Interactive Explanation and Elicitation for Multiple Criteria Decision Analysis

Vincent Mousseau & Wassila Ouerdane

Laboratoire Génie Industriel

In collaboration with : Kh. Belahcene (LGI), Ch. Labreuche (Thales) and N. Maudet (LIP6)



IRT SystemX-April 11, 2018

Contents

Motivations

Introduction to Multiple Criteria Decision Aiding Basic MCDA concepts Preference Elicitation

Explanation schemes in MCDA context Pairwise comparisons Ordinal Sorting

Future prospects and applications

- ► new regulations (eg. GDPR)
- ► raising concern in the society : making A.I. systems trustable !

Featured in mainstream press, related to prominent applications :

- automated decisions for autonomous vehicles
- ► loan agreements
- Admission Post Bac/ParcourSup



- ► new regulations (eg. GDPR)
- ► raising concern in the society : making A.I. systems trustable !

Featured in mainstream press, related to prominent applications :

- automated decisions for autonomous vehicles
- ► loan agreements
- ► Admission Post Bac/ParcourSup



- ► new regulations (eg. GDPR)
- ► raising concern in the society : making A.I. systems trustable !

Featured in mainstream press, related to prominent applications :

- automated decisions for autonomous vehicles
- loan agreements
- Admission Post Bac/ParcourSup



General Data Protection Regulation : A right to explanation?

However, in their examination of the legal status of the GDPR, Wachter et al. conclude that such a right does not exist yet. The right to explanation is only explicitly stated in a recital :

a person who has been subject to automated decision-making "should be subject to suitable safeguards, which should include specific information to the data subject and the right to obtain human intervention, to express his or her point of view, to obtain an explanation of the decision reached after such assessment and to challenge the decision "

However, recitals are not legally binding. It also appears to have been intentionally not included in the final text of the GDPR after appearing in an earlier draft.

General Data Protection Regulation : A right to explanation?

Still, Article 13 and 14 about notification duties may provide a right to be informed about the "logic involved" prior to decision

"existence of automated decision-making, including profiling [...] [and provide data subjects with] meaningful information about the logic involved, as well as the significance and the envisaged consequences of such processing."

As it stands, only provides a (limited : secret of affairs, etc.) right to obtain ex-ante explanations about the model (which they call, 'right to be informed').

Wachter et al. Why a Right to Explanation of Automated Decision-Making Does Not Exist in the General Data Protection Regulation. International Data Privacy Law, 2017.

Loi pour une république numérique

L'administration communique à la personne faisant l'objet d'une décision individuelle prise sur le fondement d'un traitement algorithmique, à la demande de celle-ci, sous une forme intelligible et sous réserve de ne pas porter atteinte à des secrets protégés par la loi, les informations suivantes :

- Le degré et le mode de contribution du traitement algorithmique à la prise de décision;
- Les données traitées et leurs sources;
- ► Les paramètres de traitement et, le cas échéant, leur pondération, appliqués à la situation de l'intéressé;
- ► Les opérations effectuées par le traitement.

Décret du 14 Mars 2017, cité et commenté dans :

Besse et al.. Loyauté des Décisions Algorithmiques. Contribution to CNIL debate, 2017.

TRANSPARENCY, INTERPRETABILITY OR EXPLAINABILITY?

According to Besse et al., a decision can be said to be :

- transparent when the algorithm/code are made available.
- interpretable when it is possible to identify the features or variables which were prominent for the decision (even sometimes quantify this importance)
- explainable when it is possible to explicitly relate the values taken by the input data and the taken decision

Besse et al.. *Loyauté des Décisions Algorithmiques*. Contribution to CNIL debate, 2017 (my translation).

Introduction to Multiple Criteria Decision Aiding

Explanation schemes in MCDA context

Future prospects and applications

TRANSPARENCY DOES NOT IMPLY EXPLAINABILITY

1 (14	mbda	· _ / / / / / / / / /
2	get	attr(
3		import_(Trueclassname_[] + []classname_[_]),
4		(). class . eq . class . name [:] +
5		(). iter (). class . name [:]
6)(
7		, (lambda , , ; (, ,))(
8		lambda , , ;
9		chr(3) + (, , //) if else
10		(lambda:).func code.co lnotab.
11		« .
12		$(((\le) +) \le ((\le) -)) + ((((\le) +)))$
13		-) <<) +) << ((<<) + (<<))) + (((<<
14		() -) << (((((<<) +)) <<) + (<<))) + (((<<))) + (((<<))) + (((<<))) + (((<<))) + (((<<))) + (((<<))) + (((<<))) + (((<<))) + (((<<))) + ((<<))) + (((<<))) + ((<<))) + ((<<)) + ((<<))) + ((<<)) + ((<<))) + ((<<)) + ((<<))) + ((<<)) + ((<<)) + ((<<))) + ((<<)) + ((<<))) + ((<<)) + ((<<))) + ((<<)) + ((<<))) + ((<<)) + ((<<))) + ((<<)) + ((<<))) + ((<<)) + ((<<))) + ((<<)) + ((<<))) + ((<<)) + ((<<))) + ((<<)) + ((<<))) + ((<<))) + ((<<)) + ((<<))) + ((<<)) + ((<<))) + ((<<)) + ((<<))) + ((<<)) + ((<<)) + ((<<))) + ((<<)) + ((<<))) + ((<<)) + ((<<))) + ((<<)) + ((<<))) + ((<<)) + ((<<))) + ((<<)) + ((<<))) + ((<<)) + ((<<))) + ((<<)) + ((<<))) + ((<<)) + ((<<))) + ((<<)) + ((<<))) + ((<<)) + ((<<))) + ((<<)) + ((<<))) + ((<<)) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<
15		<pre><<) +) << ((<<) +)) + (((<<) -) <</pre>
16		$((\ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ $
17		(1 - 1) - (1 - 1) + (1 -
18		$\ll (((((\ll) +)) \ll))) - ((((((\ll) +)) \ll) +)))$
19		$) \ll ((((\ll) +) \ll))) + (((\ll) -) \ll))$
20		(((((<<) +)) <<))) + (((<<) +) << ((<<)))) + (((<<) +) << ((<<)))) + (((<<) +) << ((<<)))) + ((<<) +) << ((<<)))) + ((<<) +) << ((<<)))) + ((<<) +) << ((<<)))) + ((<<) +) << ((<<)))) + ((<<) +) << ((<<)))) + ((<<) +) << ((<<)))) + ((<<) +) << ((<<)))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))) + ((<<))
21		_))) + (<<) + (_<<)
22)
23)	
24)(
25	*(1	ambda _,,: _(_,,))(
26		(lambda , , ;
27		[([(lambda:).func code.co nlocals])] +
28		(, , [(lambda :).func code.co nlocals:]) if else []
29),
30		lambda _:func_code.co_argcount,
31		(
32		lambda _: _,
33		lambda _,: _,
34		lambda _,,: _,
35		lambda _,,;,
36		lambda _,,,,; _,
37		lambda _,,,,,; _,
38		lambda _,,,,,,,; _,
39		lambda _, _,,,,,,,
40)
41)	
42)		

atroduction to Multiple Criteria Decision Aiding

Explanation schemes in MCDA context

Future prospects and applications

TRANSPARENCY DOES NOT IMPLY EXPLAINABILITY

getattr(().__class__.__eq__.__class__.__name__[:__] + ().__iter__().__class_.__name__[____:___])(_, (lambda _, __, __: _(_, __, __))(lambda _, __, __: chr(__ % __) + _(_, __, __ // __) if __ else 10 (lambda: _).func_code.co_lnotab, 11 << ____ 12 (((____ << ___) + _) << ((___ << ___) - ___)) + (((((___ << __) - _) << __) + _) << ((____ << ___) + (_ << _))) + (((____ _) - _) << ((((((_ << __) + _)) << __) + (_ << _))) + (((_ 14 <<) +) << ((<<) +)) + (((<<) -) << 16 ((____ << __))) + (((_ << ___) - _) << ((((___ << __) + _) << 17 _) - _)) - (_____ << (((((__ << _) - _) << _) + _)) + (_ 18 << ((((((_ << __) + _)) << _))) - (((((((_ << __) + _)) << _) + 19 _) << ((((___ << __) + _) << _))) + (((____ << __) - _) << 20 ((((((_ << __) + _)) << _))) + (((__ << __) + _) << ((___ << 21 _))) + (____ << ___) + (_ << __) 22 23 241(25 *(lambda _, __, __: _(_, __, __))(26 (lambda _, __, __: [__(__[(lambda: _).func_code.co_nlocals])] + 27 _(_, __, __[(lambda : _).func_code.co_nlocals:]) if ___ else [] 28 29 30 lambda _: _.func_code.co_argcount, 31 32 lambda _: _, 33 lambda _, __: _, 34 lambda _, __, __: _, 35 lambda _, __, __, ___; _, lambda _, __, __, ___, ___, ___; _, 36 37 lambda _, __, ___, ___, ____, ____; _, 38 39 41 42)

prints Hello World! (by Ben Kurtovic, winner of a 2017 obfuscation contest)

A panel of questions we need to answer?

- 1. what were the main factors in a decision?
- 2. would changing a given factor have changed the decision?
- 3. how to improve the decision?
- 4. why did two similar-looking cases get different conclusions, or vice-versa?
- 5. does the model indeed do what is expected?
- 6. why this decision (recommendation)?
- 7. ...

▶ ...

The explanation landscape is rich already

Examples of approaches

- ► data-based explanations (incl. counterfactuals) [Datta et al., 2016]
- ► locally faithful approximations (LIME), surrogate models [Ribeiro et al, 2016]
- ► add constraints or objective (capturing interpretability) [Sokolovska et al., 2017];
- restrict operators to argumentation schemes validated by the user. [Belahcène et al., 2017]

Datta et al.. Algorithmic transparency via quantitative input influence : Theory and experiments with learning systems. The 37th IEEE Symposium on Security and Privacy.2016.

Ribeiro et al.. "why should i trust you?" Explaining the predictions of any classifier. In ACM SIGKDD International Conference on Knowledge Discovery and Data Mining.2016.

Sokolovska et al.. The fused lasso penalty for learning interpretable medical scoring systems.2017. IJCNN.

Belahcène et al.. Explaining robust additive utility models by sequences of preference swaps. Theory and Decision. 2017.

The explanation landscape is rich already

Examples of approaches

- ► data-based explanations (incl. counterfactuals) [Datta et al., 2016]
- ► locally faithful approximations (LIME), surrogate models [Ribeiro et al, 2016]
- add constraints or objective (capturing interpretability) [Sokolovska et al., 2017];
- restrict operators to argumentation schemes validated by the user. [Belahcène et al., 2017]
- ▶ ...

An explanation (argumentation) scheme

an operator tying a *tuple of premises* (pieces of information provided or approved by the Decision Maker, or inferred during the process, and some supplementary hypotheses on the reasoning process (model's assumptions) to *a conclusion*.

Contents

Motivations

Introduction to Multiple Criteria Decision Aiding Basic MCDA concepts Preference Elicitation

Explanation schemes in MCDA context Pairwise comparisons Ordinal Sorting

Future prospects and applications

OUR CONTEXT : MULTIPLE CRITERIA DECISION AIDING



PAIRWISE COMPARISONS (CHOICE OR RANKING)



Ordinal sorting

Object	а	b	с	d	Assignment
A ₁	3	3	2.5	0	***
A_2	3	2	2.1	1	***
B_1	2	2	1.3	1	**
B_2	3	1	3.7	0	**
C_1	2	1	1.6	1	*
C_2	1	1	4.1	0	*
Х	2	2	1.1	0	?

Decision Maker

OUR CONTEXT : MULTIPLE CRITERIA DECISION AIDING

Assumes a **preference model** containing **aggregation procedures**

- ▶ mapping feature profiles to recommendations.
- extending Pareto dominance and expressed preferences.
- ► implementing a decision theoretic stance.

Dominance, Pareto-optimality

- ► Consider $a = (a_1, a_2, ..., a_n)$, $b = (b_1, b_2, ..., b_n)$, $a\Delta b$ iff $a_j \ge b_j$, $\forall j = 1..n$, one of the inequalities being strict,
- ► The dominance relation ∆ expresses unanimity among criteria in favor of one action in the comparison,
- Δ defines on A strict partial order (asymmetric and transitive),
- $\blacktriangleright~\Delta$ is usually very poor,
- ► $a \in A$ is Pareto-optimal iff $\nexists b \in A$ s.t. $b\Delta a$,

ъ.	٤	1 Å		12.5		
iv	10	u١	(2	LIC	011:	S

Introduction to Multiple Criteria Decision Aiding

Explanation schemes in MCDA context

Future prospects and applications

Pareto front

Pareto front in a discret bi-criteria problem

Pareto front

Preference information

- $\blacktriangleright\,$ To discriminate among Pareto-optimal alternatives, the dominance relation Δ is useless,
- ► Decision aiding requires to enrich ∆ by additional information called preference information,
- Preference information refers to the DM's opinions, value system, convictions ... concerning the decision problem,
- ► It is standard to distinguish :
 - Intracriterion preference information, and
 - ► Intercriteria preference information.

MCDA

Model selection

- a preference model contains aggregation procedures satisfying common properties.
- ▶ a model is selected considering decision stance, expressiveness, tractability.

Additive Utility Model

▶ preference derives from a value model

$$\exists V \, s.t. \, x \succeq y \iff V(x) \ge V(y)$$

• value is additive (i.e. $V(x) = \sum_i v_i(x_i)$)

NonCompensatory Sorting Model

► pairwise comparisons preferences

$$NCS_{S,\langle A_i \rangle}(x) = \begin{cases} \mathcal{GOOD}, \text{ if } \{i \in \mathcal{N} : x \in A_i\} \in S \\ \mathcal{BAD}, \text{ else} \end{cases}$$

MCDA

Model elicitation

• Once a model is selected, a specific decision procedure has to be determined.

▶ preference information is collected from the Decision Maker, then processed.

Approach	Summary	Pros	Cons
Complete	Standard sequence of questions	Unequivocal	Demanding
Partial	Learning from DM's statements + Loss function	Efficient	Arbitrary
Robust	Partial + Accounting for possible completions	Cautious	Indecisive
Active	Dynamically determined queries minimizing regret	Fast	Arbitrary

Contents

Motivations

Introduction to Multiple Criteria Decision Aiding Basic MCDA concepts Preference Elicitation

Explanation schemes in MCDA context Pairwise comparisons Ordinal Sorting

Future prospects and applications

Introduction to Multiple Criteria Decision Aiding

Explanation schemes in MCDA context

Future prospects and applications

PREFERENCE ELICITATION

Introduction to Multiple Criteria Decision Aiding

Explanation schemes in MCDA context

Future prospects and applications

PREFERENCE ELICITATION

Introduction to Multiple Criteria Decision Aiding ○○○○○○○○○○○○○○○● Explanation schemes in MCDA context

Future prospects and applications

PREFERENCE ELICITATION

OUR CONTEXT : MULTIPLE CRITERIA DECISION AIDING

Contents

Motivations

Introduction to Multiple Criteria Decision Aiding Basic MCDA concepts Preference Elicitation

Explanation schemes in MCDA context Pairwise comparisons Ordinal Sorting

Future prospects and applications

introduction to Multiple Criteria Decision Aiding

Explanation schemes in MCDA context

Future prospects and applications

I want to compare hotels described by 4 criteria : comfort (A), parking

(B), commute time (C), and cost (D).

I prefer :

(4*, no, 15 min, 150 \$) to (2*, yes, 45 min, 50 \$), (2*, no, 45 min, 50 \$) to (2*, yes, 15 min, 150 \$), (2*, yes, 15 min, 150 \$) to (4*, no, 45 min, 150 \$).

ANALYST

I want to know : Is (2*, no, 15 min, 50 \$) better than (4*, yes, 45 min, 150 \$)?

Assumptions :

Decision Maker

- ▶ preference derives from a value model (i.e. $\exists V \text{ s.t. } x \succeq y \iff V(x) \ge V(y)$)
- value is additive (i.e. $V(x) = \sum_i v_i(x_i)$)

Decision Maker

introduction to Multiple Criteria Decision Aiding

Explanation schemes in MCDA context

Future prospects and applications

I want to compare hotels described by 4 criteria : comfort (A), parking

(B), commute time (C), and cost (D).

I prefer :

(4*, no, 15 min, 150 \$) to (2*, yes, 45 min, 50 \$), (2*, no, 45 min, 50 \$) to (2*, yes, 15 min, 150 \$), (2*, yes, 15 min, 150 \$) to (4*, no, 45 min, 150 \$).

I want to know : Is (2*, no, 15 min, 50 \$) better than (4*, yes, 45 min, 150 \$)? ANALYST

Ordinal encoding : attribute values of interest are sorted and encoded

criterion A : $4\star$ is strong (**O**), $2\star$ is weak (**O**); criterion B : yes is strong (**O**), no is weak (**O**); criterion C : 15 min is strong (**O**), 45 min is weak (**O**); criterion D : 50 \$ is strong (**O**), 150 \$ is weak (**O**).

Streamlining the Robust Additive Value model

Explaining with sequences of preference swaps

- Assuming the complexity of preference stems from having many moving parts
- Decomposing the complexity into smaller grains by reasoning ceteris paribus

se explanations can be long, but can be kept short and computed efficiently when constraining the PI

Belahcène et al. *Explaining robust additive utility models by sequences of preference swaps.* Theory and Decision. 2017.

Contents

Motivations

Introduction to Multiple Criteria Decision Aiding Basic MCDA concepts Preference Elicitation

Explanation schemes in MCDA context

Pairwise comparisons Ordinal Sorting

Future prospects and applications

NONCOMPENSATORY SORTING PROCEDURE

Output

• A category among an ordered set $C_1 \prec \cdots \prec C_p$

Sorting rule

▶ an alternative is in category C_h or better iff it has sufficient attributes at level C_h on a coalition of criteria deemed sufficient at level C_h

History

- ▶ inspired by Electre Tri
- described and characterized in [Bouyssou & Marchant, 2007 ab]
- equivalent to the Sugeno integral model [Slowinski et al., 2002]

Particular cases

- ▶ **U** : using a Unique set of sufficient coalitions of criteria
- ▶ V : representing sufficient coalitions with a Voting model
- ▶ We call NCS models following U "U-NCS", U&V "MR-Sort" [Leroy et al., 2011]

NonCompensatory Sorting Example

Project	а	b	с	d	Category
p_1	5	6	6	5	?
p_2	3.5	1	3	9	?
p_3	7.5	2	1	3	?
p_4	2	8	2.5	7	?
p_5	3	8.5	3	8.5	?
p_6	8	4	1.5	1.5	?

*	< 4	< 3	< 2	< 2	boundary between \star and $\star\star$
**	[4,7[[3,8[[2,5[[2,8[
***	≥ 7	≥ 8	≥ 5	≥ 8	boundary between $\star\star$ and $\star\star\star$

NonCompensatory Sorting Example

1st phase : criterion-wise sorting

project	а	b	с	d	Category
p_1	**	**	***	**	?
p_2	*	*	**	***	?
p_3	***	*	*	**	?
p_4	*	***	**	**	?
p_5	*	***	**	***	?
p_6	***	**	*	*	?

*	< 4	< 3	< 2	< 2	boundary between \star and $\star\star$
**	[4,7[[3,8[[2,5[[2,8[
***	≥ 7	≥ 8	≥ 5	≥ 8	boundary between ** and ***

NonCompensatory Sorting Example

2nd phase : noncompensatory multi criteria aggregation

Insufficient coalitions

▶ Getting an overall ★★ or ★★★ requires getting ★★ or ★★★ on a sufficient coalition of criteria

Getting an overall * * * requires getting * * * on a sufficient coalition of criteria

Explanation schemes in MCDA context

NonCompensatory Sorting Example

2nd phase : noncompensatory multi criteria aggregation

Insufficient coalitions

▶ Getting an overall ★★ or ★★★ requires getting ★★ or ★★★ on a sufficient coalition of criteria

Getting an overall * * * requires getting * * * on a sufficient coalition of criteria

Explanation schemes in MCDA context

Future prospects and applications

Learning / Disaggregation of U-NCS model

Input : profiles + reference assignments

*/ **	?	?	?	?	
**/ * * *	?	?	?	?	

Sufficient coalitions

Insufficient coalitions

B Expected Outputs : set of profiles + set of sufficient coalitions.

Learning / Disaggregation of U-NCS model

- Direct elicitation with standard sequence procedures
- Computational issues with the indirect elicitation of MR Sort (learning from assignment examples) :
 - ▶ with a MIP [Leroy et al, 2011] : hardly more than toy examples
 - ► with a meta-heuristic [Sobrie et al, 2016] : learning sets from preference learning
- ► issues with knowledge representation
 - dependencies between profiles and coalitions are non-trivial
 - the profiles part seems to fall within the domain of 'logical inference'
 - the coalition part is described by linear programming
 - + need for a unified description : back to NCS (alternate solution : MR Sort + cutting planes?)

A COMPACT SAT FORMULATION

Let $\alpha:\mathbb{X}\to\{\text{GOOD},\text{BAD}\}$ an assignment. We define the boolean function $\phi_\alpha^{\textit{pairwise}}$ with variables :

- $\lambda_{i,x}$ indexed by a point of view $i \in \mathcal{N}$, and a value $x \in \mathbb{X}$,
- $\mu_{i,g,b}$ indexed by a point of view $i \in \mathcal{N}$, a *good* alternative $g \in \alpha^{-1}(\text{Good})$ and a *bad* alternative $b \in \alpha^{-1}(\text{BAD})$,

as the conjunction of clauses : $\phi^{\textit{pairwise}}_\alpha:=\phi^1_\alpha\wedge\phi^2_\alpha\wedge\phi^3_\alpha\wedge\phi^4_\alpha$

$$\begin{split} \phi_{\alpha}^{1} &:= & \bigwedge_{i \in \mathcal{N}} \bigwedge_{x' \succeq ix} & (\lambda_{i,x'} \lor \neg \lambda_{i,x}) \\ \phi_{\alpha}^{2} &:= & \bigwedge_{i \in \mathcal{N}, \ g \in \alpha^{-1}(\text{Good}), \ b \in \alpha^{-1}(\text{BAD})} & (\neg \mu_{i,g,b} \lor \neg \lambda_{i,b}) \\ \phi_{\alpha}^{3} &:= & \bigwedge_{i \in \mathcal{N}, \ g \in \alpha^{-1}(\text{Good}), \ b \in \alpha^{-1}(\text{BAD})} & (\neg \mu_{i,g,b} \lor \lambda_{i,g}) \\ \phi_{\alpha}^{4} &:= & \bigwedge_{g \in \alpha^{-1}(\text{Good}), \ b \in \alpha^{-1}(\text{BAD})} & (\bigvee_{i \in \mathcal{N}} \mu_{i,g,b}) \end{split}$$

Towards explanations for NCS $% \left({{{\rm{NCS}}}} \right)$

Situation 1 : Auditing conformity

An independent audit agency is commissioned to check that the decision on the the committee indeed comply with a publicly announced decision rule.

Situation 1 : Auditing conformity

An independent audit agency is commissioned to check that the decision on the the committee indeed comply with a publicly announced decision rule.

so computing and providing a certificate of feasibility of a SAT problem.

Situation 1 : Auditing conformity

An independent audit agency is commissioned to check that the decision on the the committee indeed comply with a publicly announced decision rule.

so computing and providing a certificate of feasibility of a SAT problem.

Situation 2 : Justifying individual decisions

A candidate, (supposedly) unsatisfied with the outcome of the process regarding his own classification, challenged the committee and asks for justification.

- necessary decisions entailed by the jurisprudence.
- Ambivalent situations.

Situation 1 : Auditing conformity

An independent audit agency is commissioned to check that the decision on the the committee indeed comply with a publicly announced decision rule.

so computing and providing a certificate of feasibility of a SAT problem.

Situation 2 : Justifying individual decisions

A candidate, (supposedly) unsatisfied with the outcome of the process regarding his own classification, challenged the committee and asks for justification.

- necessary decisions entailed by the jurisprudence.
- Ambivalent situations.

computing and providing a certificate of infeasibility (MUS)

Open issues :

- ► How do we leverage this description inside a decision process?
- ► Can we build explanations around certificates of UNSAT (MUSes)?
 - ► What is a "good" certificate?
 - Can we find a template (=argument schemes) in which they fit? (all of them? some of them?)
 - Can we compute them effectively?

FUTURE PROSPECTS AND APPLICATIONS

Open issues

- Intégration de l'explication et de l'élicitation dans un mécanisme dialectique (gestion de l'inconsitance, choix de modèle, protocole de dialogue, etc.)
 - ▶ PEPS "PULP" (S. Destercke, Heudiasyc Lip6)
 - Propale ANR 2018 IRELAND" (V. Mousseau / W. Ouerdane, LGI LIP6 -LAMSADE- IMT Atlantique)
- Encodages et méthodes SAT pour la production d'explications.
 - ▶ PEPS "SAT4EX" (N. Maudet, Lip6 CRIL)

• ...

Different application domains

- Configuration problem;
- Recommendation problem;
- Administrative decisions;

^{• ...}