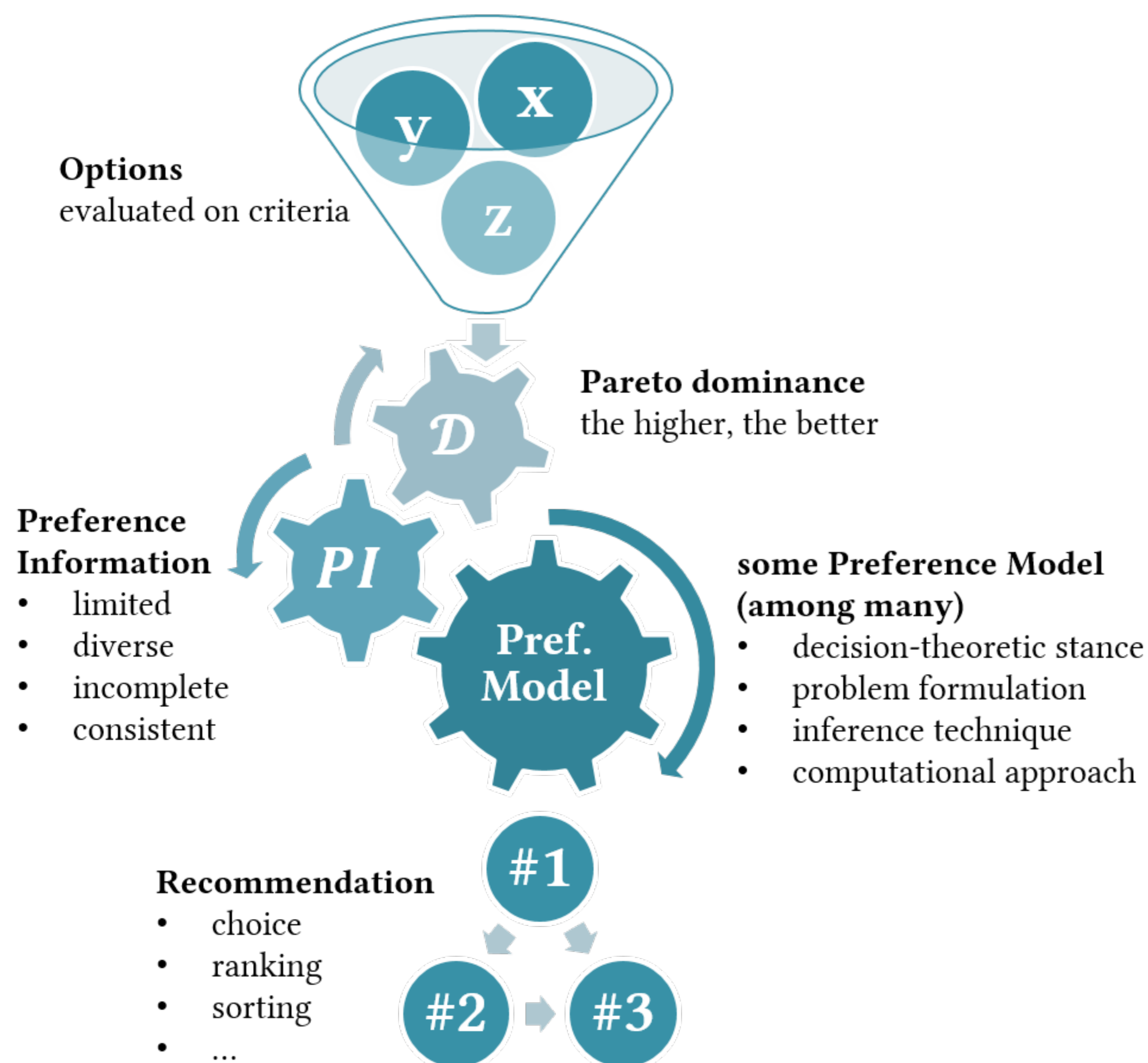


# Challenges in Explaining Decisions (and some links to Social Choice)

Khaled Belahcene<sup>1</sup>, C. Labreuche<sup>2</sup>, N. Maudet<sup>3</sup>, V. Mousseau<sup>1</sup> and W. Ouerdane<sup>1</sup>

## ■ MULTI-CRITERIA DECISION AIDING (MCDA)



## From interpretability...

- Axiomatized MCDA models claim "interpretability", but they are hardly intelligible by themselves;
- MCDA is structurally a "human in the loop" methodology and process. The Decision Maker's grasp of the stakes is crucial;
- Existing explanation frameworks, designed to complement Decision Support Systems non-specifically MCDA, are too lightweight.

## ■ LANDMARKS

1. I. Alvarez : *Explaining the result of a Decision Tree to the End-User*. ECAI 2004: 411-415  
 ☞ An explanation is more than an arbitrary trace of the decision process
2. S. Greco, V. Mousseau, and R. Slowinski. *Ordinal regression revisited: multiple criteria ranking using a set of additive value functions*. European Journal of Operational Research, 191(2):415-435, 2008.  
 ☞ Necessary and possible preference statements in an additive utility MCDA model
3. W. Ouerdane, N. Maudet, A. Tsoukiàs: *Argument Schemes and Critical Questions for Decision Aiding Process*. COMMA 2008: 285-296  
 ☞ An interactive MCDA framework based on argumentation techniques. Model selection is addressed from a user-centric perspective.
4. C. Labreuche: *A general framework for explaining the results of a multi-attribute preference model*. Artif. Intell. 175(7-8): 1410-1448 (2011)  
 ☞ A principled way of selecting arguments supporting decisions in MCDA models assigning weights to criteria
5. C. Labreuche, N. Maudet, W. Ouerdane: *Justifying Dominating Options when Preferential Information is Incomplete*. ECAI 2012: 486-491  
 ☞ An effective engine explaining necessary preference statements in a weighted Condorcet model, with duality techniques
6. Olivier Cailloux and Ulle Endriss, *Arguing about Voting Rules*, AAMAS-2016  
 ☞ Voting rules are promoted by exhibiting meaningful situations showcasing the underlying axioms
7. Spiegler, R. *Equilibrium in Justifiable Strategies: A Model of Reason-Based Choice in Extensive-Form Games*. The Review of Economic Studies, 69(3), 2002  
 ☞ In a strategic context, players are restricted to strategies they can account for, leading to a new definition of equilibrium

## Contributions

### ■ EXPLAINING ROBUST ADDITIVE UTILITY MODELS BY SEQUENCES OF PREFERENCE SWAPS, *THEORY AND DECISION*, IN PRINT

#### ► Problem statement

**PI** : a set of pairwise preference statements

**Model** : any satisfying Pareto, Transitivity and Cancellation axioms, e.g.

- any particular Additive Value model, i.e.  $x \succeq y \iff \sum V_i(x_i) \geq \sum V_i(y_i)$
- $x$  is necessarily preferred to  $y$  iff  $V(x) \geq V(y)$  for every possible Additive Value model correctly representing the PI.

**Recommendation** : a preference statement  $x \succeq y$

#### ► Proposed explanation

A sequence of options  $x = e_0 \succeq e_1 \succeq \dots \succeq e_{n-1} \succeq e_n = y$

- establishing the preference of  $x$  over  $y$  (transitivity)
- two adjacent options differ only on 1 (dominance) or 2 (trade-off) criteria

#### ► Results and Challenges

Existence ? Bound on sequence length ? Computation ?

- **Necessary Preference + binary PI** : Explanations have a term-by-term structure. Efficient algorithm for existence and actual computation. Explanations can be kept short. Proofs use PL/duality, and graph/flows techniques.

Explanations can be kept short. Proofs use PL/duality, and graph/flows techniques.

- **general case** : Open issues. We provide an example where there is no upper bound on the length of the shortest possible explanations

### ■ ACCOUNTABLE CLASSIFICATION WITHOUT FRONTIERS, DA2PL'16, SUBMITTED

#### ► Design principles favoring Accountability

No jargon. No values. No frontiers. No compensation. No inference.

- an object can not outrank any object assigned to a strictly better class;
- an object outranks objects assigned to a strictly worse class;

#### ► Implementation

- the model **observes** every pair of reference objects not assigned to the same class
- it **learns** sets of *sufficient*, *insufficient*, or undecided coalitions of criteria, accounting for monotonicity
- for a given candidate, it **recommends** every *possible* assignment not contradicting its principles
- it **explains** its recommendation with supporting statements instantiating specified *argument schemes*

Object	a	b	c	d	Assignment
A <sub>1</sub>	A	A	2.5	False	★★★
A <sub>2</sub>	A	B	2.1	True	★★★
B <sub>1</sub>	B	B	1.3	True	★★
B <sub>2</sub>	A	C	3.7	False	★★
C <sub>1</sub>	B	C	1.6	True	★
C <sub>2</sub>	C	C	4.1	False	★
Z <sub>1</sub>	B	B	1.1	False	?
Z <sub>2</sub>	B	A	1.8	False	??
Z <sub>3</sub>	A	B	1.2	False	???

	★★★	★★	★	?	?	?
	A <sub>1</sub> A <sub>2</sub>	B <sub>1</sub> B <sub>2</sub>	C <sub>1</sub> C <sub>2</sub>	Z <sub>1</sub>	Z <sub>2</sub>	Z <sub>3</sub>
A <sub>1</sub>		abc abd	abc abd	abcd	abcd	abcd
A <sub>2</sub>		abcd abd	abc abd	abcd	acd	abcd
B <sub>1</sub>	d bd		abd abd	abcd	ad	bcd
B <sub>2</sub>	acd ac		abc abd	cd	acd	acd
C <sub>1</sub>	d d	acd bd		acd	acd	cd
C <sub>2</sub>	cd c	c bcd		cd	cd	cd
Z <sub>1</sub>	d b	(ab) bd	(ab) abd			
Z <sub>2</sub>	bd bc	abc bd	(ab) abd			
Z <sub>3</sub>	(ab) (ab)	(ab) abd	(ab) abd			

For example  $Z_2$  should at least be assigned ★★, as  $Z_2$  is at least as good as  $B_1$  on every criteria except d, and abc is established as sufficient by the comparison  $A_1$  vs  $C_1$ .

## ... to Accountability.

- Accountability is the ability of a human decision maker to own a recommendation made by the system and to *transfer* this ownership
- It suits MCDA better than mere trust, transparency, or persuasiveness, and leads to actual implementation
- Explanations require in-depth understanding of the preference models, and pose interesting computational challenges
- It mixes Decision Theory, Optimization techniques, and several Artificial Intelligence approaches (e.g. knowledge representation, argumentation)

## Connections to Computational Choice

- **Structure** : MCDA and CSC are structurally close, as Choice and Ranking mirrors Voting, and Ordinal Sorting mirrors Judgment Aggregation
- **Techniques** : Explaining the **result** of a Social Choice algorithm, or the **selection of a particular procedure**, could borrow techniques and insights
- **Applications** : Accountability is particularly needed in situations addressed simultaneously by MCDA and CSC, such as committee decisions
- **Complexity** : designing a model behaving well w.r.t. Accountability, incorporating requirements for accountability in adversarial contexts, modelling the collective reconstruction of explanations in a context similar to gossip,...