

# Possibilistic Explanation

Sara Boutouhami and Aïcha Mokhtari

Institut d'Informatique, USTHB, BP 32, EL Alia, Algiers, Algeria  
e-mail: s\_boutouhami@yahoo.fr , mokhtari\_aissani@yahoo.fr

## Abstract

The philosophy literature has been struggling with the problem of defining causality. There has been extensive discussion about it. Hume taught that talk of causation was metaphysics and should be consigned to the flames (D. Hume, 1748). We present a comparison between two approaches that model this notion: The first one is the normative approach which is based on an interventionist conception of the cause (D. Kayser, A. Mokhtari, 1998). The second one is the structural-model approach that gives a new definition of actual causes, using structural equations to model counterfactuals (J. Y. Halpern, J. Pearl, 2005). Both approaches propose a modelling of the explanation using the notion of causation. Halpern and Pearl define the notion of partial explanation, so that every explanation will not be considered equally good. In order to make difference between them, the authors add a probability to the set of possible contexts (J. Y. Halpern, J. Pearl, 2005). Our purpose in this paper is to modify this definition, rather than to use a probability (quantitative modelling). We suggest to affect a degree of possibility (a more qualitative modelling) which is nearer to the human way of reasoning (P. Prade, D. Dubois, 1994). We propose to give a stratification of all possible partial explanations for a given request of an agent. The explanations in the first stratum are more possible than those belonging to the other strata.

**Keywords:** Causality, possibility, query, agent.

## 1. INTRODUCTION

Causation is a deeply intuitive and familiar relation, gripped powerfully by common sense, or so it seems. But as is typical in philosophy, deep intuitive familiarity has not led to any philosophical account of causation that is at once clean, precise, and widely agreed upon. It is safe to say that none has yet succeeded. It is also safe to say that the effort put into their development has yielded a wealth of insights into causation (N. Hall, L. A., Paul, 2004). A spring of difficulties seems to be that the notion of causality is bound to other ideas like that of explanation, of responsibility or to the problems of diagnosis where, from observations possibly imprecise and uncertain, one searches by abduction the plausible cause(s) of a given situation. This makes it's comprehension difficult (J. Y. Halpern, J. Pearl, 2005).

The paper is organized as follows. We discuss briefly in the section 2, fundamentals of the normative method of causality (A. Mokhtari, 1997). We present in section 3, the structure-based causal models, the notion of weak cause, the basic idea is

to extend the basic notion of counterfactual dependency to allow “contingent dependency” (J. Y. Halpern, J. Pearl, 2005).

The two definitions yield a plausible and elegant account of causation that handles well many of examples which have caused problems for other definitions like pre-emption and double pre-emption. They agree also with the idea that causality is not transitive. Section 4 is devoted to a comparison between the two approaches.

Thereafter, we present the definition of the partial explanation and the explanatory puissance proposed in (J. Y. Halpern, J. Pearl, 2005). The explanation proposed in the normative approach is qualitative, which is not the case in the structural approach. In order to give a qualitative character to this later, we propose to affect a degree of possibility to the definition advocated by Halpern and Pearl and then we carry out a qualitative reasoning by using the possibilistic logic based on the operation of minimum (min-based), a background of the possibilistic logic is given in this section, (P. Prade, D. Dubois, J, 1994). Finally, in section 6, we conclude and we give some perspectives of this work.

## 2. NORMATIVE METHOD OF CAUSALITY

Normative method of causality (D.Kayser, A. Mokhtari, 1998), (A. Mokhtari, 1997), extended in (M. Khelfellah, A. Mokhtari, 2001) is based on an interventionist concept of causality where an agent has the choice to perform or not an action (free will). It is based also, on the principle which stipulates that “an action may cause one or many effects”.

There are two main ideas to represent: the continuity of time and the openness of future and past. The former corresponds to the idea that things (actions, events, or others) do not happen discontinuously, the latter to the idea of non-determinism of past and future from the perspective of any particular point of time. We introduce hereafter some definitions.

Time has been explicitly defined by mean of tamed states. A *tamed state* is characterized by a subset of propositions and a date. A tamed state is a universe snapshot where the propositions of the point are true at a date of the snapshot.

Let  $T$  be the set of all tamed states, a *date* of a tamed state is defined by the following function:

$$date: T \rightarrow R$$

where  $date(t)=d$  (noted  $d_t$ ) means that  $d$  is the date tamed state  $t$ .

A *time line* is a set of tamed states which are in bijection with a set of dates. It represents a possible evolution of the universe, intended as the complete knowledge of the truth value of ‘interesting’ propositions.

A time point of this succession is supposed to respect the general principle: ‘*there is no effect without a cause*’. It is then the result of a *cause-effect* relation defined as a *causal rule*.

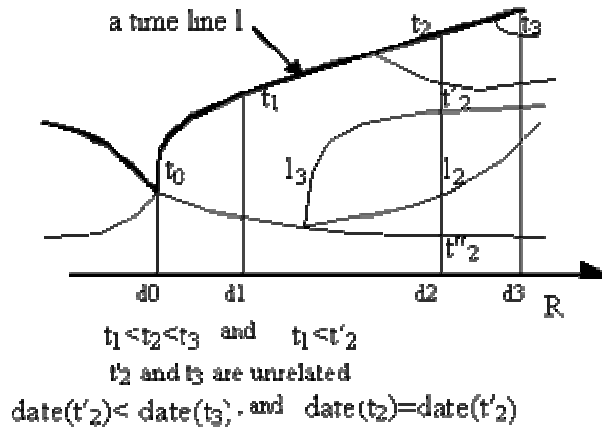


Fig 1 : Branching structure.

Let  $L$  be the set of all time lines, the tamed states of a time line are totally ordered by the precedence relation noted “ $\leq$ ”, where  $t_1 \leq t_2$  means that  $t_2$  doesn’t precede  $t_1$  whenever  $t_1 \leq t_2$  we have  $d_{t_1} \leq d_{t_2}$ .

The precedence relation expresses the principle: “no effect precedes its cause”.

The representation of the “free will” notion requires a structure of time with a branching in the future. A branching in the past is also required to examine different courses of events leading to the same situation.

Two time lines  $l_1$  and  $l_2$  coincide until a tamed state  $t$  (noted  $\text{coincide}(l_1, l_2, t)$ ) iff:

$$\forall t' \leq t, t' \in l_1 \Rightarrow t' \in l_2.$$

The set of preferred *time lines* for a tamed state  $t$  (noted  $L_p(l, d_t)$ ) is defined as a function :

$$L_p : L \times R \rightarrow 2^L$$

$$\text{Such that : } \forall l' (l' \in L_p(l, d_t) \Rightarrow \text{coincide}(l, l', t)).$$

Among the timed states, some particular ones, called *choice points*, are intended to represent states of affairs where an agent can take decision of performing or not an action. The decision of the agent is represented as a *time line* splitting into two families of futures, these two families being composed of time lines sharing the same states up to the choice point.

The proposed language is defined on two levels:

1. The first level represents *static information*. It is a plain propositional language in which:  $P$  is a set of propositions we are interested in,  $A$  is a set of actions,  $E$  is a set of effects (facts, events...), with  $A \cap E = \emptyset$  and  $P = A \cup E$ .  
A first level formula is either an action formula, or an effect formula.  
Let  $FOR(A)$  (respectively  $FOR(E)$ ) be the set of *action* (respectively *effect*) formulas.  
An *effect literal* is either an effect, or an effect negation. The set of effect literals is noted  $LIT(E)$ . It is defined as:

$$LIT(E) = E \cup \{ \neg e : e \in E \}.$$

2. The second level expresses *dynamic information* represented by formulas of the form:  $v(p, l, d)$ , which means that the formula  $p$  is true in time line  $l$  at date  $d$ .  
Causality is expressed by “normal causality” operator, noted “ $\Rightarrow$ ”.  $a \Rightarrow_e[\Delta]$  expresses that action  $a$  normally implies effect  $e$  in the delay  $\Delta$ , unless there is an occurrence of an event inhibiting the effect  $e$ . The formulation of such notion needs non-monotonic reasoning which is expressed by means of action norm and inhibiting events.

*Action norm* is defined as the set of propositions which must normally be true in order to perform the action. Formally, the norm is defined as a function:

$$Norm : A \rightarrow 2^{FOR(E)}$$

where  $norm(a)$  contains the qualifications of action  $a$  (propositions which are true unless otherwise specified when an agent considers to perform  $a$ ).

The set of the external events that can inhibit the effect of the action  $a$  is defined as a function:

$$Inhibit : LIT(E) \times A \rightarrow 2^{LIT(E)}$$

where  $e \in inhibit(e, a)$  iff  $e'$  inhibits the effect  $e$  of the action  $a$ .

The language of the normative method has certain predicates that were introduced to deal with the *ramification problem* (M.Khelfallah, A. Mokhtari, 2001):

- $occ(e, l, d)$ , with the intended meaning that the effect  $e$  has been generated in time line  $l$  at the date  $d$ . We have :  $\forall e, l, d : occ(e, l, d) \supset v(e, l, d)$ .
- $Persist(e, \delta)$ , where  $e \in LIT(E)$  and  $\delta \in R$ , means that the effect  $e$  normally persists during the delay  $\delta$ .
- Another operator called “indirect implication”, denoted “ $\dashv\rightarrow$ ”. “ $(e_1, e_2) \dashv\rightarrow e[\Delta]$ ” expresses that the occurrence of the effect  $e_1$ , in a situation in which formula  $e_2$  is true, indirectly implies the occurrence of the effect  $e$  within a delay  $\Delta$ .  
These effects are generally delayed ones, i.e., their occurrence is not necessarily instantaneous.

The normative method tackles concurrent actions without any extension. The only condition to guarantee the consistency of the obtained results is the *non existence*

of actions with opposite effects executed simultaneously in the same time line. It provides an intuitively correct analysis of the main problems encountered in the AI literature: the explanation problem, the prediction problem and the ramification problem using a simple temporal ontology (for more details see (D. Kayser, A.Mokhtari, 1998)).

### 3. STRUCTURAL APPROACH

In a recent paper (J. Y. Halpern, J. Pearl, 2005)(J. Y. Halpern, J. Pearl, 2000), Halpern and Pearl propose a definition of cause (*actual cause*) within the framework of structural causal models. Specifically, they express stories as a structural causal model (or more accurately, a causal world), and then provide a definition for when one event causes another, given this model of the story. The main idea is that a candidate C is an actual cause of an effect E when C and E have both occurred, and there exists some *counterfactual contingency* W under which E is *counterfactually dependent* on C.

Halpern and Pearl define their notion of causation within the language of structural models. Essentially, structural models are a system of equations over a set of random variables. We can divide the variables into two sets: endogenous (each of which have exactly one structural equation that determines their value) and exogenous (whose values are determined by factors outside the model, and thus have no corresponding equation). First we establish some preliminaries. We will generally use upper-case letters (e.g. X,Y) to represent random variables, and the lower-case correspondent (e.g. x, y) to represent a particular value of that variable. Dom(X) will denote the domain of a random variable X. We will use bold-face upper-case letters to represent a set of random Variables (e.g.  $\mathbf{X}, \mathbf{Y}$ ). The lower-case correspondent (e.g.  $\mathbf{x}, \mathbf{y}$ ) will represent a value assignment for the corresponding set. Formally, a signature S is a tuple  $(U, V, R)$ , where U is a set of exogenous variables, V is a set of endogenous variables, and R associates with every variable  $Y \in U \cup V$  a nonempty set  $R(Y)$  of possible values for Y (that is, the set of values over which Y ranges).

A causal model (or structural model) over signature S is a tuple  $M=(S,F)$ , where F associates with each variables  $X \in V$  a function denoted  $F_X$  such that  $F_X: (\times_{u \in U} R(U)) \times (\times_{Y \in V - \{X\}} R(Y)) \rightarrow R(X)$ .  $F_X$  determines the values of X given the values of all the other variables in  $U \cup V$ . Causal models can be depicted as a causal diagram: a directed graph whose nodes correspond to the variables in V with an edge from X to Y if  $F_Y$  depends on the value of X. Given a causal model  $M=(S,F)$ , a (possibly empty) vector  $\mathbf{X}$  of variable in V, and vectors  $\mathbf{x}$  and  $\mathbf{u}$  of values for the variables in  $\mathbf{X}$  and U, respectively, we can define a new causal model denoted  $M_{\mathbf{X} \leftarrow \mathbf{x}}$  over the signature  $S_{\mathbf{X} \leftarrow \mathbf{x}}=(U, V-\mathbf{X}, R|_{V-\mathbf{X}})$ .  $M_{\mathbf{X} \leftarrow \mathbf{x}}$  is called a submodel of M by (J. Pearl, 2000),  $R|_{V-\mathbf{X}}$  is the restriction of R to the variables in  $V-\mathbf{X}$ . Intuitively, this is the causal model that results when the variables in  $\mathbf{X}$  are set to  $\mathbf{x}$  by some external action that effects only the variables in  $\mathbf{X}$ . Formally  $M_{\mathbf{X} \leftarrow \mathbf{x}}=(S_{\mathbf{X} \leftarrow \mathbf{x}}, F^{\mathbf{X} \leftarrow \mathbf{x}})$ , where  $F_Y^{\mathbf{X} \leftarrow \mathbf{x}}$  is obtained from  $F_Y$  by setting the values of the variables in  $\mathbf{X}$  to  $\mathbf{x}$ .

#### 3.1 Syntax and Semantics:

To make the definition of actual causality precise, it is helpful to have a logic with a formal syntax. Given a signature  $S=(U, V, R)$ , a formula of the form  $X=x$ , for  $X \in V$  and

$x \in R(X)$ , is called *primitive event*. A basic *causal formula* (over  $S$ ) is one of the form  $[Y_1 \leftarrow y_1 \dots Y_k \leftarrow y_k] \phi$ , where

- $\phi$  is a Boolean combination of primitive events,
- $Y_1, \dots, Y_k$  are distinct variables in  $V$ ,
- $y_i \in R(Y_i)$ .

Such formula is abbreviated as  $[Y \leftarrow y] \phi$ . A *basic causal formula* is a Boolean combination of basic formulas. A causal formula  $\psi$  is true or false in a causal, given a context. We write  $(M, \mathbf{u}) \models \psi$  if  $\psi$  is true in the causal model  $M$  given the context  $\mathbf{u}$ .

Equipped with this background, we can now proceed to Halpern and Pearl's definition of actual cause.

### Definition 3.2

Let  $M=(U, V, F)$ , be a causal model. Let  $X \subseteq V$ ,  $X=\mathbf{x}$  is an *actual cause* of  $\phi$  if the following three conditions hold (J.Y. Halpern, J. Pearl, 2001(a)):

(AC1):  $(M, \mathbf{u}) \models X=\mathbf{x} \wedge \phi$ . (that is, both  $X=\mathbf{x}$  and  $\phi$  are true in the actual world).

(AC2): There exists a partition  $(Z, W)$  of  $V$  with  $X \subseteq V$ ,  $W \subseteq V \setminus X$  and some setting  $(\mathbf{x}', \mathbf{w}')$  of the variables in  $(X, W)$  such that if  $(M, \mathbf{u}) \models Z=z^*$  for  $Z \in Z$ , then both of the following conditions hold :

- a)  $(M, \mathbf{u}) \models [X \leftarrow \mathbf{x}', W \leftarrow \mathbf{w}'] \neg \phi$ . In worlds, changing  $(X, W)$  from  $(\mathbf{x}, \mathbf{w})$  to  $(\mathbf{x}', \mathbf{w}')$  changes  $\phi$  from true to false.
- b)  $(M, \mathbf{u}) \models [X \leftarrow \mathbf{x}, W \leftarrow \mathbf{w}', Z' \leftarrow \mathbf{z}^*] \phi$  for all subsets  $Z'$  of  $Z$ . In words, setting  $W$  to  $\mathbf{w}'$  should have no effects on  $\phi$  as long as  $X$  is kept at its current value  $\mathbf{x}$ , even if all the variables in an arbitrary subset of  $Z$  are set to their original values in the context  $\mathbf{u}$ .

(AC3):  $X$  is *minimal*; no subset of  $X$  satisfies conditions AC1 and AC2. Intuitively,  $\mathbf{x}$  is an actual cause of  $\phi$  if (AC1)  $\mathbf{x}$  and  $\phi$  are the "actual values" and (AC2) under some counterfactual contingency  $\mathbf{w}$ , the value of  $\phi$  is dependent on  $X$ , such that setting  $X$  to its actual value will ensure that  $\phi$  maintains its "actual value," even if we force all other variables in the model back to their "actual values." (AC3) is a simple minimality condition.

### 3.3. Explanation

Halpern and Pearl propose a new definition of the *explanation*, using structural equations to model counterfactuals. The definition is based on the notion of actual cause. Essentially, an explanation can be a fact that is not known to be certain but, if found to be true, would constitute an actual cause of the fact to be explained, regardless of the agent's initial uncertainty. An explanation is relative to the agent's epistemic

state, in that case, one way of describing an agent's state is by simply describing the set of the context the agent considers possible (J.Y. Halpern, J. Pearl, 2001(b)).

**Definition 3.4:**

(*Explanation*) given a structural model  $M$ ,  $\mathbf{X}=\mathbf{x}$  is an *explanation of  $\phi$  relative to a set  $K$  of contexts* if the following conditions hold (J. Y. Halpern, J. Pearl, 2005)(J. Y. Halpern, J. Pearl, 2000):

EX1:  $(M, \mathbf{u}) \models \phi$  for each  $\mathbf{u} \in K$ . (that is,  $\phi$  must hold in all contexts the agent considers possible the agent considers what he is trying to explain as an established fact).

EX2:  $\mathbf{X}=\mathbf{x}$  is a *weak cause* (without the minimal condition) of  $\phi$  in  $(M, \mathbf{u})$  for each  $\mathbf{u} \in K$  such that  $(M, \mathbf{u}) \models \mathbf{X}=\mathbf{x}$ .

EX3:  $\mathbf{X}$  is minimal; no subset of  $\mathbf{X}$  satisfies EX2.

EX4:  $(M, \mathbf{u}) \not\models \neg(\mathbf{X}=\mathbf{x})$  for some  $\mathbf{u} \in K$  and  $(M, \mathbf{u}') \models \mathbf{X}=\mathbf{x}$  for some  $\mathbf{u}' \in K$ . (This just says that the agent considers a context possible where the explanation is false, so the explanation is not known to start with, and considers a context possible where the explanation is true, so that it is not vacuous).

**Example 3.4.1:**

Suppose two arsonists lit matches in different parts of a dry forest, and both cause trees to start burning. Assume now either match by itself suffices to burn down the whole forest. We may model such a scenario in the structural-model framework as follows (see Fig 2 below):

We assume two binary background variables  $U_1$  and  $U_2$ , which determines the motivation and the state of mind of the two arsonists, where  $U_i$  is 1 iff arsonist  $i$  intends to start a fire. We then have three binary variables  $A_1$ ,  $A_2$  and  $B$  which describe the observable situation, where  $A_i$  is 1 iff arsonist  $i$  drops the match, and  $B$  is 1 iff the whole forest burns down. The causal dependencies between these variables are expressed by functions, which say that the value of  $A_i$  is given by the value of  $U_i$ , and that  $B$  is 1 iff either  $A_1$  or  $A_2$  is 1.

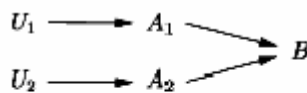


Fig 2: Causal graph

Causes and explanations for events, such as  $B = 1$  (the whole forest burns down), are defined by considering the values of variables in the above model and certain hypothetical variants. For example, arsonist 1 starting a fire is a (weak and an actual) cause of the whole forest burning down under every possible context in which arsonist 1 intends to start a fire. Moreover, arsonist 1 starting a fire is an explanation of the whole

forest burning down relative to the set of all possible contexts in which either arsonist intends to start a fire.

For more examples and extensive background on structural causal models, we refer especially to (J.Y. Halpern, J. Pearl, 2001(a)) (J.Y. Halpern, J. Pearl, 2001(b)).

## 4. COMPARISON BETWEEN THE TWO APPROACHES

The two approaches yield a plausible and elegant account of causality that handles well examples which have caused problems for other definitions like pre-emption. However, these approaches don't use the same reasoning about causality. Consequently, if we apply them to an example, we may not have the same result. We propose in this section a comparison that arises the common points (the cases which are handled by the two) and the different points.

### 4.1 Common points:

- “An explanation is relative to an agent's epistemic state”. What counts as an explanation for one agent may not count for another agent", this need is expressed in the structural approach by the set of contexts that the agent considers possible. The normative method allows also this point, by the *norm notion*. The norm is relative to the agent and to the contexts of the causal history. Example (A. Mokhtari, 1997), Serge, big smoker, is just returned in his apartment; he attempts to smoke a cigarette. To this instant an explosion destroys the apartment. It is there had a gas leak". What that will be interpreted by the insurance companies, as having caused the destructive explosion of the apartment” It is not the event «light a cigarette”, as well as it is a necessary condition to provoke the explosion. In fact, it is normal that serge smokes at his apartment. On the other hand if serge takes a walk in a petrochemical factory where it is natural to find gas in suspension in the atmosphere, and that he lights a cigarette, the cause of the explosion that will follow will be attributed to the event “light a cigarette”. This event is not here a legitimate description of the natural course of the things.
- “The relation of causality is not transitive”, the two definitions agree on this property. Halpern and Pearl give an example (see (J.Y. Halpern, J. Pearl, 2005)) to show that causality is not transitive, according to their definitions. For the normative approach, the transitivity is excluded in the definition of the relation of causality: “actions are the only elements that can change world state”.
- The notion of action is well formulated in the normative approach by the notion of choice point. The two approaches agree on the idea that only action can changes the state of world. In the structural approach, only endogenous variables can be part of a cause.
- “The *pre-emption*”: the two approaches deal with the pre-emption (where there are two potential causes of an event, one of which preempts the other). In (J.Y. Halpern, J. Pearl, 2005) a story is taken from (N. Hall, 2004) and the author give an adequate modelisation of the example. The normative approach can handle the pre-



emption cases using the notion of norm : when the first action is realised, it violates the norm (pre-condition) of the second action (the effect of an action must be true).

- “The choice of the type of variables is difficult”: as in any approach, there is no formalism that enables to choose the types of the variables. This is a task involving experts in the domain and some heuristics to make appropriate choices. In (A. Mokhtari, 1997) the partition of P into actions A and events E, it is not always obvious to decide precisely what proposition falls into what category. The same case for (J.Y. Halpern, J. Pearl, 2005) it is not easy to make the choice for a given variable if she is endogenous or exogenous see (J.Y. Halpern, J. Pearl, 2005) show the importance of the choice of variable (adding an endogenous variable corresponding to an exogenous variable can result in there being an explanation when there was none before).
- The two methods can not deal with the disjunctive causes. In (J.Y. Halpern, J. Pearl, 2005) the only reasonable definition of “ $A \vee B$  causes C” seems to be that “A causes C or B causes C”. In (A. Mokhtari, 1997) the problem is the disjunction of norm.
- “Inexplicable events are possible”; the two approaches allow inexplicable events.
- “The causality by production is transitive”: In (J.Y. Halpern, J. Pearl, 2005), the condition AC2(b) captures some of the features of production (A forced B to happen, even if  $W=w$ ). For the normative approach, this is done by applying the function closure for ramification and by exploiting the static and causal rules (only action can modify the world state) (for more details see (M. Khelfellah, A. Mokhtari. 2001)).
- The two approaches use a non-monotonic reasoning. In (A. Mokhtari, 1997) the possibility that an unexpected event prevents  $e$  from remaining true, the notion of persistence and the notion of norm. In (J.Y. Halpern, J. Pearl, 2005) is that a cause is relative to a given context in a given causal model, which means that the cause change from one context to another one.

## 4. 2. Difference

- The time plays a crucial role in the perception of the causality, and yet it has no explicit representation in structural models. (M. Hopkins, J. Pearl, 2003). While in the normative approach the time is represented explicitly as a parameter in the predicates. And when we look for the cause or the explanation, the search is based on a base of facts and effects ordered on their realisation time.
- The normative approach allows the representation of the temporal aspect with the capacity to handle:

- The duration, the realisation of an action takes time,
  - The delay, i. e., effects are not always immediate,
  - The effects of concurrent actions (the only condition is the non existence of actions with opposite effects executed simultaneously).
- The problem with using structural causal models is that the language of structural models is simply not expressive enough to capture certain intricate relationships that are important in causal reasoning. Within the structural models framework, we are only dealing with values assignments to random variables and hence cannot express the following, among others (M. Hopkins, J. Pearl, 2003):
    - A distinction between a condition and a transition,
    - A distinction between the presence and the absence of an event,
    - Time and temporal relationships.

Halpern and Pearl propose a sophisticated definition for actual causality based on structural causal models, however although this definition works on many previously problematic examples, it still does not fit with intuition on all examples. Some of these difficulties can be traced to the limited expressiveness of the structural model formulation (M. Hopkins, J. Pearl, 2003). The explanation proposed in this approach unlike the normative approach is not qualitative. To handle this problem, we propose an improvement of this definition in the next section.

## **5. POSSIBILISTIC EXPLANATION**

The explanation proposed in the latter section is not qualitative, in order to give a qualitative character to this later, we propose to affect a degree of possibility ( a more qualitative modelling) which is nearer to the human reasoning (H. Prad, D. Dubois, 1988). We propose a new definition of explanation using the possibilistic logic, as a tool for ordering the set of possible explanations. The agent's epistemic state will be represented by describing the set of the interpretations that the agent considers possible. We propose a stratification of all possible partial explanations, this stratification reflects a hierarchy of priority between partial explanations.

We give a background on the possibilistic logic, for more details see (L.A.Zadah, 1978) (D.Dubois, J.Lang, H.Prade, 1994).

### **5.1 Possibilistic logic**

The possibilistic logic based on possibility theory (L.A.Zadah, 1978) offers an uncertainty modelling framework for dealing with possibility distributions representing fuzzy information. It is for it, agent believes and the preferences are expressed by means of possibilistic logic.

Let  $L$  a finite propositional language and  $\Omega$  be the set of all propositional interpretations. Let  $\varphi, \Psi, \dots$  be propositional formulas.  $\omega \models \varphi$  means that  $\omega$  is a model of  $\varphi$ .

A possibility distribution  $\pi$  is a mapping from a set of interpretation  $\Omega$  into the unit interval  $[0, 1]$ .  $\pi(\omega)$  represents the degree of compatibility of the interpretation  $\omega$  with available pieces of information (D.Dubois, J.Lang, H.Prade, 1994).  
By convention:

- ✓  $\pi(\omega) = 0$  means that  $\omega$  is impossible to be the real world,
- ✓  $\pi(\omega) = 1$  means that  $\omega$  is totally possible to be the real world,
- ✓  $\pi(\omega) > \pi(\omega')$  means that  $\omega$  is preferred candidate to  $\omega'$  for being the real world.

A possibility distribution  $\pi$  is said to be normalized if there exists an interpretation  $\omega$  such that  $\pi(\omega) = 1$ . Given a possibility distribution  $\pi$ , two dual measures are defined :

- ✓ The possibility measure of a formula  $\varphi$ , defined by :

$$\Pi(\varphi) = \max \{ \pi(\omega) : \omega \models \varphi \text{ and } \omega \in \Omega \}$$

which evaluates the extent to which  $\varphi$  is consistent with the available beliefs expressed by  $\pi$ .

- ✓ The necessity measure of a formula  $\varphi$ , defined by :

$$N(\varphi) = 1 - \Pi(\neg\varphi) = \min \{ 1 - \pi(\omega) : \omega \models \neg\varphi \}$$

which evaluates the extent to which  $\varphi$  is entailed by the available beliefs.

A possibilistic knowledge base  $\Sigma$  is a set of weighted formulas :

$$\Sigma = \{ (\varphi_i, \alpha_i) : i = 1, \dots, n \}$$

where  $\varphi_i$  is a propositional formula and  $\alpha_i \in [0, 1]$  which represents the certainty level of  $\varphi_i$ .

Each piece of information  $(\varphi_i, \alpha_i)$  of possibilistic knowledge base can be viewed as a constraint which restricts the possibility degrees of possible interpretations (D.Dubois, J.Lang, H.Prade, 1994).

If an interpretation  $\omega$  satisfies  $\varphi_i$  then its possibility degree is equal to 1, otherwise it is equal to  $1 - \alpha_i$ .

More formally, the possibility distribution associated with a weights formula  $(\varphi_i, \alpha_i)$  is  $\forall \omega \in \Omega [1]$  :

$$\pi_{(\varphi_i, \alpha_i)}(\omega) = \begin{cases} 1 - \alpha_i & \text{if } \omega \models \neg\varphi_i \\ 1 & \text{otherwise} \end{cases}$$

More generally, the possibility distribution associated with a possibilistic knowledge base  $\Sigma$  is the result of combining possibility distributions associated with each weighted formula  $(\varphi_i, \alpha_i)$  of  $\Sigma$ , namely  $\forall \omega \in \Omega$  (D.Dubois, J.Lang, H.Prade, 1994) :

$$\pi_{\Sigma}(\omega) = \square \{ \pi_{(\varphi_i, \alpha_i)}(\omega) \mid (\varphi_i, \alpha_i) \in \Sigma \}$$

where  $\square$  is in general either equal to the minimum (min) operator (in standard possibilistic logic), or the product operator (\*).

In the rest of the paper, we only focus on the case where  $\square = \min$ . The possibilistic base  $\Sigma$  is then called qualitative possibilistic knowledge base.

### Explanation definition 5.2

Let  $\pi$  be a distribution of possibility i.e., a mapping from a set of interpretations  $\Omega$  that the agent considers possible into the interval  $[0,1]$ . Let  $\omega$  be an interpretation that the agent considers possible ( $\omega \in \Omega$ ). Given a structural model  $M$ ,  $\mathbf{X}=\mathbf{x}$  is an *explanation of  $\varphi$  relative to a set  $\Omega$  of possible interpretations* if the following conditions hold:

EX1':  $(M, \omega) \models \varphi$  for each  $\omega \in \Omega$ . (that is,  $\varphi$  must be satisfied in all interpretation the agent considers possible).

EX2':  $\mathbf{X}=\mathbf{x}$  is a *weak cause* of  $\varphi$  in  $(M, \omega)$  for each  $\omega \in \Omega$  such that  $(M, \omega) \models \mathbf{X}=\mathbf{x}$ .

EX3':  $\mathbf{X}$  is minimal; no subset of  $\mathbf{X}$  satisfies EX2.

EX4':  $(M, \omega) \models \neg(\mathbf{X}=\mathbf{x})$  for some  $\omega \in \Omega$  and  $(M, \omega') \models \mathbf{X}=\mathbf{x}$  for some  $\omega' \in \omega$ .

Not all explanations are considered equally good. Some explanations are most plausible than others. We propose to define the goodness of an explanation by introducing a degree of possibility (include priority levels between explanations).

The measure of possibility of an explanation is given by:  $\Pi(\mathbf{X}=\mathbf{x}) = \max \{ \pi(\omega) \mid \omega \models \mathbf{X}=\mathbf{x}, \omega \in \Omega \}$ .

#### Example 5.2.1:

Suppose I see that Victoria is tanned and I seek an explanation. Suppose that the causal model includes variables for “Victoria took a vacation”, “It is sunny in the Canary Islands”, “Victoria went to a tanning”. The set of  $\Omega$  includes interpretations for all settings of these variables compatible with Victoria being tanned. Note that, in particular, there is an interpretation where Victoria both went to the Canaries (and didn't get tanned there, since it wasn't sunny) and to a tanning salon. Victoria taking a vacation is not an explanation (relative to  $\Omega$ ), since there is an interpretation where Victoria went to the Canary Islands but it was not sunny, and the actual cause of her tan is the tanning salon, not the vacation. However, intuitively it is “almost” satisfied, since

it is satisfied (by every interpretation in  $\Omega$ ), in which Victoria goes to the Canaries. "Victoria went to the Canary Islands" is *partial* explanation of "Victoria being tanned". There is a situation where we can't find a complete explanation (it is inexplicable). In the next section we give our definition of the goodness of a partial explanation.

### Definition 5.3

Let  $\Omega_{X=x,\phi}$  be the largest subset  $\Omega'$  of  $\Omega$  such that  $X=x$  is an explanation of  $\phi$  relative to  $\Omega_{X=x,\phi}$  (it consists of all interpretations in  $\Omega$  except the those where  $X=x$  is true but is not a weak cause of  $\phi$ ) (S. Boutouhami, A. Mokhtari, 2005).

$\Omega' = \Omega - \{ \omega : \omega \in \Omega \mid \omega \models X=x, \omega \models \phi \text{ and } X=x \text{ is not a weak cause of } \phi \}$ .

- $X=x$  is a *partial explanation* of  $\phi$  with the *goodness*  $\Pi(\Omega_{X=x,\phi} \mid X=x) = \max \{ \pi(\omega) : \omega \models X=x, \omega \in \Omega' \}$ .
- $X=x$  is a  $\alpha$ -*partial explanation* of  $\phi$  relative to  $\pi$  and  $\Omega$ , if  $\Omega'$  exists and  $\Pi(\Omega_{X=x,\phi} \mid X=x) \geq \alpha$ .
- $X=x$  is an *partial explanation* of  $\phi$  relatively to  $\pi$  and  $\Omega$  iff  $X=x$  is a  $\alpha$ -*partial explanation* of  $\phi$ , and  $\alpha \geq 0$ .

Partial explanations will be ordered, in a set of strates  $S_{\alpha_1} \cup \dots \cup S_{\alpha_n}$  for a given request (see fig 3).

- In the strate  $S_{\alpha_1}$  we will have complete explanations if there exists,
- $X=x$  is in the strate  $S_{\alpha_i}$  if  $\Pi(\Omega_{X=x,\phi} \mid X=x) = \alpha_i$ ,
- Let  $X=x$  be a partial explanation in the strate  $S_{\alpha_i}$ , and  $Y=y$  a partial explanation in the strate  $S_{\alpha_j}$ .  $X=x$  is a partial explanation more plausible than the partial explanation  $Y=y$  if  $\Pi(\Omega_{X=x,\phi} \mid X=x) = \alpha_i > \Pi(\Omega_{Y=y,\phi} \mid Y=y) = \alpha_j$ .

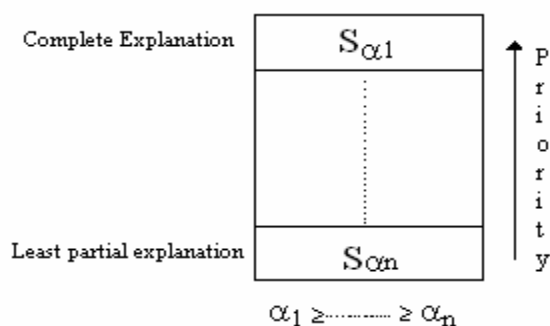


Fig 3 : A stratification of partial explanations

**Example 5.3.1:**

Suppose that two arsonists drop lit matches in different parts of a dry forest, and both cause trees to start burning. Consider the Disjunctive scenario; either match by itself suffices to burn down the whole forest. There are two possible explanations in the disjunctive scenario: arsonist 1 (A1) did it or arsonist 2 (A2) did it. Suppose that there is a variable O corresponding to the presence of oxygen. Suppose that there is an interpretation in which there is no oxygen (very unlikely). In that case, the presence of oxygen becomes a partial explanation of the fire. Nevertheless, it is an explanation with very little *explanatory power*.

Let  $\pi$  a distribution of possibility on the set  $\Omega$  of interpretations (includes  $\Omega$ ).  $\pi$  represents the agent's prior possibility before the explanandum  $\phi$  is observed or discovered. Thus,  $\pi$  is the conditioning  $\pi$  on  $\phi$  and  $\Omega$  consists of the interpretation in  $\Omega$  that satisfy  $\phi$ .  $\pi = \pi(\cdot | \phi)$ .

The usual definition of a conditional distribution of possibility is:

$$\pi(\omega | \phi) = \begin{cases} 1 & \text{if } \Pi(\phi) = \pi(\omega), \\ \pi(\omega) & \text{if } \pi(\omega) < \Pi(\phi) \text{ and } \omega \models \phi, \\ 0 & \text{else.} \end{cases}$$

Conditioner with  $\phi$  consists on a revision of degree of possibility associated to different interpretations, after having the certain information  $\phi$ . ( $\phi$  is a certain information, so interpretation that falsifies  $\phi$ , are impossible) (S. Benferhat, S. Lagrue, O. Papipi, 2003).

We propose the measure of *explanatory power* of  $\mathbf{X}=\mathbf{x}$  to be  $\Pi(\Omega_{\mathbf{X}=\mathbf{x}, \phi} | \mathbf{X}=\mathbf{x}) = \max \{ \pi(\omega) : \omega \models \mathbf{X}=\mathbf{x}, \omega \in \Omega' \}$ .

In the precedent example, the presence of oxygen is an explanation with a very low explanatory power, while the arsonist 1 has a high explanatory power.

**5.4 Algorithm of Generation of Strata**

The main idea of our algorithm is to provide a set of choice of ordred partial explanations for a given request of the agent.

Let  $\phi$  be request for which the agent seeks an explanation, let  $V$  be the set of endogenous variable, let  $X \subseteq V / Y_i, Y_i \in \phi$  the set of possible variable my formulate the explanation. For all subset  $X'$  of  $X$ , decide if there exist an attribution of values wich make it a partial explanantion, if is the case then compute  $\Pi(\Omega_{X'=X', \phi} | X'=X')$ . Ones that done, seek if there exist a strate wich will belong to, then add it to this latter ; If there isn't, create a new strate wich will have as an element this partial explanation, finally insert the new strate in the aproprate order with the existent strates. At the end of that algorithm we will have all partial explanations, this strusture facilitate the task of searching an explanation when we have a new consideration of the agent, an adaptation with the evolution of the agent beliefs.

S : the set of all partial explanation;  
V : the set of endogenous variables.  
 $\varphi$  : the fact to be explained.  
Dom(X): the set of possible attributions of X.

**Input** : {S= $\emptyset$ , V,  $\varphi$ ,  $\Omega$ , Dom(X)}

**begin**

X=V- $\{Y_j\} / \forall Y_j$ ,  $Y_j$  is a variable in  $\varphi$ .

**for all**  $X' \subseteq X$  **do**

**begin**

- Decide if there exist  $x' \in \text{Dom}(X)$ , such that  $X'=x'$  is a  $\alpha$ -partial explanation of  $\varphi$  relative to  $\Omega$ .

- **if**  $X'=x'$  a  $\alpha$ -partial explanation **then**

**begin**

Compute  $\prod(\Omega_{X'=x', \varphi} | X'=x')$  ; Let  $\alpha_i = \prod(\Omega_{X'=x', \varphi} | X'=x')$

**if** the strate  $S_{\alpha_i}$  exists **then**

Add  $\{X'=x'\}$  to the strate  $S_{\alpha_i}$ .

**else**

**begin**

Create a new strate  $S_{\alpha_i}$

Add  $\{X'=x'\}$  to the strate  $S_{\alpha_i}$ .

Insert the strate  $S_{\alpha_i}$  in the good order.

$S = S \cup S_{\alpha_i}$ .

**end.**

**end.**

**end.**

**end.**

**Out put**  $\{S = \cup S_{\alpha_i}\}$ .

The problem is to compute the set of all Subset of X, wick are partial explanation of  $\varphi$ , the problem can be reduced to the problem of computing the set of all partial explanation, computing the set of all valid formulas among a *Quantified Boolean formula (QBF)* =  $\exists B \forall C \exists D f(B, C, D)$ , where " $\exists B \forall C \exists D y$ " is a reduction of guessing some  $X' \subseteq X$  and  $x' \in \text{Dom}(x)$ , and deciding whether  $X'=x'$  is  $\alpha$ -partial explanation.

## 6. CONCLUSIONS

In this paper, we have given a comparison between two approaches that model the notion of causality. We have proposed the use of the possibilistic logic which provides a priority level between the explanations; we prefer the use of possibility instead of probability because on one hand, possibility reflects better the human reasoning, which is rather qualitative than quantitative.

We have proposed a stratification of all partial explanations for a given request. This stratification facilitates the task of searching a new explanation when we have a new consideration of the agent (an evolution of the agent beliefs).

As an immediate extension to this work, we analyse of the complexity of the reasoning process to find the stratification of the all partial possibilistic explanations

## 7. REFERENCES

- [1] S. Benferhat, S. Lagrue, O. Papini, 2003. A possibilistic handling of partially ordered information, UAI 2003, pp 29-36.
- [2] S. Boutouhami, A. Mokhtari, 2005. Possibilistic explanation, Proceeding of the Information & Communication Technologies International Symposium ICTI'S 05, Tetuan Morocco, pp 7-13, June 2005.
- [3] N. Hall, L. A., Paul, Causation and counterfactuals, Edited by John Collins, Ned Hall and L. A. Paul.,(Eds) Cloth/June 2004.
- [4] J.Y. Halpern, J. Pearl, Causes and Explanations: "A Structural-model Approach". To appear in British Journal for Philosophy of Science, 2005.
- [5] J.Y. Halpern, J. Pearl, Causes and explanations: A structural-model approach, Part I: Causes, in: Proceedings of the Seventeenth Conference in Uncertainty in Artificial Intelligence (UAI-01), Seattle, WA, pp. 194–202, 2001.
- [6] J.Y. Halpern, J. Pearl, Causes and explanations: A structural-model approach, Part II: Explanations, in: Proceedings of IJCAI-2001, pp. 27–34, Seattle, WA, 2001.
- [7] J.Y. Halpern, J. Pearl, Causes and explanations: A structural-model approach, Technical Report R-266, UCLA Cognitive Systems Lab., 2000.
- [8] M. Hopkins, J. Pearl, Clarifying the Usage of structural models for Commonsense Causal Reasoning. In Proceedings of the AAAI Spring Symposium on Logical Formalizations of Commonsense Reasoning, 2003.
- [9] Hume, D.,1748, An Enquiry concerning Human Understanding. Reprinted Open Court Press, LaSalle, IL, 1958.
- [10] D. Kayser, A. Mokhtari, Time in a causal theory. AMAI. Journal, Vol 22, pp 117-138, 1998.
- [11] M. Khelfallah, A. Mokhtari, Ramification in the Normative Method of Causality, ECSQARUS.,pp 704-713. In LNCS n° ,Toulouse, France, September 2001
- [12] A. Mokhtari, Action-Based Causal Reasoning, Applied Intelligence, The International Journal of Artificial Intelligence, Neural Networks, and Complexes Problem-Solving Technologies. Volume 7, Number 2, , pp 99-112 April 1997.
- [13] J. Pearl. Causality Models, Reasoning, and Inference. New York: Cambridge University Press, 2000.
- [14] H. Prade, D. Dubois, J. Lang, Possibilistic logic, In Handbook of Logic in Artificial Intelligence and Logic Programming, (D Gabbay, Oxford University Press, pp 439-513) , 1994.
- [15] H. Prade, D. Dubois. Possibility theory: An approach to computerized, processing of Uncertainty. Plenum Press, New York, 1988.
- [16] D. Dubois, J. Lang, H. Prade – Possibilistic logic – In handbook of logic in Artificial Intelligence and logic programming, (D. Gabbay et al., eds, 3, Oxford University Press : pages 439-513, 1994).



- [17] L. A. Zadeh- Fuzzy sets as a basis for a theory of possibility – Fuzzy Sets and Systems, 1, 3-28, 1978.