

Non-Classical Logics for Explanations in AI Systems

Emiliano Lorini
IRIT-CNRS, Toulouse University, France

9 November 2022
XAISeminar@Imperial

Logic-based approaches to explanation

- Different notions of explanation are studied in the XAI domain
 - Abductive
 - Contrastive
 - Counterfactual
- Logic-based XAI mostly based on propositional logic (PL)

- Fundamental building blocks of explanation:

- Counterfactual dependence
- Variance/invariance:

*[I]nvariance is a **modal** notion – it has to do with whether a relationship would remain stable under various hypothetical changes [Woodward 2002, p. 225].*

- Imperfect knowledge of the classifier (“black box”)
⇒ Epistemic/subjective explanation

- Beyond PL: need for more **expressive** languages

- Non-classical logics:

- Modal logic (ML)
- Conditional logic (CL)
- Epistemic logic (EL), dynamic EL (DEL)
- Deontic logic (DL)

- 1 Explanations in “white box” classifiers
- 2 Explanations in “black box” classifiers
- 3 Open problems and future extensions

- 1 Explanations in “white box” classifiers
- 2 Explanations in “black box” classifiers
- 3 Open problems and future extensions

Reasoning about “white box” classifiers

Main idea: a binary input classifier is a partition of all possible input instances in an S5 Kripke model

| Permanent job (<i>pe</i>) | > 3000 € monthly salary (<i>sa</i>) | EU citizenship (<i>eu</i>) | Loan |
|--------------------------------|--|---------------------------------|------|
| 0 | 0 | 0 | No |
| 0 | 0 | 1 | No |
| 0 | 1 | 0 | No |
| 1 | 0 | 0 | Yes |
| 0 | 1 | 1 | Yes |
| 1 | 0 | 1 | Yes |
| 1 | 1 | 0 | Yes |
| 1 | 1 | 1 | Yes |

Figure: A classifier

| States/instances | f_1 |
|------------------------|-------|
| $s_1 = \{\}$ | No |
| $s_2 = \{eu\}$ | No |
| $s_3 = \{sa\}$ | No |
| $s_4 = \{pe\}$ | Yes |
| $s_5 = \{sa, eu\}$ | Yes |
| $s_6 = \{pe, eu\}$ | Yes |
| $s_7 = \{pe, sa\}$ | Yes |
| $s_8 = \{pe, sa, eu\}$ | Yes |

Figure: Its S5 representation

- Atm_0 : countable set of atoms representing input features
- Val : finite set of classification values (or classes)

Definition (Classifier model)

A **classifier model** (CM) is a tuple $C = (S, f)$ where

- $S \subseteq 2^{Atm_0}$ is a set of input instances,
- $f : S \rightarrow Val$ is a classification function.

$$\varphi ::= p \mid \mathbf{t}(x) \mid \neg\varphi \mid \varphi \wedge \varphi \mid \Box_I \varphi$$

with p ranging over Atm_0 and x ranging over Val

$\mathbf{t}(x) \approx$ “the actual input instance is classified as x ”

$\Box_I \varphi \approx$ “the classifier *necessarily* satisfies φ ”

\approx “ φ is true for all input instances of the classifier”

Semantic interpretation wrt CM $C = (S, f)$ and $s \in S$ (**pointed CM**):

$$(C, s) \models p \iff p \in s$$

$$(C, s) \models \mathbf{t}(x) \iff f(s) = x$$

$$(C, s) \models \Box_I \varphi \iff \forall s' \in S : (C, s') \models \varphi$$

Useful “ceteris paribus” modalities

Let $X \subseteq Atm_0$ finite:

$$[X]\varphi =_{def} \bigwedge_{Y \subseteq X} ((\bigwedge_{p \in Y} \wedge \bigwedge_{p \in X \setminus Y} \neg p) \rightarrow \Box_I((\bigwedge_{p \in Y} \wedge \bigwedge_{p \in X \setminus Y} \neg p) \rightarrow \varphi))$$

We have:

$$(C, s) \models [X]\varphi \iff \forall s' \in S : \text{if } (s \cap X) = (s' \cap X) \text{ then } (C, s') \models \varphi$$

$[X]\varphi \approx$ “ φ is true all atoms in X being equal”

\approx “ φ is true regardless of the value of the atoms in $Atm_0 \setminus X$ ”

Useful “ceteris paribus” modalities

- Connection with prop. **dependence logic** [Yang & Väänänen, 2016]
- Dependence atom (“ q only depends on p_1, \dots, p_k ”):

$$\text{Dep}(p_1, \dots, p_k, q) =_{\text{def}} [\emptyset] (q \rightarrow [\{p_1, \dots, p_k\}] q) \wedge \\ [\emptyset] (\neg q \rightarrow [\{p_1, \dots, p_k\}] \neg q)$$

| | Finite (fixed) variables | Infinite variables |
|---|---------------------------------|---------------------------|
| Modalities $[X]$ are defined as abbreviations | Polynomial | NP-complete |
| Modalities $[X]$ are primitives | Polynomial | NEXPTIME-complete |

Table: Summary of complexity results

Let λ be a term (conjunction of literals):

- Prime implicant:

$$\text{PImp}(\lambda, x) =_{\text{def}} [\emptyset] \left(\lambda \rightarrow (t(x) \wedge \bigwedge_{p \in \text{Atm}(\lambda)} \langle \text{Atm}(\lambda) \setminus \{p\} \rangle \neg t(x)) \right)$$

- Abductive explanation:

$$\text{AXp}(\lambda, x) =_{\text{def}} \lambda \wedge \text{PImp}(\lambda, x)$$

- Contrastive explanation:

$$\begin{aligned} \text{CXp}(\lambda, x) =_{\text{def}} & \lambda \wedge \langle \text{Atm}_0 \setminus \text{Atm}(\lambda) \rangle \neg t(x) \wedge \\ & \bigwedge_{p \in \text{Atm}(\lambda)} [(\text{Atm}_0 \setminus \text{Atm}(\lambda)) \cup \{p\}] t(x) \end{aligned}$$

| States/instances | f_1 |
|------------------------|-------|
| $s_1 = \{\}$ | No |
| $s_2 = \{eu\}$ | No |
| $s_3 = \{sa\}$ | No |
| $s_4 = \{pe\}$ | Yes |
| $s_5 = \{sa, eu\}$ | Yes |
| $s_6 = \{pe, eu\}$ | Yes |
| $s_7 = \{pe, sa\}$ | Yes |
| $s_8 = \{pe, sa, eu\}$ | Yes |

$s_4 \models \text{AXp}(pe, \text{Yes})$

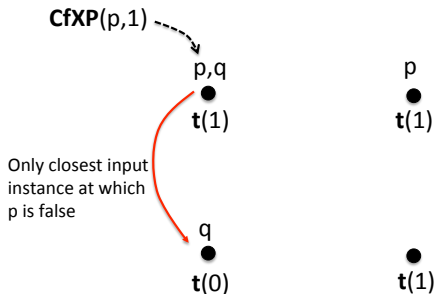
$s_2 \models \text{CXp}(\neg sa, \text{No})$

Explanations

- Counterfactual explanation:

$$\text{CfXP}(\varphi, x) =_{\text{def}} t(x) \wedge (\neg\varphi \Rightarrow \neg t(x))$$

Remark: Lewis-like conditional \Rightarrow defined as an abbreviation in finite-variable case (semantics based on Hamming dist.)



| States/instances | f_1 |
|--------------------|-------|
| $s_1=\{\}$ | No |
| $s_2=\{eu\}$ | No |
| $s_3=\{sa\}$ | No |
| $s_4=\{pe\}$ | Yes |
| $s_5=\{sa,eu\}$ | Yes |
| $s_6=\{pe,eu\}$ | Yes |
| $s_7=\{pe,sa\}$ | Yes |
| $s_8=\{pe,sa,eu\}$ | Yes |

$s_2 \models \text{CfXp}(\neg sa, No)$

Some remarkable properties

⇒ Principle of sufficient reason (PSR):

*Of everything whatsoever a cause or reason must be assigned, either for its existence, or for its non-existence.
[Spinoza, Ethics, 1p11d2]*

$$\models_{\text{Definite}} t(x) \rightarrow \bigvee_{\lambda \in \text{Term}} \text{AXp}(\lambda, x)$$

⇒ 'Atomic' CfXp and CXp coincide

$$\models \text{CXp}(l, x) \leftrightarrow \text{CfXp}(l, x) \text{ with } l \text{ a literal}$$

Local bias:

$$\text{Bias}(x) =_{\text{def}} t(x) \wedge \langle \text{NF} \rangle \neg t(x)$$

with PF the set of protected features and $\text{NF} = \text{Atm}_0 \setminus \text{PF}$

$$\models \text{Bias}(x) \leftrightarrow \bigvee_{\text{Atm}(\lambda) \subseteq \text{PF}} \text{CXp}(\lambda, x)$$

Global bias:

$$\text{GBias} =_{\text{def}} \Diamond_I \left(\bigvee_{x \in \text{Val}} \text{Bias}(x) \right)$$

Let $PF = \{eu\}$

| States/instances | f_1 |
|------------------------|-------|
| $s_1 = \{\}$ | No |
| $s_2 = \{eu\}$ | No |
| $s_3 = \{sa\}$ | No |
| $s_4 = \{pe\}$ | Yes |
| $s_5 = \{sa, eu\}$ | Yes |
| $s_6 = \{pe, eu\}$ | Yes |
| $s_7 = \{pe, sa\}$ | Yes |
| $s_8 = \{pe, sa, eu\}$ | Yes |

$s_3 \models \text{Bias}(No)$

- 1 Explanations in “white box” classifiers
- 2 Explanations in “black box” classifiers
- 3 Open problems and future extensions

From “white box” to “black box” classifiers

- Two-dimensional semantics: instance \times classifier
- Horizontal dimension \approx uncertainty about classifier's properties
- Bimodal language



Bob

| | f_1 |
|--------------------|-------|
| $s_1=\{\}$ | No |
| $s_2=\{eu\}$ | No |
| $s_3=\{sa\}$ | No |
| $s_4=\{pe\}$ | Yes |
| $s_5=\{sa,eu\}$ | Yes |
| $s_6=\{pe,eu\}$ | Yes |
| $s_7=\{pe,sa\}$ | Yes |
| $s_8=\{pe,sa,eu\}$ | Yes |

| | f_2 |
|--------------------|-------|
| $s_1=\{\}$ | No |
| $s_2=\{eu\}$ | No |
| $s_3=\{sa\}$ | No |
| $s_4=\{pe\}$ | Yes |
| $s_5=\{sa,eu\}$ | No |
| $s_6=\{pe,eu\}$ | Yes |
| $s_7=\{pe,sa\}$ | Yes |
| $s_8=\{pe,sa,eu\}$ | Yes |

$$\varphi ::= p \mid t(x) \mid \neg\varphi \mid \varphi \wedge \varphi \mid \Box_I\varphi \mid \Box_F\varphi$$

with p ranging over Atm_0 and x ranging over Val

$\Box_I\varphi \approx$ “the actual classifier *necessarily* satisfies φ ,
(regardless of the input instance)”

$\Box_F\varphi \approx$ “the actual input instance *necessarily* satisfies φ ,
(regardless of the classifier)”

Definition

A **multi-classifier model (MCM)** is a pair $\Gamma = (S, \Phi)$ where:

- $S \subseteq 2^{Atm_0}$ (set of input instances),
- $\Phi \subseteq Val^S$ (set of possible classifiers).

Semantic interpretation of formulas wrt **pointed MCM** (Γ, s, f) with $\Gamma = (S, \Phi)$ an MCM, $s \in S$ and $f \in \Phi$:

$$\begin{aligned}(\Gamma, s, f) \models p &\iff p \in s \\(\Gamma, s, f) \models t(x) &\iff f(s) = x \\(\Gamma, s, f) \models \Box_I \varphi &\iff \forall s' \in S : (\Gamma, s', f) \models \varphi \\(\Gamma, s, f) \models \Box_F \varphi &\iff \forall f' \in \Phi : (\Gamma, s, f') \models \varphi\end{aligned}$$

Two-dimensional semantics



Bob

| | f_1 | | f_2 |
|------------------------|-------|------------------------|-------|
| $s_1 = \{\}$ | No | $s_1 = \{\}$ | No |
| $s_2 = \{eu\}$ | No | $s_2 = \{eu\}$ | No |
| $s_3 = \{sa\}$ | No | $s_3 = \{sa\}$ | No |
| $s_4 = \{pe\}$ | Yes | $s_4 = \{pe\}$ | Yes |
| $s_5 = \{sa, eu\}$ | Yes | $s_5 = \{sa, eu\}$ | No |
| $s_6 = \{pe, eu\}$ | Yes | $s_6 = \{pe, eu\}$ | Yes |
| $s_7 = \{pe, sa\}$ | Yes | $s_7 = \{pe, sa\}$ | Yes |
| $s_8 = \{pe, sa, eu\}$ | Yes | $s_8 = \{pe, sa, eu\}$ | Yes |

Bob is $\{sa\}$ and (only) knows that:

- his application was unsuccessful
- necessarily not having a permanent job and not having a good salary will make loan application unsuccessful
- necessarily having a permanent job will make loan application successful

$$\begin{aligned}(s_3, f_1) \models & \Box_F t(No) \wedge \\ & \Box_F \Box_I ((\neg sa \wedge \neg pe) \rightarrow t(No)) \wedge \\ & \Box_F \Box_I (pe \rightarrow t(Yes))\end{aligned}$$

Axiomatics for Atm_0 finite

$$\blacksquare \in \{\Box_I, \Box_F\}$$

$$(\blacksquare\varphi \wedge \blacksquare(\varphi \rightarrow \psi)) \rightarrow \blacksquare\psi \quad (\mathbf{K}_{\blacksquare})$$

$$\blacksquare\varphi \rightarrow \varphi \quad (\mathbf{T}_{\blacksquare})$$

$$\blacksquare\varphi \rightarrow \blacksquare\blacksquare\varphi \quad (4_{\blacksquare})$$

$$\neg\blacksquare\varphi \rightarrow \blacksquare\neg\blacksquare\varphi \quad (5_{\blacksquare})$$

$$\Box_F\Box_I\varphi \leftrightarrow \Box_I\Box_F\varphi \quad (\mathbf{Comm})$$

$$\bigvee_{x \in Val} \mathbf{t}(x) \quad (\mathbf{AtLeast}_{\mathbf{t}(x)})$$

$$\mathbf{t}(x) \rightarrow \neg\mathbf{t}(y) \text{ if } x \neq y \quad (\mathbf{AtMost}_{\mathbf{t}(x)})$$

$$(\mathbf{cn}_{X, Atm_0} \wedge \mathbf{t}(x)) \rightarrow \Box_I(\mathbf{cn}_{X, Atm_0} \rightarrow \mathbf{t}(x)) \quad (\mathbf{Func})$$

$$p \rightarrow \Box_F p \quad (\mathbf{Indep}_{\Box_F, p})$$

$$\neg p \rightarrow \Box_F \neg p \quad (\mathbf{Indep}_{\Box_F, \neg p})$$

$$\frac{\varphi}{\blacksquare\varphi} \quad (\mathbf{Nec}_{\blacksquare})$$

\Rightarrow Satisfiability checking: **polynomial**

Axiomatics for Atm_0 infinite

$$\blacksquare \in \{\Box_I, \Box_F\}$$

$$(\blacksquare\varphi \wedge \blacksquare(\varphi \rightarrow \psi)) \rightarrow \blacksquare\psi \quad (\mathbf{K}_{\blacksquare})$$

$$\blacksquare\varphi \rightarrow \varphi \quad (\mathbf{T}_{\blacksquare})$$

$$\blacksquare\varphi \rightarrow \blacksquare\blacksquare\varphi \quad (4_{\blacksquare})$$

$$\neg\blacksquare\varphi \rightarrow \blacksquare\neg\blacksquare\varphi \quad (5_{\blacksquare})$$

$$\Box_F\Box_I\varphi \leftrightarrow \Box_I\Box_F\varphi \quad (\mathbf{Comm})$$

$$\bigvee_{x \in Val} t(x) \quad (\mathbf{AtLeast}_{t(x)})$$

$$t(x) \rightarrow \neg t(y) \text{ if } x \neq y \quad (\mathbf{AtMost}_{t(x)})$$

~~$$(\mathbf{cn}_{X, Atm_0} \wedge t(x)) \rightarrow \Box_I(\mathbf{cn}_{X, Atm_0} \rightarrow t(x)) \quad (\mathbf{Funct})$$~~

$$p \rightarrow \Box_F p \quad (\mathbf{Indep}_{\Box_F, p})$$

$$\neg p \rightarrow \Box_F \neg p \quad (\mathbf{Indep}_{\Box_F, \neg p})$$

$$\frac{\varphi}{\blacksquare\varphi} \quad (\mathbf{Nec}_{\blacksquare})$$

\Rightarrow Satisfiability checking: in NEXPTIME

Idea of the proof: polynomial reduction into satisfiability checking for product modal logic $S5^2$

From objective to subjective explanation

| | Local | Global |
|------------|-----------------------------|------------------------------|
| Objective | $\text{AXp}(\lambda, x)$ | $\text{PImp}(\lambda, x)$ |
| Subjective | $\text{SubAXp}(\lambda, x)$ | $\text{SubPImp}(\lambda, x)$ |

Table: Objective vs subjective explanation

$$\text{SubPImp}(\lambda, x) =_{\text{def}} \Box_{\text{F}} \text{PImp}(\lambda, x)$$

$$\text{SubAXp}(\lambda, x) =_{\text{def}} \Box_{\text{F}} \text{AXp}(\lambda, x)$$

A negative property



Bob

| | f_1 |
|--------------------|-------|
| $s_1=\{\}$ | No |
| $s_2=\{eu\}$ | No |
| $s_3=\{sa\}$ | No |
| $s_4=\{pe\}$ | Yes |
| $s_5=\{sa,eu\}$ | Yes |
| $s_6=\{pe,eu\}$ | Yes |
| $s_7=\{pe,sa\}$ | Yes |
| $s_8=\{pe,sa,eu\}$ | Yes |

| | f_2 |
|--------------------|-------|
| $s_1=\{\}$ | No |
| $s_2=\{eu\}$ | No |
| $s_3=\{sa\}$ | No |
| $s_4=\{pe\}$ | Yes |
| $s_5=\{sa,eu\}$ | No |
| $s_6=\{pe,eu\}$ | Yes |
| $s_7=\{pe,sa\}$ | Yes |
| $s_8=\{pe,sa,eu\}$ | Yes |

PSR principle does not hold in the “black box” setting:

$$(s_3, f_1) \models \Box_{\text{Ft}}(\text{No}) \wedge \neg \bigvee_{\lambda \in \text{Term}} \text{SubAXp}(\lambda, \text{No})$$

Extension: acquiring information about actual classifier

⇒ **Language:**

$$\varphi ::= p \mid t(x) \mid \neg\varphi \mid \varphi \wedge \varphi \mid \Box_I\varphi \mid \Box_F\varphi \mid [\varphi!]\psi$$

$[\varphi!]\psi \approx$ “ ψ holds after having discarded all classifiers that do not globally satisfy property φ ”

⇒ **Semantic interpretation** of dynamic modality $[\varphi!]$:

$$(\Gamma, s, f) \models [\varphi!]\psi \iff \text{if } (\Gamma, s, f) \models \Box_I\varphi \text{ then } (\Gamma^{\varphi!}, s, f) \models \psi$$

where $\Gamma^{\varphi!} = (S^{\varphi!}, \Phi^{\varphi!})$ is the MCM such that:

$$S^{\varphi!} = S$$

$$\Phi^{\varphi!} = \{f' \in \Phi : \forall s' \in S, (\Gamma, s', f') \models \varphi\}$$

Extension: acquiring information about actual classifier

Let $PF = \{eu\}$



Bob

| | f_1 | | f_2 |
|------------------------|-------|------------------------|-------|
| $s_1 = \{\}$ | No | $s_1 = \{\}$ | No |
| $s_2 = \{eu\}$ | No | $s_2 = \{eu\}$ | No |
| $s_3 = \{sa\}$ | No | $s_3 = \{sa\}$ | No |
| $s_4 = \{pe\}$ | Yes | $s_4 = \{pe\}$ | Yes |
| $s_5 = \{sa, eu\}$ | Yes | $s_5 = \{sa, eu\}$ | No |
| $s_6 = \{pe, eu\}$ | Yes | $s_6 = \{pe, eu\}$ | Yes |
| $s_7 = \{pe, sa\}$ | Yes | $s_7 = \{pe, sa\}$ | Yes |
| $s_8 = \{pe, sa, eu\}$ | Yes | $s_8 = \{pe, sa, eu\}$ | Yes |

GBias!



Bob

| | f_1 | | f_2 |
|------------------------|-------|------------------------|-------|
| $s_1 = \{\}$ | No | $s_1 = \{\}$ | No |
| $s_2 = \{eu\}$ | No | $s_2 = \{eu\}$ | No |
| $s_3 = \{sa\}$ | No | $s_3 = \{sa\}$ | No |
| $s_4 = \{pe\}$ | Yes | $s_4 = \{pe\}$ | Yes |
| $s_5 = \{sa, eu\}$ | Yes | $s_5 = \{sa, eu\}$ | No |
| $s_6 = \{pe, eu\}$ | Yes | $s_6 = \{pe, eu\}$ | Yes |
| $s_7 = \{pe, sa\}$ | Yes | $s_7 = \{pe, sa\}$ | Yes |
| $s_8 = \{pe, sa, eu\}$ | Yes | $s_8 = \{pe, sa, eu\}$ | Yes |

Bob learns that the classifier is biased thereby being able to conclude that unsuccess of his application is (abductively) explained by $\neg pe \wedge \neg eu$

$$(s_3, f_1) \models [\text{GBias!}] \text{SubAXp}(\neg pe \wedge \neg eu, \text{No})$$

Axiomatics for the static setting *plus* the following valid equivalences:

$$[\varphi!]p \leftrightarrow (\Box_I \varphi \rightarrow p)$$

$$[\varphi!]t(x) \leftrightarrow (\Box_I \varphi \rightarrow t(x))$$

$$[\varphi!]\neg\psi \leftrightarrow (\Box_I \varphi \rightarrow \neg[\varphi!]\psi)$$

$$[\varphi!](\psi_1 \wedge \psi_2) \leftrightarrow ([\varphi!]\psi_1 \wedge [\varphi!]\psi_2)$$

$$[\varphi!]\Box_I \psi \leftrightarrow (\Box_I \varphi \rightarrow \Box_I [\varphi!]\psi)$$

$$[\varphi!]\Box_F \psi \leftrightarrow (\Box_I \varphi \rightarrow \Box_F [\varphi!]\psi)$$

and the following rule of replacement of equivalents:

$$\frac{\varphi_1 \leftrightarrow \varphi_2}{\psi \leftrightarrow \psi[\varphi_1/\varphi_2]}$$

\Rightarrow **Decidability** via the reduction axioms

- 1 Explanations in “white box” classifiers
- 2 Explanations in “black box” classifiers
- 3 Open problems and future extensions

- Exact complexity of satisfiability checking for the logic of “black box” classifiers
- Identify interesting NP fragments:
 - Bounding modal depth
 - Single alternation of \Box_F/\Diamond_F and \Box_I/\Diamond_I modalities sufficient for defining subjective explanation
- Complexity of dynamic extension

Definition

A **multi-classifier model with ideality (MCMI)** is a triple $\Gamma = (S, \Phi, \preceq)$ with (S, Φ) an MCM and \preceq a partial preorder on Φ .

$f \preceq f'$: classifier f' is at least as good/ideal as classifier f

“Betterness” modality $[\preceq]$ interpreted wrt pointed MCMI (Γ, s, f) :

$$(\Gamma, s, f) \models [\preceq]\varphi \iff \forall f' \in \Phi, \text{ if } f \preceq f' \text{ then } (\Gamma, s, f') \models \varphi$$

Expressive power:

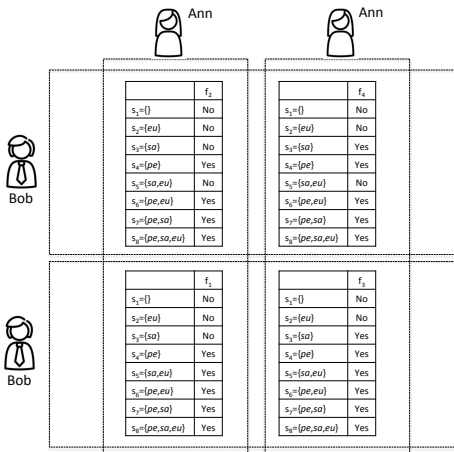
- Obligation modality:

$$\text{Oblig } \varphi =_{\text{def}} \Diamond_F [\preceq] \varphi$$

- Prohibition to have biases:

$$\text{Oblig } \neg \text{GBias}$$

Multi-agent generalization and interactive explanation



Formal semantics: multi-agent belief bases [Lorini, 2020, AIJ]

Liu, X., Lorini, E. (forthcoming). A unified logical framework for explanations in classifier systems. *Journal of Logic and Computation*.

Liu, X., Lorini, E., Rotolo, A., Sartor, G. (forthcoming). Modelling and explaining legal case-based reasoners through classifiers. Proceedings of *JURIX 2022*, IOS Press.

Liu, X., Lorini, E. (2022). A logic of “black box” classifier systems. Proceedings of *WOLLIC 2022*, Springer.

Liu, X., Lorini, E. (2021). A logic for binary classifiers and their explanation. Proceedings of *CLAR 2021*, Springer.