{sure}}$ iff $\exists s' \in S^* \mid \forall s \in S^* : s \not\vDash_{\forall}^{\neg t, \text{sure}} s'$.* [4]

If $\neg t \in T_i^{\text{sure}}$, then $t \notin T_i$: with $s'$ decisive for $\neg t$, sure, given any $s$, $s' \vartriangleright_{\exists}^{t} s$ (because $s \not\vDash_{\forall}^{\neg t, \text{sure}} s'$, see A1).

If $\neg t \in T_i^{\text{sure}}$, then $\neg t \in T_i$: from $s \not\vDash_{\forall}^{\neg t, \text{sure}} s'$ we obtain $s \not\vDash_{\forall}^{\neg t} s'$ using A2.

**Definition 3** (Clear-cut). *A situation is clear-cut iff $t \notin T_i \Rightarrow \neg t \in T_i^{\text{sure}}$.*

Here is an example of a non clear-cut situation. $S^* = \{\emptyset\}, T^* = \{t\}, \emptyset \vartriangleright_{\exists}^{t} \emptyset, \emptyset \vartriangleright_{\exists}^{\neg t, \text{sure}} \emptyset$.

## 2.1. Possible outcomes

We know $t \in T_i^{\text{sure}} \Rightarrow \neg t \notin T_i$ and $t \in T_i^{\text{sure}} \Rightarrow t \in T_i$.

Here are the remaining possibilities, considering the four pairs $(t, T_i)$, $(t, T_i^{\text{sure}})$, $(\neg t, T_i)$, $(\neg t, T_i^{\text{sure}})$. (A pair $(t, T)$ not mentioned means $t \notin T$.)

| | |
|---|---|
| **All poss** | $t \in T_i, \neg t \in T_i$ |
| **$t$ sure** | $t \in T_i, t \in T_i^{\text{sure}}$ |
| **$\neg t$ sure** | (symmetric) |
| **Unstability with $t$ possible** | $t \in T_i$ (oddity: $\neg t \notin T_i$ should imply $t \in T_i^{\text{sure}}$) |
| **Unstability with $\neg t$ possible** | (symmetric) |
| **Strong unstability** | (both oddities above) |

---

[2] A3: $[\exists s \mid (s' \vartriangleright_{\exists}^{t} s \wedge s \vartriangleright_{\exists}^{\neg t, \text{sure}} s')] \Rightarrow \exists s_2 \mid s_2 \vartriangleright_{\exists}^{\neg t} s'$. Equivalently: $s_2 \not\vDash_{\forall}^{\neg t} s', \forall s_2$ implies that for all $s$, $(s' \not\vDash_{\forall}^{t} s \vee s \not\vDash_{\forall}^{\neg t, \text{sure}} s')$. A1, A2 and A3 are equivalent to: $s \not\vDash_{\forall}^{\neg t, \text{sure}} s' \Leftrightarrow s \not\vDash_{\forall}^{\neg t} s' \wedge s' \vartriangleright_{\exists}^{t} s$.
[3] Having A1 and A2, can we have $s \not\vDash_{\forall}^{\neg t} s' \wedge s' \vartriangleright_{\exists}^{t} s \wedge s \vartriangleright_{\exists}^{\neg t, \text{sure}} s'$? Consider $(s', \neg t) \succeq (s, t) \sim (s', \neg t, \text{sure}) \succeq (s, t, \text{sure})$.
[4] Instead of accepting $\neg t$ for sure, it is tempting to define strong rejection of $t$ as follows. A proposition $t \in T^*$ is strongly rejected iff $\forall s_0 \in S^*, \exists s' \in S^* \mid s_0 \not\vDash_{\forall}^{\neg t, \text{sure}} s'$. But this is too weak: we want $s'$ to be also decisive, thus $s \not\vDash_{\forall}^{\neg t, \text{sure}} s', \forall s$. Hence the definition becomes $\exists s' \in S^* \mid s \not\vDash_{\forall}^{\neg t, \text{sure}} s'$.

# 3. Models

A model is a triple $(\rhd_\eta, \leadsto_\eta, +)$ defined as follows and satisfying the constraints as indicated here.

$\rhd_\eta$ an acyclic binary relation over $S^*$ (by which we mean that its transitive closure is irreflexive). $\leadsto_\eta \subseteq S^* \times T^*$. Define $S_\eta \subseteq S^*$ as the set of arguments used in $\rhd_\eta \cup \leadsto_\eta$. Let $+$ be defined over arguments used in the model: $s_3 + s_1 = s'$ for some $s' \in S_\eta$, for any $s_3, s_1 \in S_\eta$.

Requirements. The maximum length of a path in $\rhd_\eta$ is finite. $s_3 \rhd_\eta s_2 \rhd_\eta s_1 \Rightarrow s_2 \rhd_\eta s_3 + s_1.$[5]

Notation. Let $\leadsto_\eta^{-1}(T^*) \subseteq S_\eta$ denote the subset of arguments supporting propositions. $\rhd_\eta(s_2)$: arguments that $s_2$ attacks, $s_1 \in \rhd_\eta(s_2) \Leftrightarrow s_2 \rhd_\eta s_1$. We write $S \rhd_\eta s$ to mean that $\forall s' \in S : s' \rhd_\eta s$, and similarly for other binary relations.

Given a decision situation, define $\succ_\exists \subseteq \rhd_\eta$ as follows.

Given $s_3 \rhd_\eta s_2, t \in T^*$: $s_3 \succ_\exists^t s_2$ iff $[\exists s_1 \in \rhd_\eta(s_2) \mid (s_2 \succ_\exists^t s_1 \wedge s_2 \nsucc_\exists^t s_3 + s_1)] \vee [s_2 \leadsto t \wedge s_3 \rhd_\exists^t s_2].$ [6]

Given $s_3 \rhd_\eta s_2, t \in T^*$: $s_3 \nsucc_\exists^t s_2$ iff $[\exists s_1 \in \rhd_\eta(s_2) \mid (s_2 \succ_\exists^t s_1 \wedge s_2 \succ_\exists^t s_3 + s_1)] \vee [s_2 \leadsto_\eta t \wedge s_2 \rhd_\exists^{\neg t, \text{sure}} s_3].$ [7]

TODO this is well-defined because associate to each $s \in S_\eta$ $d(s)$, the distance to the farthest root (a root is an argument that $\rhd_\eta$-attacks nobody). Then $\succ_\exists$ is defined for all attacks from $d(.) = 1$ nodes (because those nodes attack only $d(.) = 0$ nodes), and thus is defined for all nodes 2, I suppose, …

Given $s_3 \in S_\eta, s_2 \in S_\eta$, with $\exists s_1 \in \rhd_\eta(s_2) \mid s_2 \succ_\exists^t s_1$, we have: $s_3 \succ_\exists^t s_2 \vee s_3 \nsucc_\exists^t s_2$.

Define $s_2 \succ_\exists s_1 \Leftrightarrow \exists t \in T^* \mid s_2 \succ_\exists^t s_1$.

Hence, given $s_3 \in S_\eta, s_2 \in S_\eta, s_2 \notin \leadsto_\eta^{-1}(T^*)$: $\neg(s_3 \nsucc_\exists s_2)$ iff $\forall s_1 \in \rhd_\eta(s_2) \cap \succ_\exists(s_2) : \neg(s_2 \succ_\exists s_3 + s_1)$.

---

[5] Necessary for definition of $s_3 \succ_\exists s_2$.

[6] Or $s_2 \leadsto t \wedge \neg t \notin \leadsto_\eta(S^*) \wedge s_3 \rhd_\exists^{t, \text{sure}} s_2$.

[7] Check: Given $s_4 \rhd_\eta s_3, s_3 \notin \leadsto_\eta^{-1}(T^*)$, with $s_2 \in \rhd_\eta(s_3) \Rightarrow (s_2$ and $s_4 + s_2 \rhd_\eta$-attack only root nodes and do not support any proposition): $s_4 \succ_\exists^t s_3$ iff

- $\exists s_2 \in \rhd_\eta(s_3) \mid [(s_3 \succ_\exists^t s_2) \wedge (s_3 \nsucc_\exists^t s_4 + s_2)]$ iff

- $\exists s_2 \in \rhd_\eta(s_3) \mid [(\exists s_1 \in \rhd_\eta(s_2) \mid s_2 \succ_\exists^t s_1 \wedge s_2 \nsucc_\exists^t s_3 + s_1) \wedge (\exists s_1 \in \rhd_\eta(s_4 + s_2) \mid s_4 + s_2 \succ_\exists^t s_1 \wedge s_4 + s_2 \succ_\exists^t s_3 + s_1)]$ iff

- $\exists s_2 \in \rhd_\eta(s_3) \mid [(\exists s_1 \in \rhd_\eta(s_2) \mid s_1 \leadsto t \wedge s_2 \rhd_\exists^t s_1 \wedge s_3 + s_1 \leadsto t \wedge s_3 + s_1 \rhd_\exists^{\neg t, \text{sure}} s_2) \wedge (\exists s_1 \in \rhd_\eta(s_4 + s_2) \mid s_1 \leadsto_\eta t \wedge s_4 + s_2 \rhd_\exists^t s_1 \wedge s_3 + s_1 \leadsto t \wedge s_4 + s_2 \rhd_\exists^t s_3 + s_1)]$.

# 4. Conditions

All these conditions assume that a decision situation $(\mathscr{A}, S^*, \{\rhd_{\exists}^t\}, \{\rhd_{\exists}^{t,\text{sure}}\})$ and a model $\eta = (\rhd_\eta, \leadsto_\eta, +)$ are given.

Define $S_{\text{decisive}} = S_\eta \setminus \text{im}(\succ_{\exists})$ the decisive arguments according to $\succ_{\exists}$, or $\succ_{\exists}$-decisive arguments for short: $s \in S_{\text{decisive}} \Leftrightarrow \succ_{\exists}^{-1}(s) = \varnothing$.

**Definition 4** (Reinstatement). *Given $s_3 \succ_{\exists} s_2 \succ_{\exists} s_1, s_3 \in S_{\text{decisive}}$: $\rhd_\eta(s_1) \subseteq \rhd_\eta(s_3 + s_1) \wedge \rhd_\eta^{-1}(s_3 + s_1) \subseteq \rhd_\eta^{-1}(s_1) \setminus \rhd_\eta(s_3)$.* [8] [9] [10]

**Definition 5** (Justifiable unstability). $\forall s_2 \rhd_\eta s_1 \mid s_2 \succ_{\exists} s_1, s_2 \nsucc_{\exists} s_1 : \exists s_3 \rhd_\eta s_2 \mid s_3 \succ_{\exists} s_2.$

**Definition 6** (Finite defense). *If $\succ_{\exists}^{-1}(s) \subseteq \succ_{\exists}(S_{\text{decisive}})$, then $\exists S \subseteq S_{\text{decisive}}, |S| \leq j \mid \rhd_\eta^{-1}(s) \subseteq \succ_{\exists}(S).$* [11]

Thus, if the attackers of $s$ are attacked by decisive arguments, then $j$ defenders are enough to defend $s$.

Define $R$, the reinstates relation, as follows: $s_3 R s_1$ iff $s_3 \succ_{\exists} s_2 \succ_{\exists} s_1$ (for some $s_2$), $s_3 \in S_{\text{decisive}}$. Define $S_\gamma$ as the transitive closure of $\leadsto_\eta^{-1}(T^*)$ under $R$.

**Definition 7** (Covering). $\forall s \in S_\gamma, s' \in S^* : s' \rhd_{\exists} s \Rightarrow s' \rhd_\eta s.$ [12]

**Definition 8** (Observable validity). $\forall s_2 \rhd_\eta s_1 \leadsto_\eta t : \neg(s_2 \succ_{\exists}^t s_1) \vee \exists s_3 \rhd_\eta s_2 \mid s_3 \succ_{\exists}^t s_2.$ *Furthermore, if $\neg(S_\eta \leadsto_\eta t), \forall s_1 \leadsto_\eta \neg t, s \in S^* : s_1 \rhd_{\exists}^t s.$* [13]

---

[8] TODO the condition must be $\succ_{\exists}^{-1}(s_3 + s_1) \subseteq \rhd_\eta^{-1}(s_1) \setminus \rhd_\eta(s_3)$ to allow $s_5 \rhd_\eta s_4 \rhd_\eta s_3 \rhd_\eta s_2 \rhd_\eta s_1$ and $s_5 \rhd_\eta s_4 \rhd_\eta s_3 + s_1$, considering that possibly $s_3$ is $\succ_{\exists}$-decisive. This should not invalidate the conditions, but it does currently. But it's not a problem: the model would actually not be built this way. In this scenario the argument $s_3 + s_1$ is useful only in case $s_3$ is decisive, thus $s_4 \rhd_\eta s_3 + s_1$ must not be planned. Rather $s_5 + s_3$ decisive, then $(s_5 + s_3) + s_1$. Alternatively, also $s_4 \rhd_\eta s_1$ and then no problem as well.

[9] The stronger condition mandating $\rhd_\eta^{-1}(s_3 + s_1) \subseteq \succ_{\exists}^{-1}(s_1) \setminus \rhd_\eta(s_3)$ would be more difficult to check: when some $s_2 \succ_{\exists} s_3 + s_1$, we'd need to check not only that $s_2 \rhd_\eta s_1$ but also $s_2 \succ_{\exists} s_1$.

[10] We do not mandate that $s_3 + s_1 \leadsto t$, so that the model can afford not resisting to the counter-attacks to $s_3 + s_1$ (resistance to c-a to $s_1$ suffice). We need: Obs applies to restricted supports (one per prop decided by model); Covering applies to extended supports (restricted supports plus those obtained by reinstatement). Replacement-1 applies to all and requires attack at least as large; Replacement-2 applies to restricted supports and requires no new $\rhd_{\exists}$-attacks.

[11] To satisfy Finite defense, in presence of the other conditions, suffice to limit the width of the model (TODO check). But it may be interesting to not limit it and declare that the model has specific replies to any counter-argument, but promises to use only a few rebuttals and that afterwards, the dm will stop using those kind of arguments (but we don't know in advance which ones will be chosen).

[12] Specify $\rhd_{\exists}$.

[13] If the model claims $\neg t \in T_i^{\text{sure}}$, this requires clear-cut (for that prop), so we must mandate it

# 5. Theorem

**Theorem 1** (Validity). *Given a decision situation and a model $\eta$, if all our conditions are satisfied, $\leadsto_\eta(S_\eta) \subseteq T_i$. Furthermore, if $\neg t \in \leadsto_\eta(S_\eta) \wedge t \notin \leadsto_\eta^{-1}(S_\eta)$, $\neg t \in T_i^{sure}$.*

*Proof.* $s$ is defended iff its $\succ_\exists$-attackers are $\succ_\exists$-attacked by $\succ_\exists$-decisive arguments.

First, we want to prove that $s_1$ defended implies $s_1$ replaceable by some $\succ_\exists$-decisive $s$, and if $s_1 \in S_\gamma$, then its replacer $s$ is in $S_\gamma$ as well.

By hypothesis, $\succ_\exists^{-1}(s_1) \subseteq \succ_\exists(S_{\text{decisive}})$. Thus, $\exists S \subseteq S_{\text{decisive}} \mid \rhd_\eta^{-1}(s_1) \subseteq \succ_\exists(S)$, $S$ finite [Finite defense].

Pick any $s_{3,1} \in S$ such that $s_{3,1} \succ_\exists s_2 \succ_\exists s_1$ (if there's none, $s_1 \in S_{\text{decisive}}$ and we're done). $s_{3,1} + s_1$ replaces $s_1$, and $\rhd_\eta^{-1}(s_{3,1} + s_1) \subseteq \rhd_\eta^{-1}(s_1) \setminus \rhd_\eta(s_{3,1})$ [Reinstatement]. Hence, $\rhd_\eta^{-1}(s_{3,1} + s_1) \subseteq \succ_\exists(S) \setminus \rhd_\eta(s_{3,1})$. Iterate by picking any $s_{3,2} \in S$ such that $s_{3,2} \succ_\exists s_2 \succ_\exists s_{3,1} + s_1$ (if there's none, $s_{3,1} + s_1 \in S_{\text{decisive}}$ and we're done) and obtaining $s_{3,2} + (s_{3,1} + s_1)$ replacing $s_{3,1} + s_1$ (hence, replacing $s_1$) with $\rhd_\eta^{-1}(s_{3,2} + (s_{3,1} + s_1)) \subseteq \rhd_\eta^{-1}(s_{3,1} + s_1) \setminus \rhd_\eta(s_{3,2})$. Hence, $\rhd_\eta^{-1}(s_{3,2} + (s_{3,1} + s_1)) \subseteq \succ_\exists(S) \setminus \rhd_\eta(s_{3,1}) \setminus \rhd_\eta(s_{3,2})$. Iterating in such a way over the finite set $S$ will finally yield an element that is $\succ_\exists$-decisive. The last point, $s_1 \in S_\gamma \Rightarrow s \in S_\gamma$, follows from the definition of $S_\gamma$.

Second, we want to prove that if $s_1$ not defended and has no decisive $\succ_\exists$-attackers (meaning that $\succ_\exists^{-1}(s_1) \subseteq \overline{S_{\text{decisive}}}$), then $s_1$ is $\succ_\exists$-attacked by some $s_2$ that is not defended and has no decisive $\succ_\exists$-attacker.

Consider $s_1$ not defended and having no decisive $\succ_\exists$-attackers. Because $s_1$ is not defended, by definition, it is $\succ_\exists$-attacked by some $s_2$ that has no decisive $\succ_\exists$-attacker. Because $s_1$ has no decisive attacker, $s_2$ is not decisive. If $s_2$ was defended, by the first part of this proof, it would be replaceable by a decisive argument, and $s_1$ would have a decisive attacker. Thus, $s_2$ is not defended.

Third, consider an argument $s_1 \in \leadsto_\eta^{-1}(T^*)$. It has no decisive $\rhd_\exists$-attacker: as $s_1 \in S_\gamma$, any $\rhd_\exists$-attack is a $\succ_\exists$-attack [Covering], and $s_1$ has no decisive $\succ_\exists$-attacker [Obs val]. Also, $s_1$ is defended: assume it is not, then by our second point in this proof some $s_2 \succ_\exists s_1$, with $s_2$ not defended and with no decisive $\succ_\exists$-attacker, and iterating and using finiteness of $\succ_\exists$ leads to a contradiction. Hence, by our first point in this proof, $s_1$ is replaceable by some $\succ_\exists$-decisive $s \in S_\gamma$. As $s \in S_\gamma$, any $\rhd_\exists$-attack is a $\succ_\exists$-attack [Covering], thus $s$ is $\rhd_\exists$-decisive. $\qquad\square$

---

(hopefully A3 or an equivalent such as Justifiable unstability fits). Thus we only need to prove $t \notin T_i$, for which $s_1 \rhd_\exists^t s$ suffices.

# A. Better definitions

Define $f_t$ from any unordered pair of arguments $(s, s')$ into { $t$ sure, $\neg t$ sure, both possible, unstable}. We present both arguments to $i$ and ask which proposition $i$ considers valid in her current state of mind, thus using both arguments and possibly other arguments she has in mind, where the first three are when $i$ always declares the same answer and the last one is used if $i$ sometimes wobble (because she has learnt new things in between two questioning, or just because of time passing, or for any known or unknown reason).

We can also use $\varnothing$ as one or both of the arguments.

We write $(s, s') \rightsquigarrow t$ (respectively $\neg t$, $b$, $\varnothing$) for $f_t(s, s') = t$ sure (respectively $\neg t$ sure, both possible, unstable).

To encode this data efficiently, we can define two binary relations $t^+$ and $t^-$: $(s, s') \in t^+ \Leftrightarrow (s, s')$ always supports $t$ sure or always supports both $\Leftrightarrow (s, s') \rightsquigarrow t \vee (s, s') \rightsquigarrow b$; and $(s, s') \in t^- \Leftrightarrow (s, s') \rightsquigarrow \neg t \vee (s, s') \rightsquigarrow b$.

This is indeed equivalent:

- $(s, s') \rightsquigarrow t \Leftrightarrow (s, s') \in t^+ \wedge (s, s') \notin t^-$

- $(s, s') \rightsquigarrow \neg t \Leftrightarrow (s, s') \notin t^+ \wedge (s, s') \in t^-$

- $(s, s') \rightsquigarrow b \Leftrightarrow (s, s') \in t^+ \wedge (s, s') \in t^-$

- $(s, s') \rightsquigarrow \varnothing \Leftrightarrow (s, s') \notin t^+ \wedge (s, s') \notin t^-$

We can now define $T_i^+, T_i^- \subseteq T^*$:

- $t \in T_i^+ \Leftrightarrow \exists s \mid \forall s' : [(s, s') \rightsquigarrow t] \vee [(s, s') \rightsquigarrow b]$.

- $t \in T_i^- \Leftrightarrow \exists s \mid \forall s' : [(s, s') \rightsquigarrow \neg t] \vee [(s, s') \rightsquigarrow b]$.

- $t$ strongly possible $\Leftrightarrow [\exists s \mid \forall s' : (s, s') \rightsquigarrow t] \vee [\exists s \mid \forall s' : (s, s') \rightsquigarrow b]$.

Here is a case where $t$ is possible but not strongly possible: $\exists (s_1, s) \rightsquigarrow t$, $\exists (s_2, s) \rightsquigarrow b$, $\nexists (s', s) \rightsquigarrow \neg t$, $\nexists (s', s) \rightsquigarrow \varnothing$. However this case is odd as just after showing $s_2$ to $i$ we could ask again about $s_1$ … (TODO)

Now do we have $t \in T_i^{\text{sure}} \Leftrightarrow \exists s \mid \forall s' : (s, s') \rightsquigarrow t$?

# B. Todo

- If $i$ does not consider $s$ as supporting $t$, it also works: if $t$ is not weakly acceptable by default, then any $s'$ is considered by $i$ as a better argument than $s$ in favor of certain $\neg t$, and so on. In fact, whether $\varnothing \rhd_\exists^t \varnothing$ determines whether $t$ is weakly supported by default.

- I should define $s'(\Box \rhd_\exists^t)s$ as an observable: "Assuming $s'$ would survive, do you consider $s'$ as leading to certainty of $\neg t$, even when considering $s$?". It distinguishes our knowledge and the truth: $s'(\Box \rhd_\exists^t)s \Rightarrow s' \rhd_\exists^t s$, thus, implies $\neg(s' \not\rhd_\forall^t s)$. But out of $\neg(s'(\Box \rhd_\exists^t)s)$, nothing.

- Partition (objectively) $S^*$ (or $S^* \times T^*$) into arguments in favor of $t$, sure, $\neg t$, sure, and similarly for possible. Use only one rel $\rhd_\exists$, defined on contradictory arguments only, instead of $\rhd_\exists^{t,\text{sure}}$ and others. Define $s' \rhd_\exists^t s$ equals no when $\neg(s' \rightsquigarrow \neg t, \text{sure})$, equals $\rhd_\exists$ for adequate arguments, and equals yes when $\neg(s \rightsquigarrow t, \text{possible})$ and $s' \rightsquigarrow \neg t$, sure, with probably some complications needed for the argument $\varnothing$ (and related default attitude towards $t$).

Questions: Q1. Relationship with $s \rhd_\exists^t s$?

We want to exclude: $s$ supports $p$ perhaps, attacked by $s2$ (supporting $\neg p$ sure), but then $s2$ is attacked by $s$. Exclude $s' \rhd_\exists^t s$ and $s \rhd_\exists^{\neg t,\text{sure}} s'$. Require to assume that this situation implies another argument $s_3$ "attacking" $s'$, thus, such that $s_3 + s$ is no more attacked by $s_2$.

# C. To think

Propositions weakly self-supported $T \subseteq T^*$: weakly accepted if no arg is given. Examples: $m$ = "eat miam"; $\neg b$ = "beurk is to exclude"; or, in a problem where there's no particularly good aliments, both $a$ = "eat this" and $\neg a$.

When given $(s, t)$, $i$ may say: $s$ does not survive; or: assuming $s$ survives, then $s$ supports $t$, or, assuming $s$ survives, then $s$ does not support $t$ anyway.

When given $s'$ against $s$, $i$ may say: $s'$ does not survive, or: assuming $s'$ survives, then $s'$ supports $\neg t$, …

Given $(s_2, t), (s_1, \neg t) \in D$, define $\neg(s_2 \succ_{\exists \neg t}^{\text{neg}} s_1)$ iff for some $(s, t) \in D$, where $s_1 \rhd_\exists^t s$: $s_1 \rhd_\exists^t s + s_2$. Equivalently: $s_2 \succ_{\exists \neg t}^{\text{neg}} s_1$ iff for all $s$, where $s_1 \rhd_\exists^t s$: $\neg(s_1 \rhd_\exists^t s + s_2)$. (This does not seem right: if given $s_3$ attacking $s_2$, and not given $s_4$ which would convincingly rebut $s_3$, then temporarily it may hold again that $s_1 \rhd_\exists^t s + s2$ (in the sense that $s_1 + s_3 \rhd_\exists^t s + s_2$).)

$s_2 \succ_{\exists \neg t} s_1$ can perhaps be queried directly by asking (in the context of some $s_1 \rhd_\exists^t s$): "assume $s_2$ survives, then does $s_2$ counter $s_1$?" (In the sense that $s_2$ is sufficiently convincing that $t$ holds perhaps, to cancel the argument $s_1$ according to which $\neg t$ surely holds.)

# D. Certainties

Looking for certainties. Those propositions that are in the reflexive preferences in a demanding sense: there is a strong enough reason to prefer it than its contrary.

- $s' >_\exists s$: weak attack; $s'$ renders $s$ invalid (can't be used to say that $t$ holds for sure) (assuming $s'$ survives)

- Propositions strongly self-supported: strongly accepted if no arg is given. Examples: $m$ = "eat miam"; $\neg b$ = "beurk is to exclude". We might have neither $c$ nor $\neg c$ in that set.

**Definition 9** (Sure acceptance)**.** *Define a situation* $(\mathscr{A}, S^*, \{\rhd_\exists^t\})$*. A proposition* $t \in T^*$ *is accepted as sure iff* $\exists s' \in S^* \mid \forall s \in S^* : s \not\rhd_\forall^{t,sure} s'$*.*

Assume we use rather: if $p$ is not sure, then $\neg p$ is weakly accepted (by def). Then we have never problems of inconsistency! But we could be in a situation where $p$ is not accepted as sure but nobody can tell why because it is fundamentally unstable (sometimes $p$ being accepted, sometimes not).

# E. Example about model instanciation

The general conditions are Reinstatement, Justifiable unstability, Finite defense and Covering. A general model is a model that claims it satisfies the general conditions.

TODO give up general models. In this example, $s_1$ would need to be planned as attacking sometimes $s_2$. Better consider an instanciation mechanism. An instantiated model is particular, and can be tested (especially against another one).

*Example* 1. $s_3 \rhd_\eta s_2 \rhd_\eta s_1 \rightsquigarrow_\eta t, s_2 \rightsquigarrow_\eta \neg t; s_3 + s_1 \rightsquigarrow_\eta t.$ $\triangle$

This model is compatible (meaning that it satisfies the general conditions) with the following decision situations. We describe $\rhd_\exists$ fully (no attack iff not mentioned).

- Sure of $t$: $s_1 \rightsquigarrow_\eta t$; $s_3 \rhd_\exists s_2$ (the rest is implied, for example $\forall s_4 \in S^* : s_1 \rhd_\exists^{\neg t} s_4$ because of covering).

- Sure of $t$ with reinstatement: $s_3 \rhd_\exists^t s_2 \rhd_\exists^t s_1 \rightsquigarrow_\eta t$; $s_1 + s_3 \rightsquigarrow_\eta t$; $s_3 \rhd_\exists^{\neg t} s_2$

- Sure of $\neg t$: $s_2 \rightsquigarrow_\eta \neg t$; $s_2 \rhd_\exists s_1$

- Both: $\neg(s_2 \rhd_\exists s_1), s_1 \rightsquigarrow_\eta t, \neg(s_3 \rhd_\exists s_2), s_2 \rightsquigarrow_\eta \neg t$

This situation falsifies the model. $s_4 \rhd_\exists s_1$, $s_4$ not attacked.

# F. Model certainties

Assume we define $s_1 \succ_\exists^{\neg t,\text{sure}} s_2 \Leftrightarrow s_2 \not\succ_\exists^t s_1$. Then, indeed, given $s_1 \rightsquigarrow_\eta t$, $s_2 \succ_\exists^{\neg t,\text{sure}} s_1$ $\Leftrightarrow s_2 \rhd_\exists^{\neg t,\text{sure}} s_1$. But it gives the wrong conclusion. For $s_3 \rhd_\eta s_2 \rhd_\eta s_1 \rightsquigarrow_\eta t : s_3 \succ_\exists^{t,\text{sure}} s_2$ iff $s_2 \not\succ_\exists^{\neg t} s_3$ iff $\exists s_1 \in \rhd_\eta(s_3) \mid s_3 \succ_\exists^{\neg t} s_1 \wedge s_3 \succ_\exists^{\neg t} s_2 + s_1$.

Given $s_1 \rightsquigarrow_\eta t$, define $s_2 \succ_\exists^{t,\text{sure}} s_1$ iff $s_2 \rhd_\exists^{t,\text{sure}} s_1$.

Given $s_1 \rightsquigarrow_\eta t$, define $s_2 \not\succ_\exists^{t,\text{sure}} s_1$ iff $s_1 \rhd_\exists^{\neg t} s_2$.

Given $s_3 \in S_\eta, s_2 \in S_\eta, s_2 \notin \rightsquigarrow_\eta^{-1}(T^*), t \in T$: $s_3 \succ_\exists^{t,\text{sure}} s_2$ iff $\exists s_1 \in \rhd_\eta(s_2) \mid s_2 \succ_\exists^{t,\text{sure}} s_1 \wedge s_2 \not\succ_\exists^{t,\text{sure}} s_3 + s_1$.

Given $s_3 \in S_\eta, s_2 \in S_\eta, s_2 \notin \rightsquigarrow_\eta^{-1}(T^*), t \in T$: $s_3 \not\succ_\exists^{t,\text{sure}} s_2$ iff $\exists s_1 \in \rhd_\eta(s_2) \mid s_2 \succ_\exists^{t,\text{sure}} s_1 \wedge s_2 \succ_\exists^{t,\text{sure}} s_3 + s_1$.

# G. Example about default arguments

s2 argues in favor of p against s1: s "le monde n'est pas fiable". s1 "le monde est fiable, bhl l'a dit". s2 "bhl est un clown, il s'est planté sur l'Irak". s3 "il avait raison sur l'Irak : l'Irak a des ADM". s4 "l'Irak n'a pas d'ADM, Bush l'a reconnu". Does s4 attack s3? "bhl est un clown, il s'est planté sur l'irak" + "l'irak n'a pas d'ADM, Bush l'a reconnu" VS "il avait raison sur l'Irak : l'Irak a des ADM" !

Measure problem?

# H. Alternative definitions of finite defense

Define Finite defense-$\succ_\exists$-$\succ_\exists$-$\rhd_\eta$-dec as: $\succ_\exists^{-1}(s) \subseteq \succ_\exists(S_\eta \setminus \text{im}(\rhd_\eta)) \Rightarrow \succ_\exists^{-1}(s) \subseteq \succ_\exists(S)$. Finite defense-$\succ_\exists$-$\succ_\exists$-$\rhd_\eta$-dec is insufficient to provide $T_\eta = T_i$. Define $s' \rhd_\eta^{\text{fail}} s$ iff $s' \rhd_\eta s \wedge \neg(s' \succ_\exists s)$. Consider $s_3 \succ_\exists s_2 \succ_\exists s_1, s_3' \succ_\exists s_2' \succ_\exists s_1$, and so on, and $s_4 \rhd_\eta^{\text{fail}} \{s_3, s_3', \dots\}$. Then I really need infinitely many arguments to defend $s_1$ but Finite defense-$\succ_\exists$-$\succ_\exists$-$\rhd_\eta$-dec is artificially satisfied because the antecedent fails to trigger.

Define Finite defense-$\succ_\exists$-$\succ_\exists$-subsets as: $\succ_\exists^{-1}(s) \subseteq \succ_\exists(S) \Rightarrow \succ_\exists^{-1}(s) \subseteq \succ_\exists(S')$. Finite defense-$\succ_\exists$-$\succ_\exists$-subsets is insufficient to provide $T_\eta = T_i$. This is because Reinstatement allows for new attacks in $\succ_\exists$ (it only forbids new attacks in $\rhd_\eta$), thus we can forever transform previously failing attacks to new attacks, hence always satisfying Finite defense (always finite cover of $\succ_\exists$, but infinite cover of $\rhd_\eta$) but still not converging. Consider $s_3 \succ_\exists s_2 \succ_\exists s_1, s_3' \succ_\exists s_2' \rhd_\eta^{\text{fail}} s_1$, and so on; and $s_3' \succ_\exists s_2' \succ_\exists s_3 + s_1, s_3'' \succ_\exists s_2'' \rhd_\eta^{\text{fail}} s_3 + s_1$, and so on.

Define Finite defense-$\rhd_\eta$->$_\exists$-startdec as: $\rhd_\eta^{-1}(s) \subseteq \,>_\exists(S_{\text{decisive}}) \Rightarrow \rhd_\eta^{-1}(s) \subseteq \,>_\exists(S)$. Finite defense-$\rhd_\eta$->$_\exists$-startdec is insufficient to provide $T_\eta = T_i$. Consider $s_3 >_\exists s_2 >_\exists s_1$, $s_3' >_\exists s_2' >_\exists s_1$, and so on, and $s_5 \rhd_\eta^{\text{fail}} s_4 \rhd_\eta^{\text{fail}} s_1$. Then I really need infinitely many arguments to defend $s_1$ but Finite defense-$\rhd_\eta$->$_\exists$-startdec is satisfied as there is no cover of the $\rhd_\eta$ attacks to $s_1$.

Define Finite defense-$\rhd_\eta$-$\rhd_\eta$-startdec as: $\rhd_\eta^{-1}(s) \subseteq \rhd_\eta(S_{\text{decisive}}) \Rightarrow \rhd_\eta^{-1}(s) \subseteq \rhd_\eta(S)$. Finite defense-$\rhd_\eta$-$\rhd_\eta$-startdec is insufficient to provide $T_\eta = T_i$. Consider $s_3 >_\exists s_2 >_\exists s_1$, $s_3' >_\exists s_2' >_\exists s_1$, and so on, and $s \rhd_\eta^{\text{fail}} \{s_2, s_2', \dots\}$. Then I really need infinitely many arguments to defend $s_1$ but Finite defense-$\rhd_\eta$-$\rhd_\eta$-startdec is artificially satisfied because of $s$. Define Finite defense-$\rhd_\eta$-$\rhd_\eta$-subsets as: $\rhd_\eta^{-1}(s) \subseteq \rhd_\eta(S) \Rightarrow \rhd_\eta^{-1}(s) \subseteq \rhd_\eta(S')$ (for any $S \subseteq S_{\text{decisive}}$). Finite defense-$\rhd_\eta$-$\rhd_\eta$-subsets is (rightly) non satisfied in this example.

Define Finite defense->$_\exists$->$_\exists$-subsets as: $>_\exists^{-1}(s) \subseteq \,>_\exists(S) \Rightarrow \,>_\exists^{-1}(s) \subseteq \,>_\exists(S')$. Finite defense-$\rhd_\eta$->$_\exists$-subsets $\not\Rightarrow$ Finite defense->$_\exists$->$_\exists$-subsets. Consider $s_2 \rhd_\eta^{\text{fail}} s_1$ (to be continued...)