

JIAF 2022

Step-wise Explanations for the Additive Model

Manuel Amoussou¹, Khaled Belahcene²,
Nicolas Maudet³, Vincent Mousseau¹, Wassila Ouerdane¹

¹MICS, CentraleSupélec, Université Paris-Saclay

² Université de Technologie de Compiègne, CNRS, UMR 7253 Heudiasyc

³ Sorbonne Université, CNRS, Laboratoire d'Informatique de Paris 6

What ?

Explain **methodically** the comparative statement

$$X \succsim_{\omega} Y$$

where :

\succsim_{ω} is a binary relation on $2^{[m]}$

represented by the score function $\omega : \langle \omega_i \rangle_{i \in [m]}$ with $\omega_i \in \mathbb{N}$

$[m]$ is a finite set of indivisible items/criteria

$$X \succsim_{\omega} Y \Leftrightarrow \omega(X) \geq \omega(Y)$$

$$\omega(X) = \sum_{i \in X} \omega_i \text{ and } \omega(Y) = \sum_{j \in Y} \omega_j$$

How ?

Throughout the use of **argument schemes**:

operators tying a sequence of statements, called the premise, satisfying some conditions, into another statement called the conclusion

For ?

**Elicitation procedures
improvement and
Justification of decision models
outcomes**

in Decision Theory and Theory of measurement

1. Generalities

- A motivating example

- The Argument schemes & the Explanation problem

- The Argument schemes & the Additive model properties

2. Technical aspects

- The ceteris paribus scheme

- The transitive scheme

- The covering scheme

3. Numerical experimentation

Generalities

A motivating example

Criteria

- (a) affordability: “acceptable” (+) or “expensive” (–),
- (b) Metropole shipping fee included: “yes” (+) or “no” (–),
- (c) Overseas shipping fee included: “yes” (+) or “no” (–),
- (d) mask quality: “high” (+) or “good” (–),
- (e) provider reputation: “good” (+) or “fair” (–),
- (f) washable: “yes” (+) or “no” (–),
- (g) customizable: “yes” (+) or “no” (–).

Performance table & score function

	a	b	c	d	e	f	g	Scores
W	–	+	–	–	–	+	+	204
X	–	–	+	+	+	–	–	188
Y	+	–	–	+	–	–	–	187
Z	–	–	+	–	+	–	+	166
ω	128	126	77	59	52	41	37	-

$$W \equiv \mathbf{bfg} \mid X \equiv \mathbf{cde} \mid Y \equiv \mathbf{ad} \mid Z \equiv \mathbf{ceg}$$

$$\omega(Z) = \omega(\mathbf{ceg}) = 77 + 52 + 37 = 166$$

Explain why W is the best supplier

through the *step-wise* explanations of the comparative statements

$$W \succ_{\omega} X, W \succ_{\omega} Y \text{ and } W \succ_{\omega} Z$$

The Argument schemes & the Explanation problem

An argument scheme is an operator tying a sequence of statements, called the premise, satisfying some conditions, into another statement called the conclusion.

$\omega = \{a : 128, b : 126, c : 77, d : 59, e : 52, f : 41, g : 37\}$

$\text{bfg} \succsim_{\omega} \text{cde}$

$\text{ab} \succsim_{\omega} \text{ad}$

$\text{cde} \succsim_{\omega} \text{ab}$

$\text{bfg} \succsim_{\omega} \text{cde}$

$\text{a} \succsim_{\omega} \text{c}$

$\text{bc} \succsim_{\omega} \text{ad}$

$\text{bfg} \succsim_{\omega} \text{ad}$

$\text{b} \succsim_{\omega} \text{de}$

$\text{fg} \succsim_{\omega} \text{c}$

$\text{bfg} \succsim_{\omega} \text{ceg}$

$\text{bf} \succsim_{\omega} \text{ce}$

$\text{b} \succsim_{\omega} \text{d}$

The Explanation problem

Inputs:

- The comparative statement $(A, B) \in 2^{[m]} \times 2^{[m]}$ to explain,
- The score function ω and the induced preference relation \succsim_{ω} ,
- \mathcal{A} : a set of statements belonging to \succsim_{ω} ,
- A set of argument schemes \mathcal{S} ,
- a positive integer k .

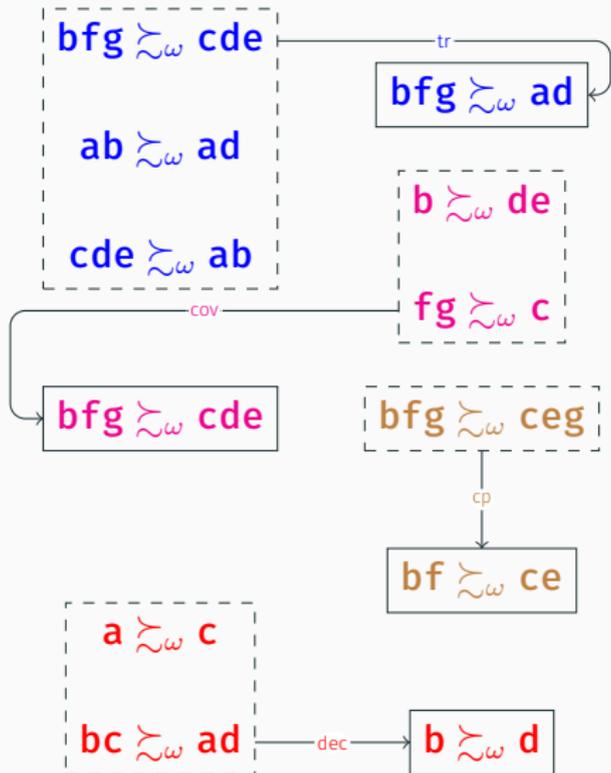
Question:

Is there a positive integer $k' \leq k$, a list of length k' of statements $[(A_1, B_1), \dots, (A_{k'}, B_{k'})]$ all belonging to \mathcal{A} , and a scheme $s \in \mathcal{S}$ such that $[(A_1, B_1), \dots, (A_{k'}, B_{k'})] \xrightarrow{s} (A, B)$?

The Argument schemes & the Explanation problem

An argument scheme is an operator tying a sequence of statements, called the premise, satisfying some conditions, into another statement called the conclusion.

$\omega = \{a : 128, b : 126, c : 77, d : 59, e : 52, f : 41, g : 37\}$



The Explanation problem

Inputs:

- The comparative statement $(A, B) \in 2^{[m]} \times 2^{[m]}$ to explain,
- The score function ω and the induced preference relation \succsim_{ω} ,
- \mathcal{A} : a set of statements belonging to \succsim_{ω} ,
- A set of argument schemes \mathcal{S} ,
- a positive integer k .

Question:

Is there a positive integer $k' \leq k$, a list of length k' of statements $[(A_1, B_1), \dots, (A_{k'}, B_{k'})]$ all belonging to \mathcal{A} , and a scheme $s \in \mathcal{S}$ such that $[(A_1, B_1), \dots, (A_{k'}, B_{k'})] \xrightarrow{s} (A, B)$?

The additive model properties

Let \succsim_ω be the binary relation induced by the score function ω .

- \succsim_ω is **transitive** : $X \succsim_\omega Y$ and $Y \succsim_\omega Z$ imply $X \succsim_\omega Z$
- \succsim_ω satisfies **1-order cancellation** : $(X \setminus Y) \succsim_\omega (Y \setminus X) \iff X \succsim_\omega Y$
- \succsim_ω satisfies **K-order cancellation** ($K \geq 2$):

$$\left. \begin{array}{l} (X^{(1)}, \dots, X^{(K)}) =_0 (Y^{(1)}, \dots, Y^{(K)}) \\ X^{(k)} \succsim_\omega Y^{(k)} \text{ for all } k < K \end{array} \right\} \Rightarrow Y^{(K)} \succsim_\omega X^{(K)}$$

$(X^{(1)}, \dots, X^{(K)}) =_0 (Y^{(1)}, \dots, Y^{(K)})$ mean that for every $i \in [m]$, $|\{k : i \in X^{(k)}\}| = |\{k : i \in Y^{(k)}\}|$ i.e. each criterion i appears as many times in the subsets $X^{(k)}$ than in the subsets $Y^{(k)}$.

Technical aspects

Generalities

- A motivating example

- The Argument schemes & the Explanation problem

- The Argument schemes & the Additive model properties

Technical aspects

- The ceteris paribus scheme

- The transitive scheme

- The covering scheme

- Numerical experimentation

The argument schemes : *Ceteris paribus* scheme

- ★ derived from the **1-order cancellation** property.
- ★ **One premise** : $(X \setminus Y) \succ_{\omega} (Y \setminus X)$
i.e. the comparison between the non common criteria of the comparative statement (X, Y) supported by ω .
- ★ **The conclusion** : $X \succ_{\omega} Y$.

e.g:

$$\omega = \{\mathbf{a} : 128, \mathbf{b} : 126, \mathbf{c} : 77, \mathbf{d} : 59, \mathbf{e} : 52, \mathbf{f} : 41, \mathbf{g} : 37\}$$

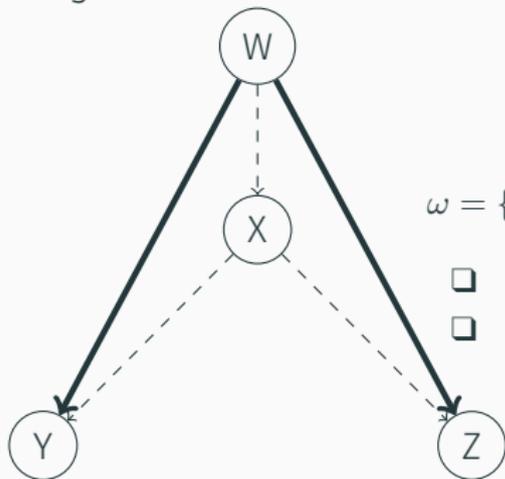
$$W \equiv \mathbf{bfg} \mid X \equiv \mathbf{cde} \mid Y \equiv \mathbf{ad} \mid Z \equiv \mathbf{ceg}$$

$$\begin{array}{lll} \square & \mathbf{d} \succ_{\omega} \mathbf{g} & \Rightarrow X \succ_{\omega} Z \\ \square & \mathbf{bf} \succ_{\omega} \mathbf{ce} & \Rightarrow W \succ_{\omega} Z \\ \square & \mathbf{bfg} \succ_{\omega} \mathbf{cde} & \Rightarrow W \succ_{\omega} X \end{array}$$

The argument schemes : *Transitive scheme*

- ★ derived from the **transitive** property.
- ★ **Premises** : $X^{(1)} \succsim_{\omega} Y^{(1)}, \dots, X^{(k)} \succsim_{\omega} Y^{(k)}$
with $X^{(1)} = X, X^{(j)} = Y^{(j-1)}$ for all $j \geq 2$ and $Y^{(k)} = Y$.
- ★ **The conclusion** : $X \succsim_{\omega} Y$.

e.g:



$$\omega = \{\mathbf{a} : 128, \mathbf{b} : 126, \mathbf{c} : 77, \mathbf{d} : 59, \mathbf{e} : 52, \mathbf{f} : 41, \mathbf{g} : 37\}$$

$$\begin{array}{l} \square \quad W \succsim_{\omega} X \quad \wedge \quad X \succsim_{\omega} Y \quad \Rightarrow \quad W \succsim_{\omega} Y \\ \square \quad W \succsim_{\omega} X \quad \wedge \quad X \succsim_{\omega} Z \quad \Rightarrow \quad W \succsim_{\omega} Z \end{array}$$

The argument schemes : *Covering scheme*

- ★ derived from (a particular case of) the ***K-order cancellation*** property.

$$\left. \begin{array}{l} (X^{(1)}, \dots, X^{(K-1)}, Y^{(K)}) =_0 (Y^{(1)}, \dots, Y^{(K-1)}, X^{(K)}) \\ X^{(k)} \succ_{\omega} Y^{(k)} \text{ for all } k < K \end{array} \right\} \Rightarrow X^{(K)} \succ_{\omega} Y^{(K)}$$

- ★ **Premises** : $X^{(k)} \succ_{\omega} Y^{(k)}$ for all $k < K$ with
 - the subsets of criteria $X^{(1)}, \dots, X^{(K-1)}$ partitioning $X^{(K)}$.
 - the subsets of criteria $Y^{(1)}, \dots, Y^{(K-1)}$ partitioning $Y^{(K)}$.
 - $X \setminus Y = X^{(K)}$ and $Y \setminus X = Y^{(K)}$

- ★ **The conclusion** : $X \succ_{\omega} Y$.

e.g: $\omega = \{\mathbf{a} : 128, \mathbf{b} : 126, \mathbf{c} : 77, \mathbf{d} : 59, \mathbf{e} : 52, \mathbf{f} : 41, \mathbf{g} : 37\}; W \equiv \mathbf{bfg} \mid X \equiv \mathbf{cde}$

$$(\mathbf{b}, \mathbf{fg}, \overbrace{\mathbf{cde}}^X) =_0 (\mathbf{de}, \mathbf{c}, \overbrace{\mathbf{bfg}}^W)$$

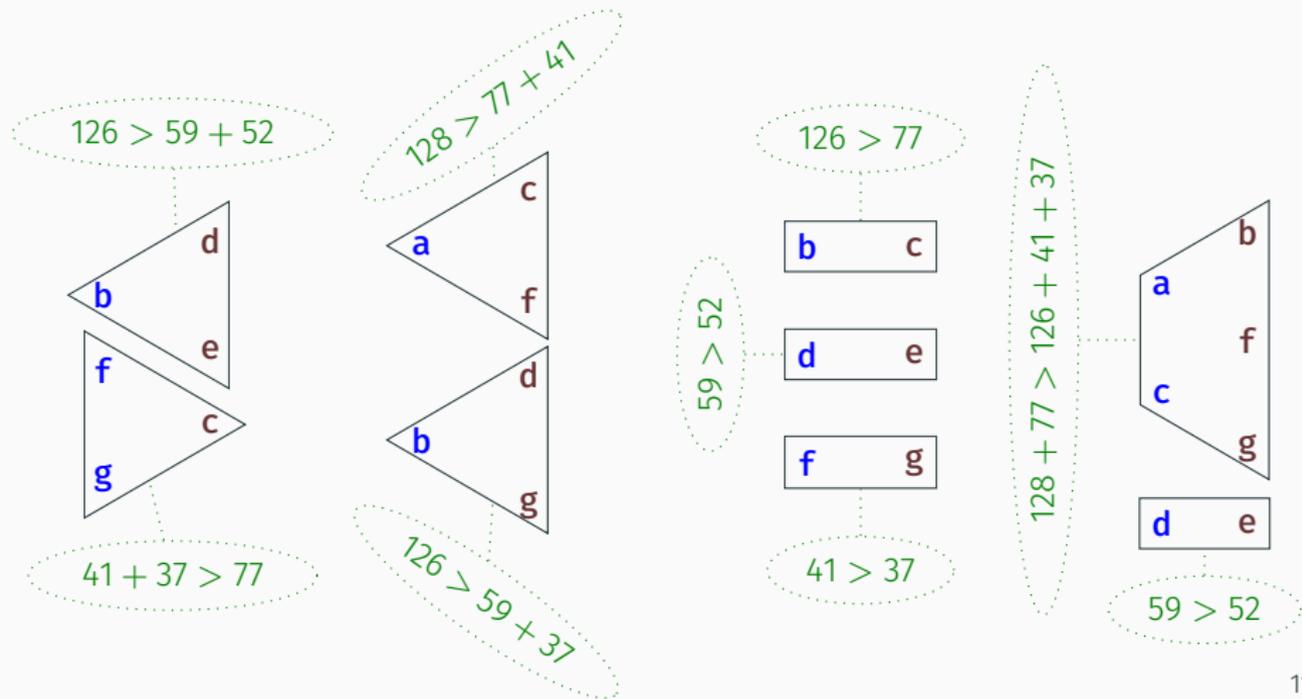
$$\mathbf{b} \succ_{\omega} \mathbf{de}$$

$$\mathbf{fg} \succ_{\omega} \mathbf{c}$$

The argument schemes : *Covering scheme*

$\omega = \{a : 128, b : 126, c : 77, d : 59, e : 52, f : 41, g : 37\}$

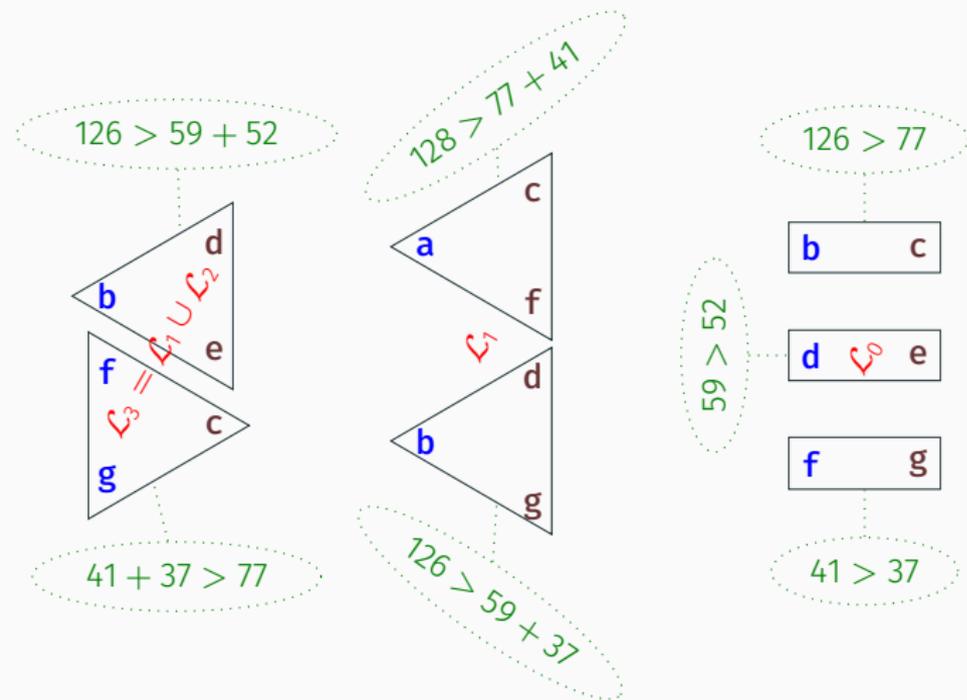
pro criteria vs. con criteria



The argument schemes : *Covering scheme*

$\omega = \{a : 128, b : 126, c : 77, d : 59, e : 52, f : 41, g : 37\}$

pro criteria vs. con criteria



Decomposition Languages

- $\mathcal{L}_0 = \Delta(1, 1)$
- $\mathcal{L}_1 = \Delta(1, m)$
- $\mathcal{L}_2 = \Delta(m, 1)$
- $\mathcal{L}_3 = \mathcal{L}_1 \cup \mathcal{L}_2$

The chosen decomposition languages



- + Cognitively easy to grasp.
- + Easily scriptable in a natural language.
- Not complete.

$$\omega = \{\mathbf{a} : 128, \mathbf{b} : 126, \mathbf{c} : 77, \mathbf{d} : 59, \mathbf{e} : 52, \mathbf{f} : 41, \mathbf{g} : 37\}$$

$$\mathbf{cf} \stackrel{\omega}{\sim} \mathbf{de}$$

is not decomposable !!!

Covering scheme: Computational complexity

Instance:

- A set **PRO** of pro criteria and a set **CON** of con criteria
- The preference model representation ω

\mathcal{L}_0

Question: Is there a bijective function f from **CON** to **PRO** such that for all $j \in \text{CON}$, $\omega(f(j)) \geq \omega(j)$?

Complexity class: P

\mathcal{L}_1

Question: Is there a function f from **PRO** to 2^{CON} such that :

1. $f(i) \cap f(i') = \emptyset$ if $i \neq i'$ with $i, i' \in \text{PRO}$
2. $\bigcup_{i \in \text{PRO}} f(i) = \text{CON}$
3. $\omega(i) \geq \sum_{j \in f(i)} \omega(j)$ for all $i \in \text{PRO}$

Complexity class: NP-hard

\mathcal{L}_2

Question: Is there a application g from **CON** to $2^{\text{PRO} \setminus \{\emptyset\}}$ such that :

1. $g(j) \cap g(j') = \emptyset$ if $j \neq j'$ with $j, j' \in \text{CON}$
2. $\sum_{i \in g(j)} \omega(i) \geq \omega(j)$ for all $j \in \text{CON}$

Complexity class: NP-hard

\mathcal{L}_3

Question: Are there two disjoint subsets **PRO**¹ and **PRO**² of **PRO**, two disjoint subsets **CON**¹ and **CON**² of **CON**, a function f from **PRO**¹ to 2^{CON^1} and a application g from **CON**² to $2^{\text{PRO}^2 \setminus \{\emptyset\}}$ such that:

1. **PRO**¹ \cup **PRO**² = **PRO** and **CON**¹ \cup **CON**² = **CON**
2. $f(i) \cap f(i') = \emptyset$ if $i \neq i'$ with $i, i' \in \text{PRO}^1$
3. $\bigcup_{i \in \text{PRO}^1} f(i) = \text{CON}^1$
4. $\omega(i) \geq \sum_{j \in f(i)} \omega(j)$ for all $i \in \text{PRO}^1$
5. $g(j) \cap g(j') = \emptyset$ if $j \neq j'$ with $j, j' \in \text{CON}^2$
6. $\sum_{i \in g(j)} \omega(i) \geq \omega(j)$ for all $j \in \text{CON}^2$

Complexity class: NP-hard

\mathcal{L}_3 —Covering scheme: Computation using ILP

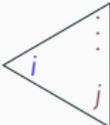
Inputs

- The score function $\omega : \langle \omega_i \rangle_{i \in [m]}$.
- The pair (X, Y) to explain.

Notations

- *Pro arguments set* : $(X, Y)^+$
- *Con arguments set* : $(X, Y)^-$

(Binary) Variables

$$s_{ij}^1 = \begin{cases} 1 & \text{if } j \in f(i). \\ 0 & \text{Otherwise.} \end{cases}$$


$$s_{ij}^2 = \begin{cases} 1 & \text{if } i \in g(j). \\ 0 & \text{Otherwise.} \end{cases}$$


Constraints

- *Syntactic constraints*
 - For each *pro argument* i

$$s_{ij}^1 + \sum_{j' \in (X, Y)^-} s_{ij'}^2 \leq 1 \quad \forall j \in (X, Y)^-$$

- For each *con argument* j

$$\sum_{i' \in (X, Y)^+} s_{i'j}^1 + s_{ij}^2 \leq 1 \quad \forall i \in (X, Y)^+$$

- ω -compatibility constraints
 - For each *pro argument* i

$$\omega_i \geq \sum_{j \in (X, Y)^-} \omega_j s_{ij}^1$$

- For each *con argument* j

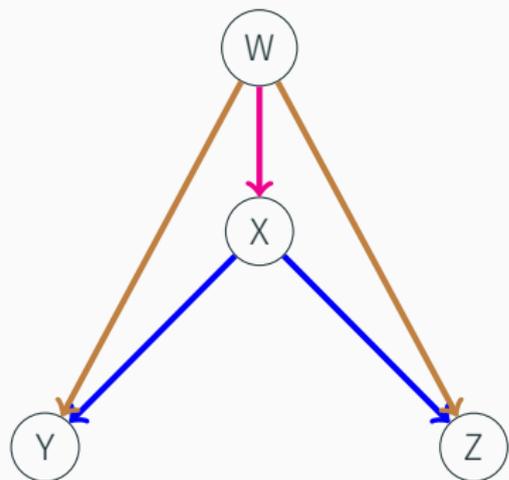
$$\sum_{i \in (X, Y)^+} (s_{ij}^1 + s_{ij}^2) \omega_i \geq \omega_j$$

Back to the motivating example

Combining Ceteris Paribus, Transitive and Covering schemes

$$\omega = \{a : 128, b : 126, c : 77, d : 59, e : 52, f : 41, g : 37\}$$

$$W \equiv bfg \mid X \equiv cde \mid Y \equiv ad \mid Z \equiv ceg$$



\mathcal{L}_3 -Covering Scheme

$$\star \quad b \underset{\omega}{\sim} de \quad \wedge \quad fg \underset{\omega}{\sim} c \quad \Rightarrow \quad W \underset{\omega}{\sim} X$$

Transitive Scheme

$$\star \quad W \underset{\omega}{\sim} X \quad \wedge \quad X \underset{\omega}{\sim} Y \quad \Rightarrow \quad W \underset{\omega}{\sim} Y$$

$$\star \quad W \underset{\omega}{\sim} X \quad \wedge \quad X \underset{\omega}{\sim} Z \quad \Rightarrow \quad W \underset{\omega}{\sim} Z$$

Ceteris paribus Scheme

$$\star \quad ce \underset{\omega}{\sim} a \quad \Rightarrow \quad X \underset{\omega}{\sim} Y$$

$$\star \quad d \underset{\omega}{\sim} g \quad \Rightarrow \quad X \underset{\omega}{\sim} Z$$

Numerical experimentation

Numerical Experiments

Instance: (ω, \mathbb{A})

- $|\mathbb{A}| = 10$
- $|\omega| = m \in [6; 15]$

Sample: 500.000 instances

- 500 sets \mathbb{A}
- 1000 score functions ω

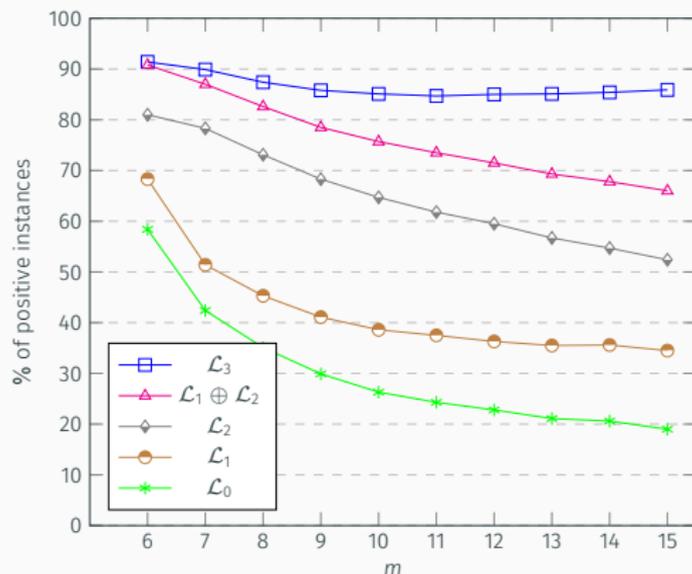
\mathbb{A} set characteristics:

- no Pareto dominance
- significance of all criteria

Argument Scheme used:

- Covering scheme

Percentage of positive instances



An instance is positive if and only if all “direct” statements (x^*, y) are \mathcal{L} -covering explainable.

m	6	7	8	9	10	11	12	13	14	15
\mathcal{L}_3	91.4%	89.9%	87.4%	85.8%	85.1%	84.7%	85.0%	85.1%	85.4%	85.9%
$\mathcal{L}_1 \oplus \mathcal{L}_2$	90.8%	87.0%	82.6%	78.5%	75.7%	73.5%	71.5%	69.3%	67.8%	66.0%
\mathcal{L}_1	81.0%	78.3%	73.1%	68.3%	64.7%	61.8%	59.1%	56.7%	54.7%	52.4%
\mathcal{L}_2	68.3%	51.4%	45.3%	41.1%	38.6%	37.5%	36.3%	35.5%	35.6%	34.5%
\mathcal{L}_0	58.4%	42.4%	35.1%	29.9%	26.3%	24.3%	22.8%	21.1%	20.6%	19.0%

Merci!

-  Amgoud, L. and Prade, H. (2009).
Using arguments for making and explaining decisions.
Artificial Intelligence, 173(3):413–436.
-  Belahcene, K., Chevaleyre, Y., Labreuche, C., Maudet, N., Mousseau, V., and Ouerdane, W. (2018).
Accountable approval sorting.
In *Proceedings IJCAI'18*, pages 70–76.
-  Belahcene, K., Labreuche, C., Maudet, N., Mousseau, V., and Ouerdane, W. (2017).
Explaining robust additive utility models by sequences of preference swaps.
Theory and Decision, 82(2):151–183.



Belahcene, K., Labreuche, C., Maudet, N., Mousseau, V., and Ouerdane, W. (2019).

Comparing options with argument schemes powered by cancellation.

In *Proc. of IJCAI'19*, pages 1537–1543.



Belton, V. and Stewart, T. J. (2002).

Multiple criteria decision analysis - an integrated approach.

Springer.



Bouyssou, D., Marchant, T., Pirlot, M., Tsoukias, A., and Vincke, P. (2006).

***Evaluation and Decision models with multiple criteria: stepping stones for the analyst*, volume 86.**

Springer.

-  Greco, S., Mousseau, V., and Słowiński, R. (2008).
Ordinal regression revisited: Multiple criteria ranking using a set of additive value functions.
EJOR, 191(2):416–436.
-  Hammond, J., Keeney, R., and Raiffa, H. (1998).
Even swaps: A rational method for making trade-offs.
Harvard business review, 76:137–8, 143.
-  Labreuche, C. and Fossier, S. (2018).
Explaining multi-criteria decision aiding models with an extended shapley value.
In *Proc. of IJCAI'18*, pages 331–339.



Labreuche, C., Maudet, N., and Ouerdane, W. (2011).

Minimal and complete explanations for critical multi-attribute decisions.

In *Proc. of ADT*, pages 121–134, Piscataway, NJ, USA.



Labreuche, C., Maudet, N., and Ouerdane, W. (2012).

Justifying dominating options when preferences are incomplete.

In *Proceedings of ECAI*, pages 486–491, Montpellier, France.