

2018-12

## A Qualitative Investigation of the Degree of Explainability of Defeasible Argumentation and Non-monotonic Fuzzy Reasoning

Lucas Rizzo  
lucas.rizzo@tudublin.ie

Luca Longo  
*Technological University Dublin*, luca.longo@tudublin.ie

Follow this and additional works at: <https://arrow.tudublin.ie/scschcomcon>

---

### Recommended Citation

Rizzo, L. & Longo, L. (2018). A Qualitative Investigation of the Degree of Explainability of Defeasible Argumentation and Non-monotonic Fuzzy Reasoning. *26th AIAI Irish Conference on Artificial Intelligence and Cognitive Science*. pp. 138-149. doi:10.21427/tby8-8z04

This Conference Paper is brought to you for free and open access by the School of Computing at ARROW@TU Dublin. It has been accepted for inclusion in Conference papers by an authorized administrator of ARROW@TU Dublin. For more information, please contact [yvonne.desmond@tudublin.ie](mailto:yvonne.desmond@tudublin.ie), [arrow.admin@tudublin.ie](mailto:arrow.admin@tudublin.ie), [brian.widdis@tudublin.ie](mailto:brian.widdis@tudublin.ie).



This work is licensed under a [Creative Commons Attribution-Noncommercial-Share Alike 3.0 License](https://creativecommons.org/licenses/by-nc-sa/3.0/)

# A Qualitative Investigation of the Degree of Explainability of Defeasible Argumentation and Non-monotonic Fuzzy Reasoning

Lucas Rizzo and Luca Longo\*

The ADAPT global centre of excellence for digital content and media innovation  
School of Computing, Dublin Institute of Technology, Dublin, Ireland  
lucas.rizzo@mydit.ie, luca.longo@dit.ie\*

**Abstract.** Defeasible argumentation has advanced as a solid theoretical research discipline for inference under uncertainty. Scholars have predominantly focused on the construction of argument-based models for demonstrating non-monotonic reasoning adopting the notions of arguments and conflicts. However, they have marginally attempted to examine the degree of explainability that this approach can offer to explain inferences to humans in real-world applications. Model explanations are extremely important in areas such as medical diagnosis because they can increase human trustworthiness towards automatic inferences. In this research, the inferential processes of defeasible argumentation and non-monotonic fuzzy reasoning are meticulously described, exploited and qualitatively compared. A number of properties have been selected for such a comparison including understandability, simulatability, algorithmic transparency, post-hoc interpretability, computational complexity and extensibility. Findings show how defeasible argumentation can lead to the construction of inferential non-monotonic models with a higher degree of explainability compared to those built with fuzzy reasoning.

**Keywords:** Defeasible Argumentation, Non-monotonic Reasoning, Fuzzy Reasoning, Argumentation Theory, Explainable Artificial Intelligence

## 1 Introduction

Knowledge-driven approaches have been extensively used in the field of Artificial Intelligence (AI) for producing inferential models of reasoning. Among them, fuzzy reasoning [21] and defeasible argumentation [4] possess a higher explanatory capacity when compared to other reasoning approaches for dealing with partial, vague and conflicting information [2, 20]. This is because, intuitively, the inferences that can be produced by these approaches can be better understood by humans, due to the fact that they deal and manipulate knowledge provided by experts preserving their natural language. However, to the best of our knowledge, no empirical investigation of their explanatory capacity has been made so far. Model explainability is essential for its adoption and usage. The lower the model explanatory capacity, the lower the degree of trust posed by humans

towards their inferences. Medical diagnosis and autonomous driving are examples of application areas where this often occur. In these areas, humans need to fully understand model functioning in order to trust its inferences. In the field of Artificial Intelligence a number of properties have been proposed for evaluating the degree of explainability of inferential models. Some of these include model extensibility [11], its simulatability and its post-hoc interpretability [12]. The aim of this research is to qualitatively analyse the explanatory capacity of non-monotonic fuzzy reasoning and defeasible argumentation. A detailed step-by-step description of their inferential mechanisms is described and contrasted according to a selection of properties from the literature. Both these inferential mechanisms are exploited by adopting a knowledge-base provided by an expert in the field of biomarkers. This knowledge-base is composed by a set of rules which are brought together and evaluated to predict the mortality risk of elderly individuals. In detail, the research question investigated is: *“How do the explanatory capacity provided by defeasible argumentation and non-monotonic fuzzy reasoning relate qualitatively?”*

The remainder of this paper is organised as it follows: Section 2 firstly outlines defeasible argumentation and non-monotonic fuzzy reasoning. Secondly, it introduces related work on Explainable Artificial Intelligence (XAI) presenting a number of properties useful for assessing model explainability. The design of a comparative research study and the inferential processes of defeasible argumentation and non-monotonic fuzzy reasoning are detailed in Section 3. Section 4 provides a qualitative comparison of the selected properties followed by a discussion, while Section 5 concludes the research study.

## 2 Related work

Defeasible (non-monotonic) reasoning has emerged as a solid theoretical approach within AI for modeling non-monotonic activities under fragmented, ambiguous and conflicting knowledge. In a non-monotonic reasoning process, conclusions do not necessarily increase monotonically, but instead they can be withdrawn as new information arises [14]. A particular type of defeasible reasoning is argumentation, built upon the notions of arguments and their conflicts [13, 2]. Defeasible argumentation provides the basis for the development of computational models of arguments. Such development starts with the definition of the internal structure of arguments to the resolution of their conflicts and final accrual towards a rational conclusion.

Another type of non-monotonic reasoning can be achieved by employing fuzzy logic and reasoning. This allows the creation of computational models with a robust representation of linguistic information provided by domain experts by employing the notion of degree of truth. Fuzzy reasoning consists of a fuzzification module, responsible for assigning to each proposition or linguistic fuzzy term, provided by an expert, a degree of truth; an inference engine accountable for firing rules and aggregating fuzzy terms; and a defuzzification module, which translates this aggregation using the original natural language employed in the

underlying reasoning [15]. This robustness to deal with vagueness of information have led to 50 years of research endeavour, with a plethora of applications in many domains. However, in order to deal with non-monotonic information, the classical fuzzification-engine-defuzzification process has to be extended with a non-monotonic layer. Unfortunately not many research studies exist for this purpose. For example, in [6], an average function is proposed for aggregating conclusions from conflicting rules, while in [10] a reduction of non-monotonic rules is suggested by means of a rule base compression method. In this study, the approach proposed in [20] is selected. It employs the use of Possibility Theory [7] as a way of dealing with conflicting rules. In a nutshell, truth values are represented by the notions of *possibility* and *necessity*. These indicate respectively the extent to which data fail to refute its truth and the extent it supports its truth.

Previous studies have attempted to analyse the inferential capacity of defeasible argumentation in the context of other approaches of quantitative reasoning under uncertainty [17–19]. However, so far, such analysis has been brought forward only by means of predictive accuracy. It has been demonstrated that the evaluation of predictive accuracy alone might not be sufficient for a model to be employed and trusted by domain experts. For instance, in [5] a model was trained to predict the probability of death from pneumonia and inferred less risk to patients who also had asthma. However, asthma is, in fact, a predictor of higher risk of death. The inference reflected a pattern of lower risk in the training data as a consequence of the more intrusive treatment received by asthmatic patients. Hence, if we expect defeasible argumentation to be trusted and understood by domain experts it is also necessary to situate its explanatory capacity in relation to other similar reasoning approaches. The literature on Explainable Artificial Intelligence is vast and it contains several properties for explainability analysis [11, 1, 12]. Six of these were selected and considered relevant to the knowledge-driven approaches under scrutiny. Some of them were initially defined in the machine learning context, but we believe they can be borrowed for the analysis of reasoning approaches. Table 1 lists their definitions.

Table 1: Properties for explainability, their definitions and sources.

| Property   | Definition  | Source   |
|--|---|----------|
| Understandability/<br>Post-hoc<br>Interpretability | Capacity of understanding the inferential process behind a model in order to trust and adopt it as a decision supporting tool / Capacity of extracting information from a constructed model and the degree of elucidation of its inferences | [1]/[12] |
| Simulatability                                     | Capacity of a human to step through every calculation required to produce a prediction in a reasonable time by employing input and parameters   | [12]     |
| Extendibility                                      | The easiness of an inferential system to accommodate new input parameters and new output classes.   | [11]     |
| Computational<br>Complexity                        | Complexity of the algorithms employed in the inferential process (computational time needed to produce an inference)  | [11]     |
| Algorithmic<br>transparency                        | Degree of application of the inferential process to new domains   | [12]     |

### 3 Design and methodology

In order to investigate the explanatory capacity provided by defeasible argumentation and non-monotonic fuzzy reasoning, a knowledge-base was selected and operationalized employing two mechanisms for non-monotonic reasoning: defeasible argumentation and non-monotonic fuzzy reasoning. This knowledge-base was produced by a clinician. The reasoning models built upon it aimed at predicting the risk of mortality in elderly individuals by using information related to their biomarkers. The first inferential approach, defeasible argumentation, is structured over 5 layers as in [13]: 1) definition of the structure of arguments, 2) and their conflicts, 3) their evaluation 4) the definition of their dialectical status, 5) their final accrual. The second approach, non-monotonic fuzzy reasoning, is composed of three main parts: 1) a fuzzification module, 2) an inference engine and 3) a defuzzification module. Fig. 1 summarises the design of the research.

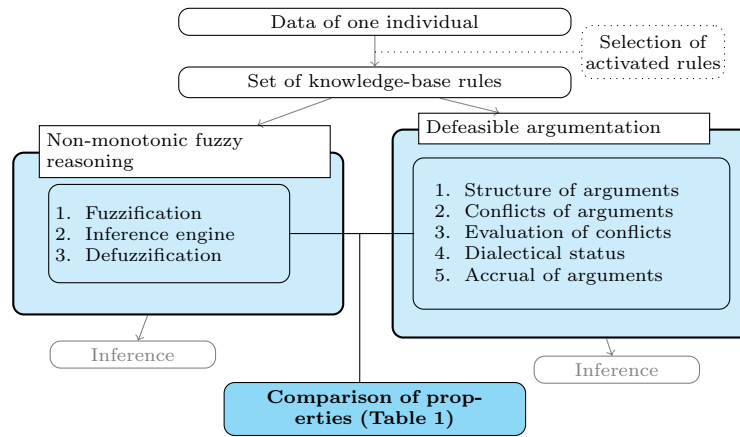


Fig. 1: Design of the comparative research study.

#### 3.1 Data and knowledge-base

Fifty-one biomarkers were described by a clinician and their association with mortality risk levels was provided through the use of ‘IF *premises* THEN *risk-level*’. Some biomarkers were described by natural language terms such as *low* or *high*. This applies also to risk levels (*no*, *low*, *medium*, *high* and *extremely high*). Numerical ranges had to be defined for these terms and were used in different ways within the defeasible argumentation and fuzzy reasoning approaches. Contradictions among biomarkers were also made explicit as rules of the form ‘IF *premises* THEN *conclusion*’. Eventually, some preferences among biomarkers were provided. A contradiction refers to a situation in which some biomarker should not be logically employed, while a preference occurs when a biomarker should be used instead of another biomarker. Since the full knowledge-base contains many rules, contradictions and preferences, it cannot be presented in this

paper but it can be accessed online<sup>1</sup>. A dataset<sup>2</sup> was obtained in a primary health care European hospital and the survival status of the 93 patients was recorded 5 years after data collection. One random individual was picked for a detailed analysis and the associated data can be seen in Table 2. From this information, a set of rules, contradictions and preferences was activated as shown in Table 3. Note that rules, contradictions and preferences activation depend on the patient’s data. A rule designed for female will not be activated for males. A contradiction is not evaluated if its premises or conclusion are not activated. Similarly, a preference is evaluated if both its terms are activated.

Table 2: Data about the biomarkers associated to an elderly. A full description can be found online<sup>1</sup>.

|               |                 |               |               |             |               |                |              |                |
|---------------|-----------------|---------------|---------------|-------------|---------------|----------------|--------------|----------------|
| <b>Age</b>    | <b>Sex</b>      | <b>Hypert</b> | <b>DM</b>     | <b>Fglu</b> | <b>HbA1c</b>  | <b>Chol</b>    | <b>HDL</b>   | <b>statins</b> |
| 60            | female          | high          | no            | 5.3         | 4.17          | 8.7            | 2.06         | no             |
| <b>CVD</b>    | <b>BMI</b>      | <b>w\h</b>    | <b>skinf</b>  | <b>COPB</b> | <b>allerd</b> | <b>draller</b> | <b>analg</b> | <b>derm</b>    |
| no            | 26.68           | 0.88          | 32            | no          | no            | no             | no           | no             |
| <b>OSP</b>    | <b>Psy</b>      | <b>MMS</b>    | <b>CMV</b>    | <b>EBV</b>  | <b>HPA</b>    | <b>LE</b>      | <b>MO</b>    | <b>NEU</b>     |
| ?             | no              | 26            | 2.6           | 170         | 10.4          | 6.94           | 11.7         | 28.8           |
| <b>CRP</b>    | <b>E</b>        | <b>HB</b>     | <b>HTC</b>    | <b>MCV</b>  | <b>FE</b>     | <b>ALB</b>     | <b>Clear</b> | <b>HOMCIS</b>  |
| 3.8           | 4.42            | 140           | 0.41          | 93.2        | 23.6          | 47.7           | 2.11         | 7.9            |
| <b>VitB12</b> | <b>FOLNA</b>    | <b>INS</b>    | <b>CORTIS</b> | <b>PRL</b>  | <b>TSH</b>    | <b>FT3</b>     | <b>FT4</b>   | <b>GAMA</b>    |
| 445           | 37.1            | 8.6           | 470.8         | 86.1        | 0.491         | 5.57           | 12.3         | 12.6           |
| <b>IGE</b>    | <b>anticoag</b> | <b>neo</b>    | <b>Ly</b>     | <b>RF</b>   | <b>ANA</b>    | <b>Death</b>   |              |                |
| 46.2          | yes             | no            | 53.6          | 9           | 36.8          | no             |              |                |

Table 3: Activated rules, contradictions and preferences from data on Table 2.

| Rules   |                     |                                 |                               |
|---|---------------------|---------------------------------|-------------------------------|
| Premises                                      | Risk                | Contradictions                  |                               |
| HDL <i>high</i> ( $> 1.0$ )                   | no risk             | <b>Premises</b>                 | <b>Conclusion</b>             |
| ANA <i>high</i> ( $> 32$ )                    | low risk            | <i>no CVD</i>                   | <i>no Anticoag</i>            |
| w/h <i>high</i> ( $> 0.8$ ) AND <i>female</i> | low risk            | INS <i>low</i> ( $\leq 12.26$ ) | w/h <i>low</i> ( $\leq 0.8$ ) |
| Age $\in [60, 65]$                            | low risk            |                                 |                               |
| Hypert <i>yes</i>                             | extremely high risk |                                 |                               |
| HbA1c <i>high</i> ( $> 3.8$ )                 | low risk            |                                 |                               |
| Anticoag <i>yes</i>                           | medium risk         |                                 |                               |
| Chol <i>high</i> ( $\geq 6.19$ )              | extremely high risk |                                 |                               |
| MO <i>high</i> ( $> 8.6$ )                    | medium risk         |                                 |                               |
| CRP $> 3$                                     | high risk           |                                 |                               |
| Ly <i>high</i> ( $> 40$ )                     | medium risk         |                                 |                               |
| LE $> 6.5$ AND <i>female</i>                  | medium risk         |                                 |                               |
| FE <i>high</i> ( $> 18$ )                     | low risk            |                                 |                               |
| BMI <i>medium</i> ( $\in [26, 29]$ )          | medium risk         |                                 |                               |
|   |                     | <b>Preferences</b>              |                               |
|   |                     | CRP $>$ LE                      |                               |
|   |                     | CRP $>$ ANA                     |                               |
|   |                     | w\h $>$ BMI                     |                               |
|   |                     | Hypert $>$ Age                  |                               |
|   |                     | MO $>$ LE                       |                               |
|   |                     | LY $>$ LE                       |                               |

### 3.2 Non-monotonic fuzzy reasoning inference

**Fuzzification module** Rules in the form “*IF ... THEN ...*” and contradictions rules were constructed from data in Table 3 and depicted in Fig. 2-A on page 7. Afterwards, fuzzy membership functions (FMF) were defined for linguistic variables such as BMI *low* (low body mass index) and FE *high* (high serum iron). Each category of risk had an associated FMF (Fig. 2-B) with input in the range  $[0, 100] \in \mathbb{R}$ . Because of that the input variables (biomarkers) had to be normalised for the same range according to their possible minimum and maximum values.

<sup>1</sup> <http://dx.doi.org/10.6084/m9.figshare.7028480>

<sup>2</sup> <https://doi.org/10.6084/m9.figshare.7028516.v1>

Fig. 2-C depicts examples of FMFs for *FE high* and *FE low*. Not all biomarkers had a fuzzy representation provided by the domain expert and were incorporated into the fuzzy inference as crisp variables (membership degree always 0 or 1). For the case under analysis (picked patient), the crisp variables are HDL, Hypert, Anticoag, MO, CRP and LE. Due to space limitations, not all FMFs are shown here but they can be accessed online<sup>1</sup>.

**Inference engine** For each linguistic term provided by the domain expert, and used within rules and exceptions, its membership degree have to be computed by evaluating the associated membership function with a given input (from table 2). Once each membership degree of each linguistic term in the premises of a rule has been computed, then also a degree of truth for that rule can be computed. This can be done by employing some fuzzy operators OR and AND. The ones selected here are: *Zadeh*<sup>3</sup>, *Product*<sup>4</sup> and *Lukasiewicz*<sup>5</sup> (Fig. 2-D). Eventually, contradictions, which in fuzzy reasoning define non-monotonicity, have to be evaluated. This evaluation can be done using Possibility Theory, as proposed by [20] for fuzzy reasoning with rule-based systems. In this case truth values are represented by *possibility* (Pos) and *necessity* (Nec) as defined on Section 2. The Nec of a proposition is treated here as its membership grade and the Pos is always 1 for all propositions. Under these circumstances ( $\text{Pos} \geq \text{Nec}$ ), the effect on the necessity of a proposition  $A$  ( $\text{Nec}(A)$ ) by a set of  $n$  propositions  $Q$  which refute  $A$  is derivable in [20] and given by:

$$\text{Nec}(A) = \min(\text{Nec}(A), \neg\text{Nec}(Q_1), \dots, \neg\text{Nec}(Q_n)) \quad (1)$$

where  $\neg\text{Nec}(Q) = 1 - \text{Nec}(Q)$ . In addition, an order of precedence has to be defined when applying equation 1. In this study, contrarily to usual fuzzy control systems, the reasoning is done in a single step with all the activated rules fired at once. Nonetheless, it is possible to organise exceptions in a tree structure in which the consequent of an exception is the antecedent of the next exception. Fig. 2-E illustrates this structure which allows equation 1 to be applied from the roots to the leaves. The updated truth values of those rules subject to refutation by other rules are listed in Fig. 2-F. The last step of the inference engine is to aggregate all the truth values of the membership functions associated to each risk category (grouped by the same category), by using the fuzzy-OR operator (as per figure 2-G). The output of this can be graphically represented (Fig. 2-H).

**Defuzzification module** A single defuzzified scalar which represents the final mortality risk inferred has to be computed. Two common methods are selected: *mean of max* and *centroid*. The former returns the average of all  $x$  coordinates (mortality risks) whose respective  $y$  coordinates (membership grades) are maximum in the graphical representation (Fig. 2-H). The latter returns the coordinates of the centre of gravity of the same graphical representation (the  $x$  coordinate is the final scalar). Fig. 2-I lists all the final inferences produced for the patient under analysis.

<sup>3</sup> Given propositions  $a, b$ , then fuzzy-AND and fuzzy-OR are “ $\min(a, b)$ ”, “ $\max(a, b)$ ”.

<sup>4</sup> Product’s fuzzy-AND and fuzzy-OR are respectively “ $a \times b$ ” and “ $a + b - a \times b$ ”.

<sup>5</sup> Lukasiewicz’s fuzzy-AND and fuzzy-OR are “ $\max(a + b - 1, 0)$ ” and “ $\min(a + b, 1)$ ”.

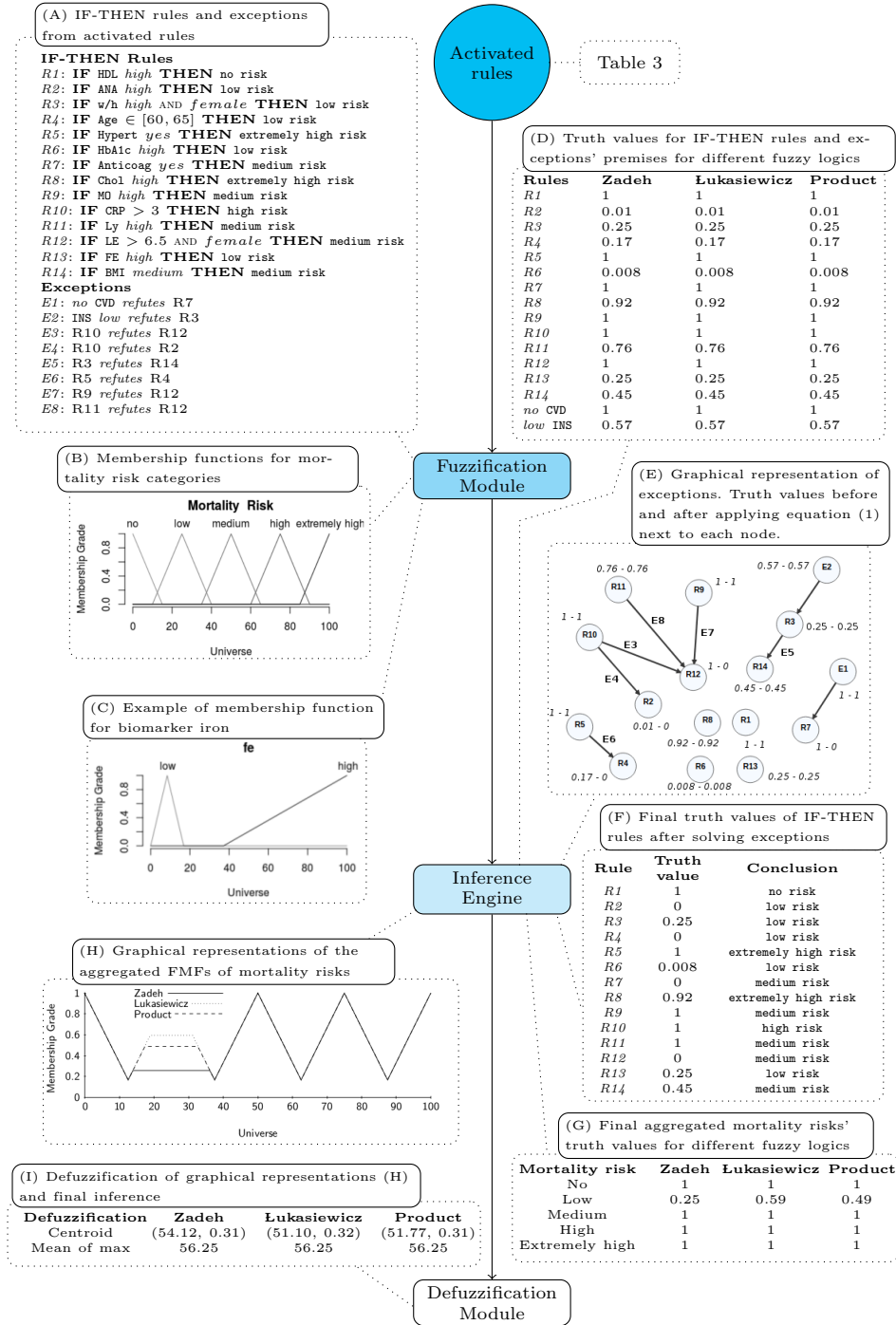


Fig. 2: An illustration of the non-monotonic fuzzy reasoning process for the selected elderly patient. The order of operations is from A to I.



### 3.3 Defeasible argumentation inference

**Layer 1 - Definition of the internal structure of arguments** The first step of a defeasible argumentation process is to define a set of *arguments*. Internally these are generally composed by a set of premises and a conclusion derivable by applying an inference rule  $\rightarrow$ . A typical version of this is known as forecast argument in which, from a set of premises, a conclusion can be reasonably forecasted. Examples can be found in Table 3 (left) where premises reasonably forecast a degree of risk of mortality (as also listed in Fig. 3-A). Note that, in contrast to fuzzy rules, the natural language linguistic terms associated to the premises are not quantitatively exploited. Instead, the premises are evaluated true or not if input values are within certain ranges.

**Layer 2 - Definition of the conflicts of arguments** Given a set of forecast arguments, the next step for modelling an underlying knowledge-base, is to define the conflicts between arguments. The goal is to evaluate potential inconsistencies and identify invalid arguments through the notion of attack (conflict). In this research, the notion of *undercutting attack* [16] is employed for the resolution of conflicts. It defines an exception, where the application of the knowledge carried in some argument is no longer allowed. It is formed by a set of premises and an undercutting inference  $\Rightarrow$  to another argument. Examples of undercutting attacks, derived from Table 3 (right), are in Fig. 3-B. All the designed arguments and attacks can now be seen as an argumentation framework (Fig. 3-C).

**Layer 3 - Evaluation of the conflicts of arguments** After conflicts formalisation, these can be evaluated using different approaches such as considering the strength of attacks or the notion of preferentiality of arguments [9]. Alternatively, as in this study, conflicts follow a binary relation, that means, if two arguments (attacker and attacked) are activated, the conflict between them is fully considered.

**Layer 4 - Definition of the dialectical status of arguments** Given an argumentation framework and a notion of conflict, it is necessary to define the set of defeated arguments. An argument A is *defeated* by B if there is a valid attack from A to B. A well-known approach has been proposed by [8] in the form of acceptability semantics. A semantics is an algorithm designed to produce a set of acceptable and conflict-free arguments, called *extensions*. Note that the internal structure of arguments is not considered at this stage. Well-known examples are the *grounded* and the *preferred* semantics. In this study, only the former algorithm is illustrated (Fig. 3-D). Fig. 3-E depicts its computed extension.

**Layer 5 - Accrual of acceptable arguments** Having a set of acceptable forecast arguments, it is necessary to accrue them in case a final inference is required. If no quantity can be associated to an argument, then the conclusion supported by the highest number of arguments could be chosen as final inference. In case arguments can be quantitatively evaluated (they carry a value as in this study), then several approaches can be used, including the selection of measures of central tendency such as average (used in this study). Fig. 3-F illustrates the value associated to each argument and the final inference which is their average.

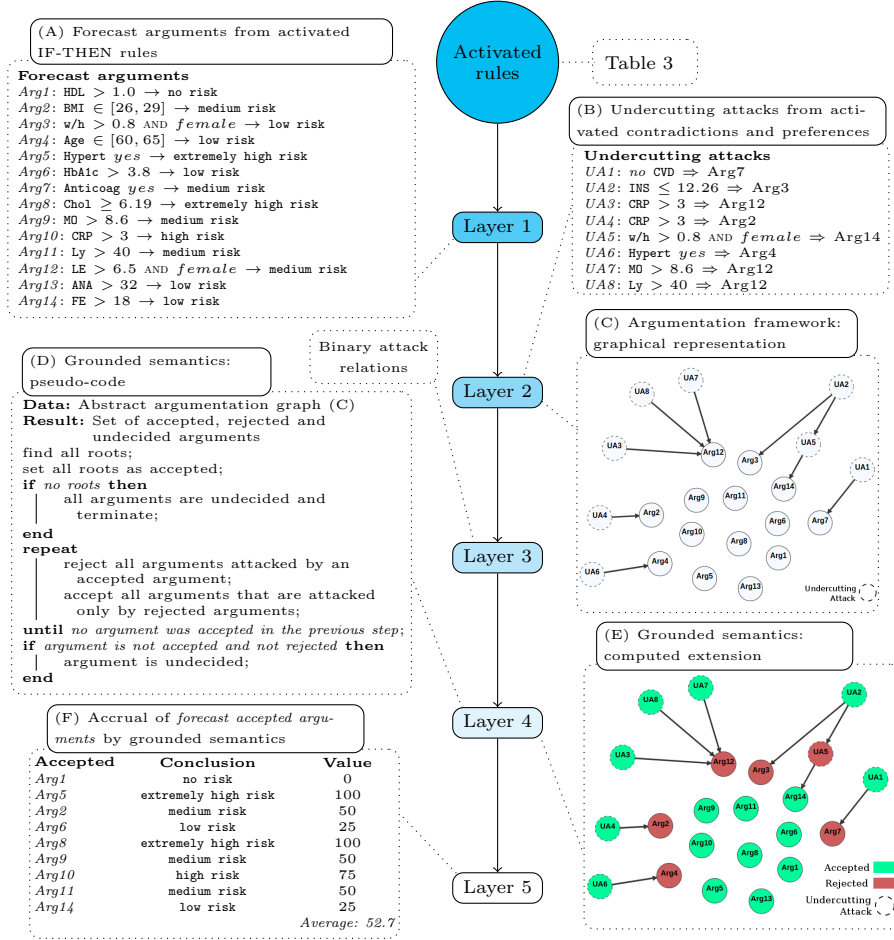


Fig. 3: An illustration of the defeasible argumentation process for an elderly (order from A to F).

## 4 Comparison and discussion

A comparative qualitative analysis of the explanatory capacity of the defeasible argumentation and non-monotonic fuzzy reasoning processes is performed by using the properties listed in Table 1 (Section 2).

### Understandability/Post-hoc Interpretability

- *Non-monotonic fuzzy reasoning* - The inferential process is aligned to the expert’s knowledge and natural language for most of its parts, which makes it generally intuitively understandable by humans. However, this does not apply for some parts, such as the normalisation of the input values, the selection of fuzzy logic and the defuzzification mechanism. Some mathematical reasoning is required to select suitable parameters of these parts.
- *Defeasible argumentation* - The initial reasoning steps (layers 1-3) are built upon the same natural language terms provided by the domain expert in the

knowledge-base. In layer 4 the grounded semantics was selected. This particular semantics is not a complex algorithm to understand: intuitively, an argument is only rejected if it is attacked by an accepted argument. In layer 5, the accrual of accepted arguments can be done by an intuitive measure of central tendency (average here). In case more complex (less intuitive) semantics, such as preferred [8] or ranking-based [3], are employed, then the understandability of the inferential process might be compromised.

#### **Simulatability**

- *Non-monotonic fuzzy reasoning* - Practical applications built upon a small number of simple membership functions could support simulatability. However, with more complex membership functions, a domain expert is not likely able to step through their calculation with high precision and in a reasonable time. Similarly, this applies to the calculations required within the defuzzification unit (example, computation of the centroid).
- *Defeasible argumentation* - Reasonably, an expert could perform the calculations behind all the steps of the inferential process. However, this would be significantly impacted by the number of arguments in the knowledge-base, the complexity of selected acceptability semantics and the accrual strategy.

#### **Extendibility**

- *Non-monotonic fuzzy reasoning* - New rules can be added/updated in the light of new information. However, fuzzy membership functions have to be defined, demanding further effort, not common in human reasoning.
- *Defeasible argumentation* - New arguments can be constructed from new information and easily plugged-in the knowledge-base. They follow the same structure (premise to conclusions) which does not require the definitions of mathematical functions and is close to the way humans reason.

#### **Computational complexity**

- *Non-monotonic fuzzy reasoning* - The full inferential process, in the worst case, is linear in the number of rules.
- *Defeasible argumentation* - Layers 3 and 5 are linear in the number of arguments and attacks relations. However, for layer 4 (application of acceptability semantics for the computation of the dialectical status of arguments), complexity can range from linear (example the grounded semantics) to exponential (example the preferred semantics) [8].

#### **Algorithmic transparency**

- *Non-monotonic fuzzy reasoning* - The inferential process can be applied across different domains. A knowledge-base is a formalisation of a reasoning activity for a specific underlying domain, thus it can be re-used or extended provided the new domains are similar. However, it is important to highlight that traditional fuzzy reasoning has not been designed for application in those domains requiring non-monotonic reasoning. In fact, in this study, the traditional fuzzy reasoning process has been extended through the incorporation of Possibility Theory in order to deal with non-monotonicity.
- *Defeasible argumentation* - The inferential process can be applied across different domains. By nature, defeasible argumentation is suitable for appli-

cation in domains requiring non-monotonic reasoning activities. However, in the absence of conflicts, the inferential process can still be applied as it is. The analysis of the two reasoning approaches suggests that defeasible argumentation might lead to explanations that are more suitable to understand for humans, both for a domain expert and a lay person. In fact, through the comparison performed above, on one hand, without some comprehension of fuzzy logic and its membership functions, the understandability/post-hoc interpretability, simulatability of non-monotonic fuzzy reasoning and the extendibility of its models is compromised. On the other hand, defeasible argumentation tends to use the same natural language terms, provided by the domain expert, throughout the whole inferential process, except in the conflict resolution layer (semantics). Semantics vary in computational complexity (linear or exponential in the number of arguments), allowing fuzzy reasoning to offer an equal or lower complexity, since its fuzzification-engine-defuzzification layers are always linear in the number of rules. However, Possibility Theory always requires the specification of a precedence order of exceptions in the inference engine of fuzzy reasoning. Contrarily to acceptability semantics that do not require any precedence order of attacks for solving conflicts, thus it has a higher algorithmic transparency.

## 5 Conclusion and future work

Despite theoretical advances in defeasible argumentation, to the best of our knowledge, there is lack of research devoted to the examination of the degree of explainability that this reasoning approach can offer to illustrate inferences to humans in real-world applications. Therefore, this research focused on a qualitative comparison of the degree of explainability of defeasible argumentation and non-monotonic fuzzy reasoning in a real-world setting: prediction of mortality of elderly people by using biomarkers. The inferential processes behind the two selected reasoning techniques were meticulously illustrated and exploited. The comparison was performed using six properties for explainability extracted from the literature. A qualitative discussion of these properties show how defeasible argumentation has a greater potential for tackling the problem of explainability of reasoning activities under uncertainty, partial and conflictual information. The contribution of this study is to situate defeasible argumentation among similar approaches for reasoning under uncertainty in terms of degree of explainability.

## Acknowledgments

Lucas Middeldorf Rizzo would like to thank CNPq (Conselho Nacional de Desenvolvimento Científico e Tecnológico) for his Science Without Borders scholarship, proc n. 232822/2014-0.

## References

1. Allahyari, H., Lavesson, N.: User-oriented assessment of classification model understandability. In: 11th scandinavian conference on Artificial intelligence (2011)

2. Bench-Capon, T.J., Dunne, P.E.: Argumentation in artificial intelligence. *Artificial intelligence* 171(10-15), 619–641 (2007)
3. Bonzon, E., Delobelle, J., Konieczny, S., Maudet, N.: A comparative study of ranking-based semantics for abstract argumentation. In: *AAAI*. pp. 914–920 (2016)
4. Bryant, D., Krause, P.: A review of current defeasible reasoning implementations. *The Knowledge Engineering Review* 23(3), 227–260 (2008)
5. Caruana, R., Lou, Y., Gehrke, J., Koch, P., Sturm, M., Elhadad, N.: Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission. In: *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. pp. 1721–1730. ACM (2015)
6. Castro, J.L., Trillas, E., Zurita, J.M.: Non-monotonic fuzzy reasoning. *Fuzzy Sets and Systems* 94(2), 217–225 (1998)
7. Dubois, D., Prade, H.: Possibility theory: qualitative and quantitative aspects. In: *Quantified representation of uncertainty and imprecision*, pp. 169–226 (1998)
8. Dung, P.M.: On the acceptability of arguments and its fundamental role in non-monotonic reasoning, logic programming and n-person games. *Artificial intelligence* 77(2), 321–358 (1995)
9. García, D., Simari, G.: Strong and weak forms of abstract argument defense. *Computational Models of Argument: Proceedings of COMMA 2008* 172, 216 (2008)
10. Gegov, A., Gobalakrishnan, N., Sanders, D.: Rule base compression in fuzzy systems by filtration of non-monotonic rules. *Journal of Intelligent & Fuzzy Systems* 27(4), 2029–2043 (2014)
11. Giraud-Carrier, C.: Beyond predictive accuracy: what. In: *Proceedings of the ECML-98 Workshop on Upgrading Learning to Meta-Level: Model Selection and Data Transformation*. pp. 78–85 (1998)
12. Lipton, Z.C.: The mythos of model interpretability. *Queue* 16(3), 30:31–30:57 (2018)
13. Longo, L.: Argumentation for knowledge representation, conflict resolution, defeasible inference and its integration with machine learning. In: *Machine Learning for Health Informatics*, pp. 183–208. Springer (2016)
14. Longo, L., Kane, B., Hederman, L.: Argumentation theory in health care. In: *Computer-Based Medical Systems, 25th Int. Symposium on*. pp. 1–6. IEEE (2012)
15. Passino, K.M., Yurkovich, S., Reinfrank, M.: Fuzzy control
16. Prakken, H.: An abstract framework for argumentation with structured arguments. *Argument and Computation* 1(2), 93–124 (2010)
17. Rizzo, L., Longo, L.: Representing and inferring mental workload via defeasible reasoning: a comparison with the nasa task load index and the workload profile. In: *1st Workshop on Advances In Argumentation In Artificial Intelligence*. pp. 126–140 (2017)
18. Rizzo, L., Majnaric, L., Dondio, P., Longo, L.: An investigation of argumentation theory for the prediction of survival in elderly using biomarkers. In: *Int. Conf. on Artificial Intelligence Applications and Innovations*. pp. 385–397. Springer (2018)
19. Rizzo, L., Majnaric, L., Longo, L.: A comparative study of defeasible argumentation and non-monotonic fuzzy reasoning for elderly survival prediction using biomarkers. In: *AI\*IA 2018 - Advances in Artificial Intelligence - XVIIth Int. Conference of the Italian Association for Artificial Intelligence*. pp. 197–209 (2018)
20. Siler, W., Buckley, J.J.: Fuzzy expert systems and fuzzy reasoning (2005)
21. Zadeh, L.A., et al.: Fuzzy sets. *Information and control* 8(3), 338–353 (1965)