

# 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)



CentraleSupélec

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

# MOTIVATIONS

- ▶ 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

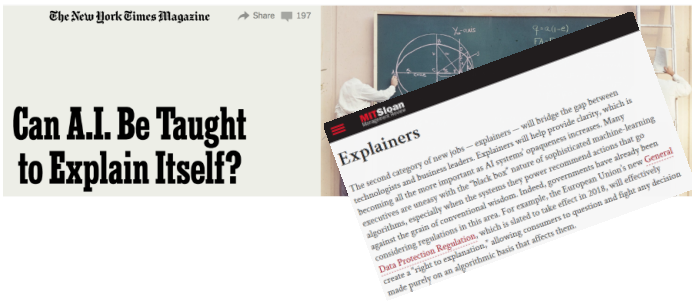


# MOTIVATIONS

- ▶ 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

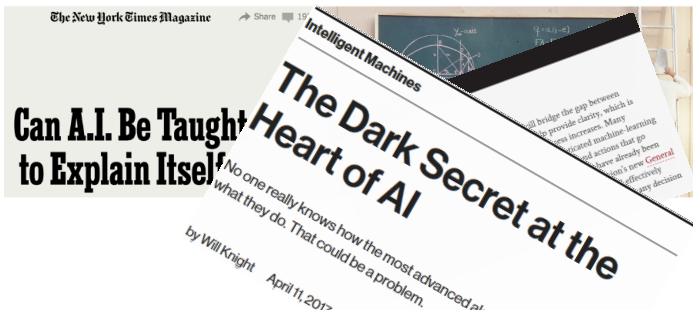


# MOTIVATIONS

- ▶ 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).







## 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



Decision  
Maker

A **performance table**, describing several actions according to various criteria - the higher the better

A **decision problem** : Is action A better than action B? Is action C good enough?

Sparse **preferences** between some actions



## PAIRWISE COMPARISONS (CHOICE OR RANKING)

I want to compare hotels described by 4 criteria :

$A$ -	comfort :	(4 <sup>*</sup> )	$A \succ a$	(2 <sup>*</sup> )
$B$ -	restaurant :	(presence)	$B \succ b$	(absence)
$C$ -	commute time :	(15 min)	$C \succ c$	(45 min)
$D$ -	cost :	(50 \$)	$D \succ d$	(150 \$)

I **prefer**  $[AbCd]$  to  $[aBcD]$ ,  $[abcD]$   
to  $[ABcd]$  and  $[aBCd]$  to  $[Abcd]$

I want to know :  
Is  $[abCD]$  better than  $[ABcd]$ ?



Decision  
Maker

# ORDINAL SORTING

Object	a	b	c	d	Assignment
A <sub>1</sub>	3	3	2.5	0	***
A <sub>2</sub>	3	2	2.1	1	***
B <sub>1</sub>	2	2	1.3	1	**
B <sub>2</sub>	3	1	3.7	0	**
C <sub>1</sub>	2	1	1.6	1	*
C <sub>2</sub>	1	1	4.1	0	*
X	2	2	1.1	0	?



**Decision  
Maker**

**What class should I assign to X?**

# OUR CONTEXT : MULTIPLE CRITERIA DECISION AIDING



Decision  
Maker

A **performance table**, describing several actions according to various criteria - the higher the better

A **decision problem** : Is action A better than action B? Is action C good enough?

Sparse **preferences** between some actions

A **recommendation**



**ANALYST**

# OUR CONTEXT : MULTIPLE CRITERIA DECISION AIDING



**Decision  
Maker**

A **performance table**, describing several actions according to various criteria - the higher the better

A **decision problem** : Is action A better than action B? Is action C good enough?

Sparse **preferences** between some actions

**A recommendation**

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

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



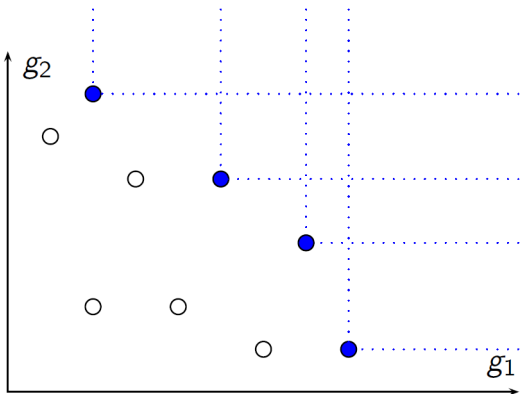
**Analyst**



# DOMINANCE, PARETO-OPTIMALITY

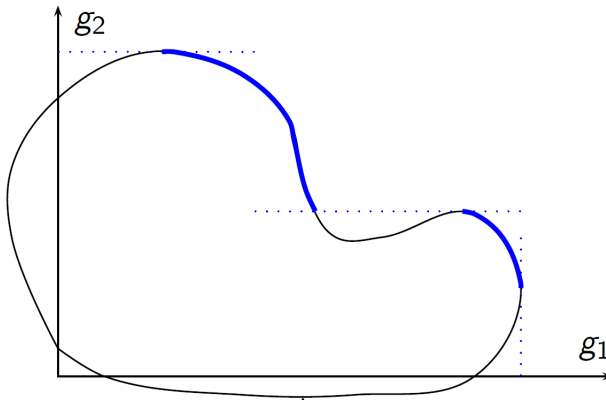
- ▶ Consider  $a = (a_1, a_2, \dots, a_n)$ ,  $b = (b_1, b_2, \dots, b_n)$ ,  
 $a \Delta b$  iff  $a_j \geq b_j, \forall j = 1..n$ , one of the inequalities being strict,
- ▶ The dominance relation  $\Delta$  expresses unanimity among criteria in favor of one action in the comparison,
- ▶  $\Delta$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

# PARETO FRONT

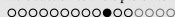


Pareto front in a discret bi-criteria problem

# PARETO FRONT



Pareto front in a continuous bi-criteria problem



# PREFERENCE INFORMATION

- ▶ To discriminate among Pareto-optimal alternatives, the dominance relation  $\Delta$  is useless,
- ▶ Decision aiding requires to enrich  $\Delta$  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

$$\exists V \text{ s.t. } x \succsim y \iff V(x) \geq V(y)$$

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

### NonCompensatory Sorting Model

- ▶ pairwise comparisons preferences

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

# 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**

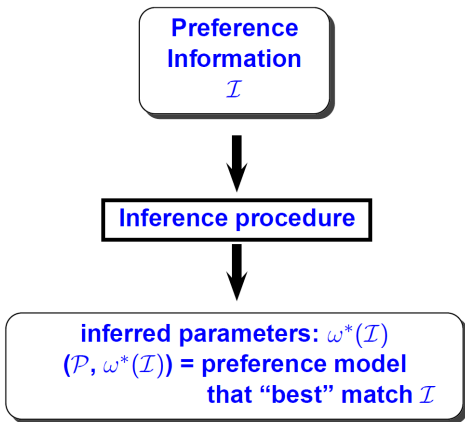
Explanation schemes in MCDA context

Pairwise comparisons

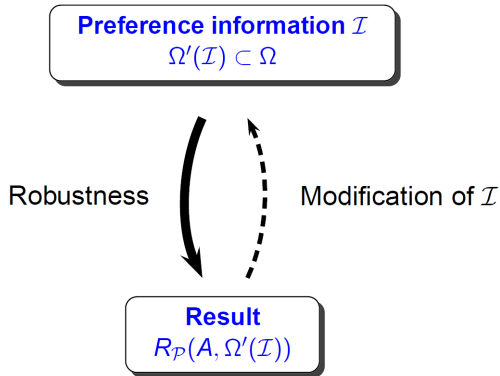
Ordinal Sorting

Future prospects and applications

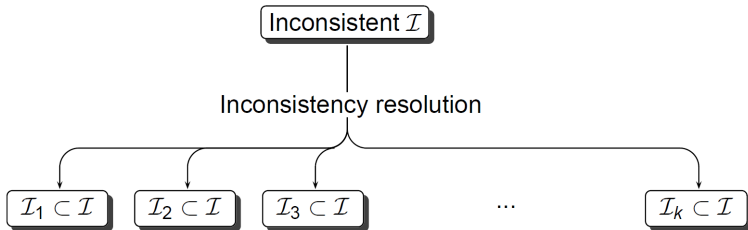
# PREFERENCE ELICITATION



# PREFERENCE ELICITATION



# PREFERENCE ELICITATION



# OUR CONTEXT : MULTIPLE CRITERIA DECISION AIDING

A **performance table**, describing several actions according to various criteria - the higher the better

A **decision problem** : Is action A better than action B? Is action C good enough?

Sparse **preferences** between some actions

A **recommendation**

Why?

A Recommendation Supported with an **explanation** :

Expressed Preference + Model Features  $\Rightarrow$  Recommendation



Decision  
Maker



ANALYST

# 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



?



Decision Maker

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 \$) bet-  
ter than (4★, yes, 45 min, 150 \$)?

### Assumptions :

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

?



Decision Maker

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 \$) bet-  
ter than (4★, yes, 45 min, 150 \$)?

**Ordinal encoding** : attribute values of interest are sorted and encoded

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

# STREAMLINING THE ROBUST ADDITIVE VALUE MODEL



I prefer  $\overset{A}{\odot}\overset{B}{\odot}\overset{C}{\odot}\overset{D}{\odot}$  to  $\overset{A}{\odot}\overset{B}{\odot}\overset{C}{\odot}\overset{D}{\odot}$ ,  $\overset{A}{\odot}\overset{B}{\odot}\overset{C}{\odot}\overset{D}{\odot}$  to  $\overset{A}{\odot}\overset{B}{\odot}\overset{C}{\odot}\overset{D}{\odot}$  and  $\overset{A}{\odot}\overset{B}{\odot}\overset{C}{\odot}\overset{D}{\odot}$  to  $\overset{A}{\odot}\overset{B}{\odot}\overset{C}{\odot}\overset{D}{\odot}$

I want to know : Is  $\overset{A}{\odot}\overset{B}{\odot}\overset{C}{\odot}\overset{D}{\odot}$  better than  $\overset{A}{\odot}\overset{B}{\odot}\overset{C}{\odot}\overset{D}{\odot}$ ?

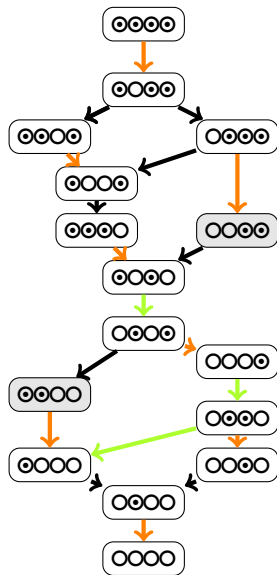
IT IS !

## Knowledge representation

- ▶ comparative pairwise statements are represented as inequations between elementary score differences
- ▶ Knowledge Base : Preference Information ( $\searrow$ ) + Pareto dominance ( $\searrow$ )

## Inference

- ▶ a comparative statement holds iff it is a conical combination of statements in the KB
- ▶ a finite number of inferred statements ( $\searrow$ ) are computed by Linear Programming



# 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*

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

I prefer [○○○○] to [○●○○], [○○○○] to [○○○○] and [○○○○] to [○○○○]  
I want to know : is [○○○○] better than [○○○○]?



Decision Maker

- IT IS ! HERE IS WHY :
1. [○, ○, ●, ○] IS BETTER THAN [●, ○, ●, ○]  
BECAUSE, EVERYTHING ELSE BEING EQUAL,  
[○, B, C, ●] (2★ FOR 50 \$) IS BETTER THAN  
[●, B, C, ○] (4★ FOR 150 \$).
  2. [●, ○, ●, ○] IS BETTER THAN [○, ●, ○, ○]  
BECAUSE, YOU TOLD ME SO!
  3. [●, ○, ●, ○] IS BETTER THAN [●, ●, ○, ○]  
BECAUSE, EVERYTHING ELSE BEING EQUAL,  
[A, ○, ●, D] (NO PARKING, 15 MIN COMMUTE) IS  
BETTER THAN [A, ●, ○, D] (PARKING, 45 MIN)



ANALYST

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  $C_1 \prec \dots \prec C_p$

## 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	a	b	c	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 ★ and ★★
★★	[4,7[	[3,8[	[2,5[	[2,8[	
★★★	≥ 7	≥ 8	≥ 5	≥ 8	boundary between ★★ and ★★★

# NONCOMPENSATORY SORTING EXAMPLE

## 1<sup>st</sup> phase : criterion-wise sorting

project	a	b	c	d	Category
$p_1$	**	**	***	**	?
$p_2$	*	*	**	***	?
$p_3$	***	*	*	**	?
$p_4$	*	***	**	**	?
$p_5$	*	***	**	***	?
$p_6$	***	**	*	*	?

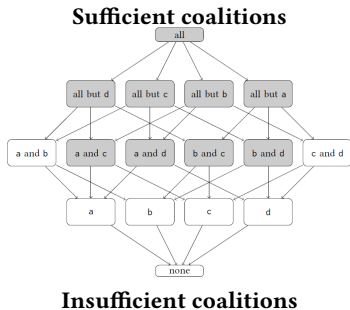
*	< 4	< 3	< 2	< 2	boundary between * and **
**	[4,7[	[3,8[	[2,5[	[2,8[	
***	≥ 7	≥ 8	≥ 5	≥ 8	boundary between ** and ***



# NONCOMPENSATORY SORTING EXAMPLE

2<sup>nd</sup> phase : noncompensatory multi criteria aggregation

project	a	b	c	d	Category
<i>p</i> <sub>1</sub>	**	**	***	**	?
<i>p</i> <sub>2</sub>	*	*	**	***	?
<i>p</i> <sub>3</sub>	***	*	*	**	?
<i>p</i> <sub>4</sub>	*	***	**	**	?
<i>p</i> <sub>5</sub>	*	***	**	***	?
<i>p</i> <sub>6</sub>	***	**	*	*	?

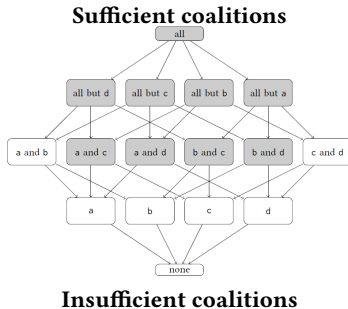


- ▶ Getting an overall \*\* or \*\*\* requires getting \*\* or \*\*\* on a sufficient coalition of criteria
- ▶ Getting an overall \*\*\* requires getting \*\*\* on a sufficient coalition of criteria

# NONCOMPENSATORY SORTING EXAMPLE

2<sup>nd</sup> phase : noncompensatory multi criteria aggregation

project	a	b	c	d	Category
<i>p</i> <sub>1</sub>	**	**	***	**	**
<i>p</i> <sub>2</sub>	*	*	**	***	*
<i>p</i> <sub>3</sub>	***	*	*	**	**
<i>p</i> <sub>4</sub>	*	***	**	***	**
<i>p</i> <sub>5</sub>	*	***	**	***	***
<i>p</i> <sub>6</sub>	***	**	*	*	*



- ▶ Getting an overall \*\* or \*\*\* requires getting \*\* or \*\*\* on a sufficient coalition of criteria
- ▶ Getting an overall \*\*\* requires getting \*\*\* on a sufficient coalition of criteria

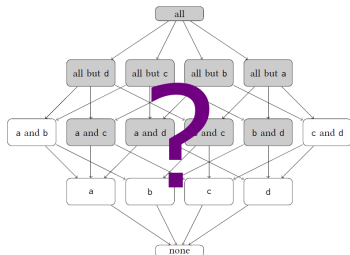
# LEARNING / DISAGGREGATION OF U-NCS MODEL

☞ **Input : profiles + reference assignments**

model	a	b	c	d	Category
$m_1$	16 973	29	2.66	2.5	★★
$m_2$	18 342	30.7	2.33	3	★
$m_3$	15 335	30.2	2	2.5	★★
$m_4$	18 971	28	2.33	2	★★
$m_5$	17 537	28.3	2.33	2.75	★★★
$m_6$	15 131	29.7	1.66	1.75	★

★/★★	?	?	?	?	
★★/★★★	?	?	?	?	

**Sufficient coalitions**



**Insufficient coalitions**

☞ **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  $\alpha : \mathbb{X} \rightarrow \{\text{GOOD}, \text{BAD}\}$  an assignment. We define the boolean function  $\phi_\alpha^{\text{pairwise}}$  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 :  $\phi_\alpha^{\text{pairwise}} := \phi_\alpha^1 \wedge \phi_\alpha^2 \wedge \phi_\alpha^3 \wedge \phi_\alpha^4$

$$\phi_\alpha^1 := \bigwedge_{i \in \mathcal{N}} \bigwedge_{x' \succsim_i x} (\lambda_{i,x'} \vee \neg \lambda_{i,x})$$

$$\phi_\alpha^2 := \bigwedge_{i \in \mathcal{N}, g \in \alpha^{-1}(\text{GOOD}), b \in \alpha^{-1}(\text{BAD})} (\neg \mu_{i,g,b} \vee \neg \lambda_{i,b})$$

$$\phi_\alpha^3 := \bigwedge_{i \in \mathcal{N}, g \in \alpha^{-1}(\text{GOOD}), b \in \alpha^{-1}(\text{BAD})} (\neg \mu_{i,g,b} \vee \lambda_{i,g})$$

$$\phi_\alpha^4 := \bigwedge_{g \in \alpha^{-1}(\text{GOOD}), b \in \alpha^{-1}(\text{BAD})} (\bigvee_{i \in \mathcal{N}} \mu_{i,g,b})$$

# TOWARDS EXPLANATIONS FOR NCS

## 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.

# TOWARDS EXPLANATIONS FOR NCS

## 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.

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

# TOWARDS EXPLANATIONS FOR NCS

## 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.

☞ 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.



# TOWARDS EXPLANATIONS FOR NCS

## 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.

☞ 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)

# TOWARDS EXPLANATIONS FOR NCS

## 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'inconsistance, 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 ;
- ▶ ...