

Decision Making under Weakly Structured Information

with Applications to Robust Statistics and Machine Learning

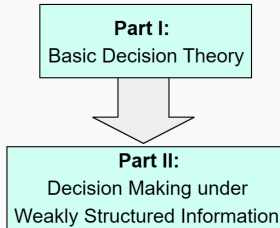
SIPTA Seminar, April 30, 2024

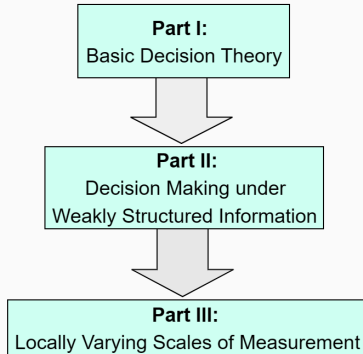
Christoph Jansen, Department of Statistics, LMU Munich

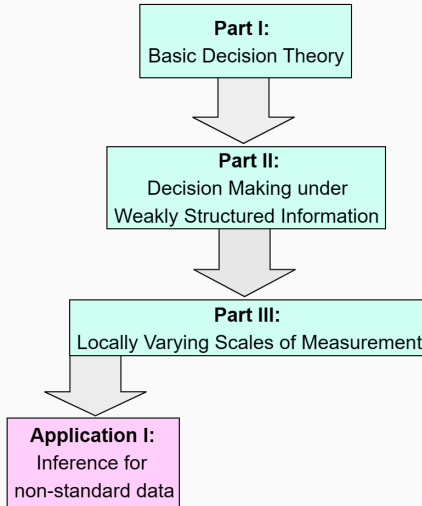
Introducing Myself

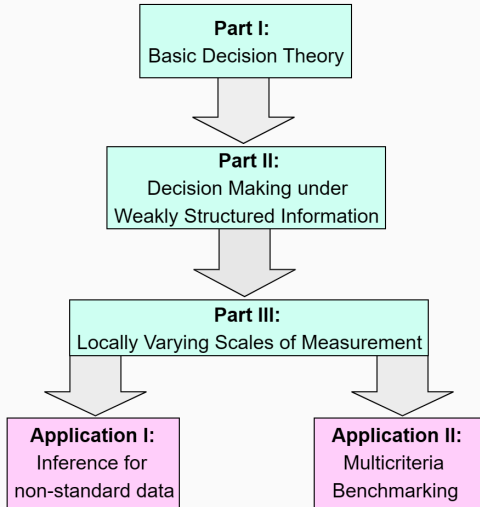
- Christoph Jansen
- Member of Thomas Augustin's *Foundations of Statistics* group
- Affiliated with the Department of Statistics of LMU Munich
- About to finish my habilitation in statistics
- Assistant professor at Lancaster University Leipzig starting June 2024
- First contact with SIPTA at the 2014 Summerschool in Montpellier

Part I:
Basic Decision Theory









Basic Decision Theory

Classical Decision Theory

Informal description of the model:

- An **agent** has to choose among different **acts** X from a set \mathcal{G} .
- The **consequence** that choosing X yields depends on which **state of nature** s from a set S is the true one.

Classical Decision Theory

Informal description of the model:

- An agent has to choose among different acts X from a set \mathcal{G} .
- The consequence that choosing X yields depends on which state of nature s from a set S is the true one.

Formal description of the model:

- Let A denote some non-empty set of consequences.
- Each act X corresponds to a mapping $X : S \rightarrow A$.
- The set \mathcal{G} is a subset of $A^S = \{X : S \rightarrow A\}$.

Classical Decision Theory

Informal description of the model:

- An agent has to choose among different acts X from a set \mathcal{G} .
- The consequence that choosing X yields depends on which state of nature s from a set S is the true one.

Formal description of the model:

- Let A denote some non-empty set of consequences.
- Each act X corresponds to a mapping $X : S \rightarrow A$.
- The set \mathcal{G} is a subset of $A^S = \{X : S \rightarrow A\}$.

Goal: Determining a **choice function**

$$ch : 2^{\mathcal{G}} \rightarrow 2^{\mathcal{G}} \text{ with } ch(\mathcal{D}) \subseteq \mathcal{D} \text{ for all } \mathcal{D} \in 2^{\mathcal{G}}$$

that best possibly utilizes the available information.

Interpreting Choice Functions

Depending on the quality of the underlying information, the choice sets $ch(\mathcal{D})$ can be given two different interpretations:

Interpreting Choice Functions

Depending on the quality of the underlying information, the choice sets $ch(\mathcal{D})$ can be given two different interpretations:

Strong interpretation:

$ch(\mathcal{D})$ is the set of equally optimal acts from \mathcal{D} .
The agent is indifferent between these acts.

Interpreting Choice Functions

Depending on the quality of the underlying information, the choice sets $ch(\mathcal{D})$ can be given two different interpretations:

Strong interpretation:

$ch(\mathcal{D})$ is the set of equally optimal acts from \mathcal{D} .
The agent is indifferent between these acts.

Weak interpretation:

$ch(\mathcal{D})$ is the set of all non-neglectable acts from \mathcal{D} given the information.
These acts are incomparable for the agent.

Interpreting Choice Functions

Depending on the quality of the underlying information, the **choice sets** $ch(\mathcal{D})$ can be given two different **interpretations**:

Strong interpretation:

$ch(\mathcal{D})$ is the set of equally optimal acts from \mathcal{D} .
The agent is **indifferent** between these acts.

Weak interpretation:

$ch(\mathcal{D})$ is the set of all non-neglectable acts from \mathcal{D} given the information.
These acts are **incomparable** for the agent.

Obvious comment:

If only **weakly structured information** is available, we often have to work with weakly interpretable choice functions.

Two Classical Choice Functions under Risk

Expected utility:

If a probability π on S and a cardinal scale $u : A \rightarrow [0, 1]$ are available, set

$$ch_{u,\pi}(\mathcal{D}) = \left\{ Y \in \mathcal{D} : \mathbb{E}_\pi(u \circ Y) \geq \mathbb{E}_\pi(u \circ X) \text{ for all } X \in \mathcal{D} \right\},$$

and choose that acts from \mathcal{D} that **maximize expected utility**.

Two Classical Choice Functions under Risk

Expected utility:

If a probability π on S and a cardinal scale $u : A \rightarrow [0, 1]$ are available, set

$$ch_{u,\pi}(\mathcal{D}) = \left\{ Y \in \mathcal{D} : \mathbb{E}_\pi(u \circ Y) \geq \mathbb{E}_\pi(u \circ X) \text{ for all } X \in \mathcal{D} \right\},$$

and choose that acts from \mathcal{D} that **maximize expected utility**.

First-Order Stochastic Dominance:

If a probability π on S and a preorder \succsim on A are available, set

$$ch_{\succsim,\pi}(\mathcal{D}) = \left\{ Y : \nexists X \text{ s.t. } \begin{array}{l} \mathbb{E}_\pi(u \circ X) \geq \mathbb{E}_\pi(u \circ Y) \text{ for all } u \in \mathcal{U}_{\succsim} \\ \mathbb{E}_\pi(u \circ X) > \mathbb{E}_\pi(u \circ Y) \text{ for some } u \in \mathcal{U}_{\succsim} \end{array} \right\}$$

where \mathcal{U}_{\succsim} is the set of all \succsim -isotone $u : A \rightarrow [0, 1]$. Choose acts that are **not excluded by every compatible EU-maximizer**.

A Toy Example

An agent wants to invest in **exactly one** of the stocks in $\mathcal{G} = \{X_1, X_2, X_3\}$.

The consequence depends on the true economic scenario from $S = \{s_1, s_2, s_3\}$.

Suppose we have the following consequence table, where $A = \{a_1, \dots, a_9\}$:

	s_1	s_2	s_3
X_1	a_1	a_4	a_7
X_2	a_2	a_5	a_8
X_3	a_3	a_6	a_9

Moreover, assume π is the uniform distribution on S .

A Toy Example, continued: Cardinal preferences

Assume, the agent's preferences allow for a description via a **cardinal** utility $u : A \rightarrow \mathbb{R}$ on A (i.e., u is unique up to plts).

A Toy Example, continued: Cardinal preferences

Assume, the agent's preferences allow for a description via a **cardinal** utility $u : A \rightarrow \mathbb{R}$ on A (i.e., u is unique up to plts).

⇒ Consequence table can be transformed in utility table, e.g.:

	s_1	s_2	s_3
$u \circ X_1$	6000	3000	-2000
$u \circ X_2$	8000	1000	-3000
$u \circ X_3$	5000	4000	0

Due to uniqueness of u , maximizing expected utility is well-defined.

A Toy Example, continued: Cardinal preferences

Assume, the agent's preferences allow for a description via a **cardinal** utility $u : A \rightarrow \mathbb{R}$ on A (i.e., u is unique up to plts).

⇒ Consequence table can be transformed in utility table, e.g.:

	s_1	s_2	s_3
$u \circ X_1$	6000	3000	-2000
$u \circ X_2$	8000	1000	-3000
$u \circ X_3$	5000	4000	0

Due to uniqueness of u , maximizing expected utility is well-defined.

⇒ We can apply

$$\text{ch}_{u,\pi}(\mathcal{G}) = \operatorname{argmax}_{X \in \mathcal{G}} \mathbb{E}_\pi(u \circ X) = \{X_3\}$$

A Toy Example, continued: Cardinal preferences

Assume, the agent's preferences allow for a description via a **cardinal** utility $u : A \rightarrow \mathbb{R}$ on A (i.e., u is unique up to plts).

⇒ Consequence table can be transformed in utility table, e.g.:

	s_1	s_2	s_3
$u \circ X_1$	6000	3000	-2000
$u \circ X_2$	8000	1000	-3000
$u \circ X_3$	5000	4000	0

Due to uniqueness of u , maximizing expected utility is well-defined.

⇒ We can apply

$$\text{ch}_{u,\pi}(\mathcal{G}) = \operatorname{argmax}_{X \in \mathcal{G}} \mathbb{E}_\pi(u \circ X) = \{X_3\}$$

⇒ X_3 is the unique optimal stock. (**Strong interpretation!**)

A Toy Example, continued: (P)ordinal preferences

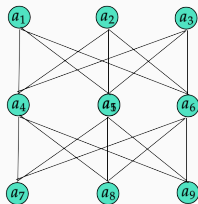
Now assume, the preferences allow only for a preorder \succsim on A .

A Toy Example, continued: (P)ordinal preferences

Now assume, the preferences allow only for a **preorder** \succsim on A .

Then, the situation looks for instance like this:

	S_1	S_2	S_3
X_1	a_1	a_4	a_7
X_2	a_2	a_5	a_8
X_3	a_3	a_6	a_9



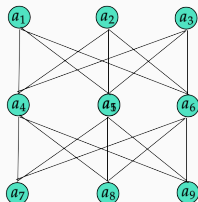
Every \succsim -isotone function $u : A \rightarrow \mathbb{R}$ is a compatible scale.

A Toy Example, continued: (P)ordinal preferences

Now assume, the preferences allow only for a **preorder** \succsim on A .

Then, the situation looks for instance like this:

	S_1	S_2	S_3
X_1	a_1	a_4	a_7
X_2	a_2	a_5	a_8
X_3	a_3	a_6	a_9



Every \succsim -isotone function $u : A \rightarrow \mathbb{R}$ is a compatible scale.

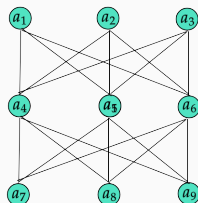
\Rightarrow Maximizing expected utility is not-well-defined!

A Toy Example, continued: (P)ordinal preferences

Now assume, the preferences allow only for a preorder \succsim on A .

Then, the situation looks for instance like this:

	S_1	S_2	S_3
X_1	a_1	a_4	a_7
X_2	a_2	a_5	a_8
X_3	a_3	a_6	a_9



Every \succsim -isotone function $u : A \rightarrow \mathbb{R}$ is a compatible scale.

\Rightarrow Maximizing expected utility is not-well-defined!

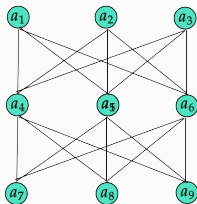
\Rightarrow We can still apply FSD, but now $ch_{\succsim, \pi}(\mathcal{G}) = \mathcal{G}$, since for every stock there exists a compatible scale making it the unique EU-maximizer.

A Toy Example, continued: (P)ordinal preferences

Now assume, the preferences allow only for a **preorder** \succsim on A .

Then, the situation looks for instance like this:

	S_1	S_2	S_3
X_1	a_1	a_4	a_7
X_2	a_2	a_5	a_8
X_3	a_3	a_6	a_9



Every \succsim -isotone function $u : A \rightarrow \mathbb{R}$ is a compatible scale.

\Rightarrow Maximizing expected utility is not-well-defined!

\Rightarrow We can still apply FSD, but now $\text{ch}_{\succsim, \pi}(\mathcal{G}) = \mathcal{G}$, since for every stock there exists a compatible scale making it the unique EU-maximizer.

\Rightarrow None of the stocks can be excluded. (**Weak interpretation!**)

Weakly structured Information

Common Assumptions in Classic Decision Theory

Classical assumptions:

- (I) The agent's preferences among the elements of A are characterized by a **cardinal utility function** $u : A \rightarrow \mathbb{R}$.
- (II) The uncertainty among the states from S is described by some **classical probability measure** π .

Recall:

Expected utility rule $ch_{u,\pi}(\cdot)$ relies on both (I) and (II).

Stochastic dominance rule $ch_{\succsim,\pi}(\cdot)$ relies on (II) but not on (I).

Challenging the Classical Assumptions

Problem: Both (I) and (II) rely on **strong axiomatic assumptions**.

(e.g., [von Neumann et al., 1944, Savage, 1954]))

Together, these assumptions explicitly dismiss:

- Purely **ordinal** or **partial** preferences.
(e.g., [Seidenfeld et al., 1995, Nau, 2006]))
- Agents with **partial probabilistic** beliefs.
(e.g., [Levi, 1974, Walley, 1991, Kikuti et al., 2011]))
- Problems of **group decision making**.
(e.g., [Bacharach, 1975, Bradley, 2019]))

These are **highly relevant situations** to investigate!

Weakly structured Information

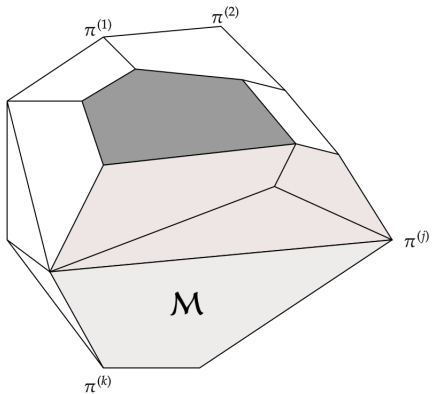
Two different sources of complexity:

Weakly structured Information

Two different sources of complexity:

- (I)' **Imprecise probabilistic models:** If it isn't possible to specify **one** probability on S , we still can work with the *set* \mathcal{M} of all probabilities **compatible with the information**.

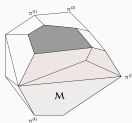
Weakly structured Information



Weakly structured Information

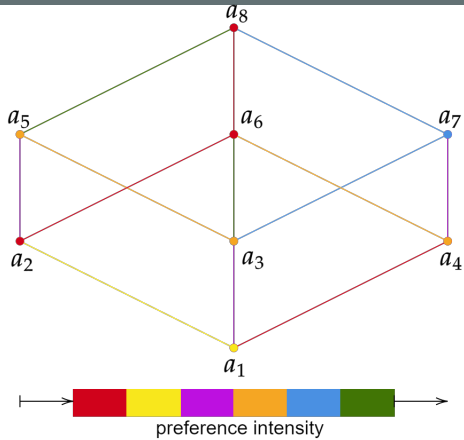
Two different sources of complexity:

- (I)' **Imprecise probabilistic models:** If it isn't possible to specify **one** probability on S , we still can work with the *set* \mathcal{M} of all probabilities **compatible with the information**.



- (II)' **Complexly ordered consequences:** A cardinal utility demands the agent to satisfy **very restrictive axioms**. If these are too restrictive, we still can work with the *set* \mathcal{U} of all utilities **compatible with the information**.

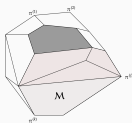
Weakly structured Information



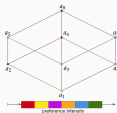
Weakly structured Information

Two different sources of complexity:

- (I)' **Imprecise probabilistic models:** If it isn't possible to specify **one** probability on S , we still can work with the *set* \mathcal{M} of all probabilities **compatible with the information**.



- (II)' **Complexly ordered consequences:** A cardinal utility demands the agent to satisfy **very restrictive axioms**. If these are too restrictive, we still can work with the *set* \mathcal{U} of all utilities **compatible with the information**.



Modelling \mathcal{U} : Preference Systems

Notation: Binary relation R has **strict part** P_R and **indifference part** I_R .

Preference system & Consistency

Let A denote a set of consequences. Let further

- $R_1 \subseteq A \times A$ be a binary relation on A
- $R_2 \subseteq R_1 \times R_1$ be a binary relation on R_1

The triplet $\mathcal{A} = [A, R_1, R_2]$ is called a **preference system** on A .

We call \mathcal{A} **consistent** if there is $u : A \rightarrow [0, 1]$ with for all $a, b, c, d \in A$:

$$(a, b) \in R_1 \Rightarrow u(a) \geq u(b) \quad (\text{with } = \text{ iff } \in I_{R_1}).$$

$$((a, b), (c, d)) \in R_2 \Rightarrow u(a) - u(b) \geq u(c) - u(d) \quad (\text{with } = \text{ iff } \in I_{R_2}).$$

The set of all representations u of \mathcal{A} is denoted by $\mathcal{U}_{\mathcal{A}}$.

Interpretation of the components of \mathcal{A} :

- $(a, b) \in R_1$: “*a is at least as desirable as b*”
- $((a, b), (c, d)) \in R_2$: “*exchanging b by a is at least as desirable as d by c*”

Modelling \mathcal{U} : Preference Systems

Notation: Binary relation R has **strict part** P_R and **indifference part** I_R .

Preference system & Consistency

Let A denote a set of consequences. Let further

$R_1 \subseteq A \times A$ be a binary relation on A

$R_2 \subseteq R_1 \times R_1$ be a binary relation on R_1

The triplet $\mathcal{A} = [A, R_1, R_2]$ is called a preference system on A .

We call \mathcal{A} **consistent** if there is $u : A \rightarrow [0, 1]$ with for all $a, b, c, d \in A$:

- $(a, b) \in R_1 \Rightarrow u(a) \geq u(b)$ (with $=$ iff $\in I_{R_1}$).
- $((a, b), (c, d)) \in R_2 \Rightarrow u(a) - u(b) \geq u(c) - u(d)$ (with $=$ iff $\in I_{R_2}$).

The set of all representations u of \mathcal{A} is denoted by $\mathcal{U}_{\mathcal{A}}$.

Preference System: Toy Example

Suppose, an agent is looking for a new job.

Preference System: Toy Example

Suppose, an agent is looking for a new job.

Jobs are compared w.r.t. **salary after tax** and **additional benefits**.

Preference System: Toy Example

Suppose, an agent is looking for a new job.

Jobs are compared w.r.t. salary after tax and additional benefits.

The set of potential benefits $\mathcal{X} = \{b_1, \dots, b_5\}$ is given by

b_1	b_2	b_3	b_4	b_5
overtime premium	child care	advanced training	promotion prospects	flexible hours

Preference System: Toy Example

Suppose, an agent is looking for a new job.

Jobs are compared w.r.t. salary after tax and additional benefits.

The set of potential benefits $\mathcal{X} = \{b_1, \dots, b_5\}$ is given by

b_1	b_2	b_3	b_4	b_5
overtime premium	child care	advanced training	promotion prospects	flexible hours

Then $A \subseteq \mathbb{R}^+ \times 2^{\mathcal{X}}$ and $\mathcal{A} = [A, R_1, R_2]$ is given by

Preference System: Toy Example

Suppose, an agent is looking for a new job.

Jobs are compared w.r.t. salary after tax and additional benefits.

The set of potential benefits $\mathcal{X} = \{b_1, \dots, b_5\}$ is given by

b_1	b_2	b_3	b_4	b_5
overtime premium	child care	advanced training	promotion prospects	flexible hours

Then $A \subseteq \mathbb{R}^+ \times 2^{\mathcal{X}}$ and $\mathcal{A} = [A, R_1, R_2]$ is given by

$$R_1 = \left\{ ((y_1, B_1), (y_2, B_2)) : y_1 \geq y_2 \wedge B_1 \supseteq B_2 \right\}$$

Preference System: Toy Example

Suppose, an agent is looking for a new job.

Jobs are compared w.r.t. salary after tax and additional benefits.

The set of potential benefits $\mathcal{X} = \{b_1, \dots, b_5\}$ is given by

b_1	b_2	b_3	b_4	b_5
overtime premium	child care	advanced training	promotion prospects	flexible hours

Then $A \subseteq \mathbb{R}^+ \times 2^{\mathcal{X}}$ and $\mathcal{A} = [A, R_1, R_2]$ is given by

$$R_1 = \left\{ \left((y_1, B_1), (y_2, B_2) \right) : y_1 \geq y_2 \wedge B_1 \supseteq B_2 \right\}$$

$$R_2 = \left\{ \left(\left((y_1, B_1), (y_2, B_2) \right), \left((y_3, B_3), (y_4, B_4) \right) \right) : \begin{array}{l} y_1 - y_2 \geq y_3 - y_4 \wedge \\ B_1 \supseteq B_3 \supseteq B_4 \supseteq B_2 \end{array} \right\}$$

Preference System: Toy Example

Suppose, an agent is looking for a new job.

Jobs are compared w.r.t. salary after tax and additional benefits.

The set of potential benefits $\mathcal{X} = \{b_1, \dots, b_5\}$ is given by

b_1	b_2	b_3	b_4	b_5
overtime premium	child care	advanced training	promotion prospects	flexible hours

Then $A \subseteq \mathbb{R}^+ \times 2^{\mathcal{X}}$ and $\mathcal{A} = [A, R_1, R_2]$ is given by

$$R_1 = \left\{ \left((y_1, B_1), (y_2, B_2) \right) : y_1 \geq y_2 \wedge B_1 \supseteq B_2 \right\}$$

$$R_2 = \left\{ \left(\left((y_1, B_1), (y_2, B_2) \right), \left((y_3, B_3), (y_4, B_4) \right) \right) : \begin{array}{l} y_1 - y_2 \geq y_3 - y_4 \wedge \\ B_1 \supseteq B_3 \supseteq B_4 \supseteq B_2 \end{array} \right\}$$

" $(y_1, B_1) \leftarrow (y_2, B_2)$ is more desirable than $(y_3, B_3) \leftarrow (y_4, B_4)$ if it yields a higher increase of salary and a superset increase of benefits."

Preference System: Toy Example

Suppose, an agent is looking for a new job.

Jobs are compared w.r.t. salary after tax and additional benefits.

The set of potential benefits $\mathcal{X} = \{b_1, \dots, b_5\}$ is given by

b_1	b_2	b_3	b_4	b_5
overtime premium	child care	advanced training	promotion prospects	flexible hours

Then $A \subseteq \mathbb{R}^+ \times 2^{\mathcal{X}}$ and $\mathcal{A} = [A, R_1, R_2]$ is given by

$$R_1 = \left\{ \left((y_1, B_1), (y_2, B_2) \right) : y_1 \geq y_2 \wedge B_1 \supseteq B_2 \right\}$$

$$R_2 = \left\{ \left(\left((y_1, B_1), (y_2, B_2) \right), \left((y_3, B_3), (y_4, B_4) \right) \right) : \begin{array}{l} y_1 - y_2 \geq y_3 - y_4 \wedge \\ B_1 \supseteq B_3 \supseteq B_4 \supseteq B_2 \end{array} \right\}$$

*" $(y_1, B_1) \leftarrow (y_2, B_2)$ is more desirable than $(y_3, B_3) \leftarrow (y_4, B_4)$ if it yields a **higher increase of salary** and a **superset increase of benefits**."*

Preference System: Toy Example

Suppose, an agent is looking for a new job.

Jobs are compared w.r.t. salary after tax and additional benefits.

The set of potential benefits $\mathcal{X} = \{b_1, \dots, b_5\}$ is given by

b_1	b_2	b_3	b_4	b_5
overtime premium	child care	advanced training	promotion prospects	flexible hours

Then $A \subseteq \mathbb{R}^+ \times 2^{\mathcal{X}}$ and $\mathcal{A} = [A, R_1, R_2]$ is given by

$$R_1 = \left\{ \left((y_1, B_1), (y_2, B_2) \right) : y_1 \geq y_2 \wedge B_1 \supseteq B_2 \right\}$$

$$R_2 = \left\{ \left(\left((y_1, B_1), (y_2, B_2) \right), \left((y_3, B_3), (y_4, B_4) \right) \right) : \begin{array}{l} y_1 - y_2 \geq y_3 - y_4 \wedge \\ B_1 \supseteq B_3 \supseteq B_4 \supseteq B_2 \end{array} \right\}$$

“(y_1, B_1) \leftarrow (y_2, B_2) is more desirable than (y_3, B_3) \leftarrow (y_4, B_4) if it yields a higher increase of salary and a **superset increase of benefits.**”

Modelling \mathcal{M} : Credal sets

Credal set

The uncertainty among the elements of S is described by a polyhedral *credal set* of probability measures of the form

$$\mathcal{M} = \left\{ \pi \in \mathcal{P} : \underline{b}_\ell \leq \mathbb{E}_\pi(f_\ell) \leq \bar{b}_\ell \text{ for } \ell = 1, \dots, r \right\}$$

where \mathcal{P} is the set of all probability measures on $(S, \sigma(S))$ and

- $f_1, \dots, f_r : S \rightarrow \mathbb{R}$ are real-valued mappings and
- $\underline{b}_\ell \leq \bar{b}_\ell, \ell = 1, \dots, r$, are lower and upper expectation bounds.

Modelling \mathcal{M} : Credal sets

Credal set

The uncertainty among the elements of S is described by a polyhedral *credal set* of probability measures of the form

$$\mathcal{M} = \left\{ \pi \in \mathcal{P} : \underline{b}_\ell \leq \mathbb{E}_\pi(f_\ell) \leq \bar{b}_\ell \text{ for } \ell = 1, \dots, r \right\}$$

where \mathcal{P} is the set of all probability measures on $(S, \sigma(S))$ and

- $f_1, \dots, f_r : S \rightarrow \mathbb{R}$ are real-valued mappings and
- $\underline{b}_\ell \leq \bar{b}_\ell, \ell = 1, \dots, r$, are lower and upper expectation bounds.

Special cases: *Classical probability – Probability intervals – Interval probability – Linear partial information – (Finitely generated) Lower previsions*

Generalized Choice Functions and Elicitation

Choice functions for decision making based on the sets \mathcal{U}_A and \mathcal{M} as well as efficient computation algorithms have been developed in:



Information Science (2018)

Information-efficient procedures for eliciting optimal decisions according to these criteria are discussed in:



Information Science (2018)

Generalized Stochastic Dominance

Today, we focus on only one choice function from these papers, based on:

Generalized Stochastic Dominance Relation (GSD-Relation)

Let $\mathcal{A} = [A, R_1, R_2]$ be consistent and \mathcal{M} a credal set on (S, \mathcal{S}) .

For $X, Y \in \mathcal{F}_{(\mathcal{A}, S)}$,¹ we say that Y is $(\mathcal{A}, \mathcal{M})$ -dominated by X if

$$\mathbb{E}_\pi(u \circ X) \geq \mathbb{E}_\pi(u \circ Y)$$

for all $u \in \mathcal{U}_A$ and $\pi \in \mathcal{M}$. The induced relation is denoted by $\geq_{(\mathcal{A}, \mathcal{M})}$ and called Generalized Stochastic Dominance Relation (GSD-Relation).

¹ $\mathcal{F}_{(\mathcal{A}, S)} := \{X \in A^S : u \circ X \text{ is } \mathcal{S}\text{-}\mathcal{B}_{\mathbb{R}}([0, 1])\text{-measurable for all } u \in \mathcal{U}_A\}$.

Generalized Stochastic Dominance

Today, we focus on only one choice function from these papers, based on:

Generalized Stochastic Dominance Relation (GSD-Relation)

Let $\mathcal{A} = [A, R_1, R_2]$ be consistent and \mathcal{M} a credal set on (S, \mathcal{S}) .

For $X, Y \in \mathcal{F}_{(\mathcal{A}, S)}$,¹ we say that Y is $(\mathcal{A}, \mathcal{M})$ -dominated by X if

$$\mathbb{E}_\pi(u \circ X) \geq \mathbb{E}_\pi(u \circ Y)$$

for all $u \in \mathcal{U}_A$ and $\pi \in \mathcal{M}$. The induced relation is denoted by $\succeq_{(\mathcal{A}, \mathcal{M})}$ and called Generalized Stochastic Dominance Relation (**GSD-Relation**).

The GSD-relation now directly induces the **GSD choice function** by setting

$$\text{ch}_{\mathcal{A}, \mathcal{M}}(\mathcal{D}) := \left\{ X \in \mathcal{D} : \nexists Y \in \mathcal{D} \text{ such that } (Y, X) \in \succ_{(\mathcal{A}, \mathcal{M})} \right\}$$

¹ $\mathcal{F}_{(\mathcal{A}, S)} := \left\{ X \in A^S : u \circ X \text{ is } \mathcal{S}\text{-}\mathcal{B}_{\mathbb{R}}([0, 1])\text{-measurable for all } u \in \mathcal{U}_A \right\}$.

Some Special Cases of GSD

The GSD-relation $\geq_{(\mathcal{A}, \mathcal{M})}$ has some prominent special cases.

For ...

Some Special Cases of GSD

The GSD-relation $\succeq_{(\mathcal{A}, \mathcal{M})}$ has some prominent special cases.

For ...

- ... and $\mathcal{M} = \{\pi\}$ and R_2 trivial

→ Reduction to (first-order) **stochastic dominance**

(see, e.g., [Mosler and Scarsini, 1991]))

Some Special Cases of GSD

The GSD-relation $\succeq_{(\mathcal{A}, \mathcal{M})}$ has some prominent special cases.

For ...

- ... and $\mathcal{M} = \{\pi\}$ and R_2 trivial
→ Reduction to (first-order) **stochastic dominance**
(see, e.g., [Mosler and Scarsini, 1991]))
- ... and $\mathcal{M} = \{\pi\}$ and R_1 and R_2 guaranteeing utility unique up to plts
→ Reduction to comparing **expected utilities**.
(see, e.g., [Krantz et al., 1971]))

Some Special Cases of GSD

The GSD-relation $\succeq_{(\mathcal{A}, \mathcal{M})}$ has some prominent special cases.

For ...

- ... and $\mathcal{M} = \{\pi\}$ and R_2 trivial
→ Reduction to (first-order) **stochastic dominance**
(see, e.g., [Mosler and Scarsini, 1991]))
- ... and $\mathcal{M} = \{\pi\}$ and R_1 and R_2 guaranteeing utility unique up to plts
→ Reduction to comparing **expected utilities**.
(see, e.g., [Krantz et al., 1971]))
- ... \mathcal{M} non-trivial and R_1 and R_2 guaranteeing utility unique up to plts
→ Reduction to **Bewley dominance**.
(see, e.g., [Bewley, 2002, Troffaes, 2007, Etner et al., 2012]))

Locally Varying Scales of Measurement

Group and collaborators

Most of the following is **joint work** with (in alphabetic order):

- Thomas Augustin,
- Hannah Blocher,
- Malte Nalenz,
- Julian Rodemann,
- Georg Schollmeyer,

and mainly based on the following three papers: [Jansen et al. \(2023\)](#)

C. Jansen, G. Schollmeyer, H. Blocher, J. Rodemann and T. Augustin (2023): **Robust statistical comparison of random variables with locally varying scale of measurement**. In: Proceedings of the Thirty-Ninth Conference on Uncertainty in Artificial Intelligence (**UAI 2023**). Proceedings of Machine Learning Research, vol. 216.

C. Jansen, M. Nalenz, G. Schollmeyer and T. Augustin (2023): **Statistical comparisons of classifiers by generalized stochastic dominance**. Journal of Machine Learning Research (**JMLR**), 24 (231): 1 - 37.

C. Jansen, G. Schollmeyer, J. Rodemann, H. Blocher and T. Augustin (2024): **Statistical multicriteria benchmarking via the GSD-front**. Under review.

Motivation

1.) Statistical methods are usually tailored for data situations that can be clearly assigned to a **standard scale of measurement**.

2.) **Non-standard data** can often not clearly be assigned to a standard scale.

1.)+2.) \Rightarrow Statistical methods are often not well-suited for analyzing non-standard data!

Idea: Use the notion of a preference system to model data with scales of measurement which not correspond to one of these extreme poles.

Standard Scales of Measurement

Consider some random variable $X : \Omega \rightarrow A$ mapping to some set A .

Standard Scales of Measurement

Consider some random variable $X : \Omega \rightarrow A$ mapping to some set A .

- If A is structured only by a **preorder** \succsim , we call A of **ordinal** scale.
 - \Rightarrow The set \mathcal{U}_{all} of all \succsim -isotone **candidate scales** $u : A \rightarrow \mathbb{R}$ as a **whole** represents the structural information on A .
 - \Rightarrow Any analysis of the variable X should be **invariant** under the choice of the candidate scale $u \in \mathcal{U}_{all}$.

Standard Scales of Measurement

Consider some random variable $X : \Omega \rightarrow A$ mapping to some set A .

- If A is structured only by a **preorder** \succsim , we call A of **ordinal** scale.
 - ⇒ The set \mathcal{U}_{all} of all \succsim -isotone **candidate scales** $u : A \rightarrow \mathbb{R}$ as a **whole** represents the structural information on A .
 - ⇒ Any analysis of the variable X should be **invariant** under the choice of the candidate scale $u \in \mathcal{U}_{all}$.
- If the order on A is induced by some **metric** d , we call A of **cardinal** scale.
 - ⇒ There exists a scale $u^* : A \rightarrow \mathbb{R}$ that is **unique** (up to irrelevant trasfos).
 - ⇒ Any analysis of the variable X can be based on u^* **alone**.

Preference Systems in Statistics

Question: What if the structure on A does not belong to either extreme pole?

In other words: What if the structuredness of A **varies along its subsets**?

Preference Systems in Statistics

Question: What if the structure on A does not belong to either extreme pole?

In other words: What if the structuredness of A **varies along its subsets**?

A preference system $\mathcal{A} = [A, R_1, R_2]$ helps to formalize this intuition:

- R_1 formalizes the available **ordinal information**, i.e. information about the arrangement of the elements of A .
- R_2 describes the available **cardinal information**, i.e. pairs standing in relation are ordered with respect to the intensity of the relation.
- A is **locally almost cardinal** on subsets where R_1 and R_2 are very dense.
- A is **locally at most ordinal** on subsets where R_2 is sparse or even empty.

Regularization and Preference Systems

Opportunity: Preference systems offer a nice way for **regularization** by excluding those $u \in \mathcal{U}_A$ that are too extreme (in some sense).

Regularization and Preference Systems

Opportunity: Preference systems offer a nice way for **regularization** by excluding those $u \in \mathcal{U}_{\mathcal{A}}$ that are too extreme (in some sense).

Simple idea: If \mathcal{A} has R_1 -minimal/maximal elements a_*, a^* , define

$$\mathcal{N}_{\mathcal{A}} := \left\{ u \in \mathcal{U}_{\mathcal{A}} : u(a_*) = 0 \wedge u(a^*) = 1 \right\}$$

$$\mathcal{N}_{\mathcal{A}}^{\delta} := \left\{ u \in \mathcal{N}_{\mathcal{A}} : u(c) - u(d) - u(e) + u(f) \geq \delta \quad \forall ((c, d), (e, f)) \in P_{R_2} \right\}$$

Regularization and Preference Systems

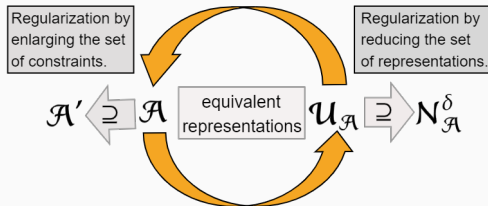
Opportunity: Preference systems offer a nice way for **regularization** by excluding those $u \in \mathcal{U}_{\mathcal{A}}$ that are too extreme (in some sense).

Simple idea: If \mathcal{A} has R_1 -minimal/maximal elements a_*, a^* , define

$$\mathcal{N}_{\mathcal{A}} := \{u \in \mathcal{U}_{\mathcal{A}} : u(a_*) = 0 \wedge u(a^*) = 1\}$$

$$\mathcal{N}_{\mathcal{A}}^{\delta} := \{u \in \mathcal{N}_{\mathcal{A}} : u(c) - u(d) - u(e) + u(f) \geq \delta \quad \forall ((c, d), (e, f)) \in P_{R_2}\}$$

Two ways for regularization:



Random Variables Mapping Into Preference Systems

Goal: We now want to address the problem of comparing random variables $X, Y : \Omega \rightarrow A$ that map into a preference system.

Challenge: We have epistemic uncertainty in form of

- **Approximation uncertainty:** Only samples of the considered variables (rather than π itself) are available.
- **Model uncertainty:** The weakly structured order information makes a set of candidate scales compatible with the structure on A .

Addressing Model Uncertainty via GSD

Idea: Weaken $\succsim_{E(u)}$ to a **preorder** by demanding expectation dominance for all scales u compatible with the preference system \mathcal{A} .

\Rightarrow This idea leads to a "precise" version of GSD.

Recall:

Precise GSD

Let \mathcal{A} be consistent and π be a probability measure on (S, \mathcal{S}) .

For $X, Y \in \mathcal{F}_{(\mathcal{A}, S)}$, we call Y $(\mathcal{A}, \{\pi\})$ -dominated by X if

$$\mathbb{E}_\pi(u \circ X) \geq \mathbb{E}_\pi(u \circ Y)$$

for all $u \in \mathcal{U}_{\mathcal{A}}$. This induces preorder $R_{(\mathcal{A}, \pi)}$ on $\mathcal{F}_{(\mathcal{A}, \{\pi\})}$ which is called the **precise GSD-relation**.

Obviously, precise GSD is invariant under the scale.

Addressing Approximation Uncertainty

Practical Problem: Usually, we do not know π but only *i.i.d.* samples $\mathbf{X} = (X_1, \dots, X_n)$ and $\mathbf{Y} = (Y_1, \dots, Y_m)$ of X and Y are available.

Approach: Perform a statistical test for GSD.

Ideal Hypotheses:

$$H_0^{id} : (X, Y) \notin R_{(\mathcal{A}, \pi)} \quad \text{vs.} \quad H_1^{id} : (X, Y) \in R_{(\mathcal{A}, \pi)}$$

Pragmatic Hypotheses:

$$H_0 : (Y, X) \in R_{(\mathcal{A}, \pi)} \quad \text{vs.} \quad H_1 : (Y, X) \notin R_{(\mathcal{A}, \pi)}$$

Addition: To mitigate the effect of the reversed hypotheses, we can additionally test with the variables X and Y in reversed roles.

The Choice of the Test Statistic

Observation: It holds $(X, Y) \in R_{(\mathcal{A}, \pi)}$ if and only if

$$D(X, Y) := \inf_{u \in \mathcal{N}_{\mathcal{A}}} (\mathbb{E}_{\pi}(u \circ X) - \mathbb{E}_{\pi}(u \circ Y)) \geq 0.$$

Consequence: A natural test statistic is the empirical version of $D(X, Y)$, i.e.,

$$d_{X, Y} : \Omega \rightarrow \mathbb{R}$$

$$\omega \mapsto \inf_{u \in \mathcal{N}_{\mathcal{A}_{\omega}}} \sum_{z \in (\mathbf{XY})_{\omega}} u(z) \cdot (\hat{\pi}_X^{\omega}(\{z\}) - \hat{\pi}_Y^{\omega}(\{z\}))$$

with, for $\omega \in \Omega$ fixed,

- $\hat{\pi}_X^{\omega}$ and $\hat{\pi}_Y^{\omega}$ the observed empirical image measures of X and Y ,
- $(\mathbf{XY})_{\omega} = \{X_i(\omega) : i \leq n\} \cup \{Y_i(\omega) : i \leq m\} \cup \{a_*, a^*\}$, and
- \mathcal{A}_{ω} the subsystem of \mathcal{A} restricted to $(\mathbf{XY})_{\omega}$.

Regularization of the Test Statistic

Observation: $d_{X,Y}$ cannot measure **extent** of GSD in the sample. Thus, $d_{X,Y}$ may be too little sensitive.

Idea: Regularize $d_{X,Y}$ so that it can also account for the **extent** of GSD.

Formally: The regularized test statistic looks as follows:

$$d_{X,Y}^\varepsilon : \Omega \rightarrow \mathbb{R}$$
$$\omega \mapsto \inf_{u \in \mathcal{N}_{\mathcal{A}_\omega}^{\delta_\varepsilon(\omega)}} \sum_{z \in (XY)_\omega} u(z) \cdot (\hat{\pi}_X^\omega(\{z\}) - \hat{\pi}_Y^\omega(\{z\}))$$

with $\varepsilon \in [0, 1]$ and

$$\delta_\varepsilon(\omega) := \varepsilon \cdot \sup\{\xi : \mathcal{N}_{\mathcal{A}_\omega}^\xi \neq \emptyset\}.$$

Computation: Both test statistics $d_{X,Y}$ and $d_{X,Y}^\varepsilon$ can be computed by solving one single **linear programming problem**.

A Permutation Test

Assumption: We made observations of the i.i.d. variables, i.e., we observed:

$$\mathbf{x} := (X_1, \dots, X_n) := (X_1(\omega_0), \dots, X_n(\omega_0))$$

$$\mathbf{y} := (Y_1, \dots, Y_m) := (Y_1(\omega_0), \dots, Y_m(\omega_0))$$

A Permutation Test

Assumption: We made observations of the i.i.d. variables, i.e., we observed:

$$\mathbf{x} := (X_1, \dots, X_n) := (X_1(\omega_0), \dots, X_n(\omega_0))$$

$$\mathbf{y} := (Y_1, \dots, Y_m) := (Y_1(\omega_0), \dots, Y_m(\omega_0))$$

Good News: As the worst case of the null hypothesis H_0 is $\pi_X = \pi_Y$, performing a **permutation test** is a valid level α test.

The resampling scheme then looks:

Step 1: Pool data sample: $\mathbf{w} := (w_1, \dots, w_{n+m}) := (x_1, \dots, x_n, y_1, \dots, y_m)$

Step 2: Take all $k := \binom{n+m}{n}$ index sets $I \subseteq \{1, \dots, n+m\}$ of size n . Compute $d_{X,Y}$ resp. $d_{X,Y}^\varepsilon$ for $(w_i)_{i \in I}$ and $(w_i)_{i \in \{1, \dots, n+m\} \setminus I}$ instead of \mathbf{x}/\mathbf{y} to get d_I resp. d_I^ε .

Step 3: Sort all d_I resp. d_I^ε in increasing order to get $d_{(1)}, \dots, d_{(k)}$ resp. $d_{(1)}^\varepsilon, \dots, d_{(k)}^\varepsilon$.

Step 4: Reject H_0 if $d_{X,Y}(\omega_0)$ resp. $d_{X,Y}^\varepsilon(\omega_0)$ is greater than $d_{(\ell)}$ resp. $d_{(\ell)}^\varepsilon$, with $\ell := \lceil (1 - \alpha) \cdot k \rceil$ and α the significance level.

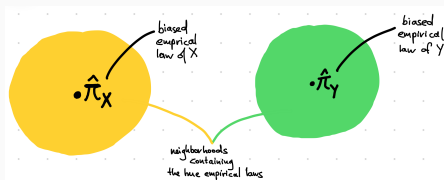
Credal Sets For Robustification

Rough Idea: Use **credal sets** to robustify the permutation test to small deviations from the *i.i.d.* assumption.

Credal Sets For Robustification

Rough Idea: Use **credal sets** to robustify the permutation test to small deviations from the *i.i.d.* assumption.

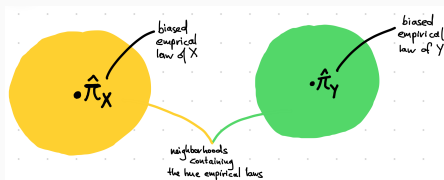
More concrete: We allow our samples to be (potentially) **biased** in the sense that we only assume the **true empirical laws** to lie in some credal neighborhoods \mathcal{M}_X and \mathcal{M}_Y around the **biased empirical laws**.



Credal Sets For Robustification

Rough Idea: Use **credal sets** to robustify the permutation test to small deviations from the *i.i.d.* assumption.

More concrete: We allow our samples to be (potentially) **biased** in the sense that we only assume the **true empirical laws** to lie in some credal neighborhoods \mathcal{M}_X and \mathcal{M}_Y around the **biased empirical laws**.



Adapted Resampling Scheme: Replace

- $d_{X,Y}^{\varepsilon}(\omega_0)$ by $\inf_{(\pi_1, \pi_2) \in \mathcal{M}_X^{\omega_0} \times \mathcal{M}_Y^{\omega_0}} \tilde{d}_{X,Y}^{\varepsilon}(\omega_0)$
- $d_I^{\varepsilon}(\omega_0)$ by $\sup_{(\pi_1, \pi_2) \in \mathcal{M}_X^{\omega_0} \times \mathcal{M}_Y^{\omega_0}} \tilde{d}_I^{\varepsilon}(\omega_0)$

Results in: Valid (yet conservative) statistical test!

A special class of credal sets with a very intuitive interpretation are

γ -contamination models

For $\omega \in \Omega$, $\gamma \in [0, 1]$, and $Z \in \{X, Y\}$ fixed, we set

$$\mathcal{M}_Z^\omega = \left\{ \pi : \pi \geq (1 - \gamma) \cdot \hat{\pi}_Z^\omega \right\}$$

or equivalently

$$\mathcal{M}_Z^\omega = \left\{ \gamma \cdot \nu + (1 - \gamma) \cdot \hat{\pi}_Z^{\omega_0} : \nu \text{ probability measure} \right\}.$$

γ -Contamination Model

A special class of credal sets with a very intuitive interpretation are

γ -contamination models

For $\omega \in \Omega$, $\gamma \in [0, 1]$, and $Z \in \{X, Y\}$ fixed, we set

$$\mathcal{M}_Z^\omega = \left\{ \pi : \pi \geq (1 - \gamma) \cdot \hat{\pi}_Z^\omega \right\}$$

or equivalently

$$\mathcal{M}_Z^\omega = \left\{ \gamma \cdot \nu + (1 - \gamma) \cdot \hat{\pi}_Z^{\omega_0} : \nu \text{ probability measure} \right\}.$$

The observed **p-values** of the robustified test can then be computed as a function of the contamination size γ :

$$f_\varepsilon(\gamma) := 1 - \frac{1}{N} \cdot \sum_{I \in \mathcal{I}_N} \mathbb{1} \left\{ d_{X,Y}^\varepsilon(\omega_0) - d_I^\varepsilon > \frac{2\gamma}{(1-\gamma)} \right\}$$

Application I

Spaces with Differently Scaled Dimensions (SDSDs)

Situation: Consider an r -dimensional space $A \subseteq \mathbb{R}^r$ and assume that

- the first $0 \leq z \leq r$ dimensions are of cardinal scale and
- the remaining dimensions are purely ordinal.

Question: How can we utilize the cardinal dimensions without making unjustified assumptions about the ordinal ones?

Idea: Utilize the cardinal information only on those parts of A where there is no possible conflict with the ordinal information.

Formalization: Consider A to be a subsystem of $\mathcal{P} = [\mathbb{R}^r, R_1^*, R_2^*]$, where

$$R_1^* = \{(x, y) : x_j \geq y_j \forall j \leq r\}$$
$$R_2^* = \left\{ ((x, y), (x', y')) : \begin{array}{l} x_j - y_j \geq x'_j - y'_j \quad \forall j \leq z \\ x_j \geq x'_j \geq y'_j \geq y_j \quad \forall j > z \end{array} \right\}.$$

A Characterization Theorem in SDSDs

For the special case of A being a multidimensional space with differently scaled dimensions, the GSD-relation can be neatly characterized.

Theorem

Let $X = (\Delta_1, \dots, \Delta_r), Y = (\Lambda_1, \dots, \Lambda_r) \in \mathcal{F}_{(\mathcal{P}, \pi)}$. Then:

- i) \mathcal{P} is consistent.
- ii) If $z = 0$, then $R_{(\mathcal{P}, \pi)}$ equals (first-order) stochastic dominance w.r.t. π and R_1^* (short: $\text{FSD}(R_1^*, \pi)$).
- iii) If $(X, Y) \in R_{(\mathcal{P}, \pi)}$ and $\Delta_j, \Lambda_j \in \mathcal{L}^1(\Omega, \mathcal{S}_1, \pi)$ for all $j = 1, \dots, r$, then
 - I. $\mathbb{E}_\pi(\Delta_j) \geq \mathbb{E}_\pi(\Lambda_j)$ for all $j = 1, \dots, r$, and
 - II. $(\Delta_j, \Lambda_j) \in \text{FSD}(\geq, \pi)$ for all $j = z + 1, \dots, r$.

If all components of X are jointly independent and all components of Y are jointly independent, I. and II. imply $(X, Y) \in R_{(\mathcal{P}, \pi)}$.

Multidimensional Poverty Analysis

Capability Approach: Poverty is a multidimensional concept with more facets than just income or wealth ([Sen, 1985]).

Exemplary operationalization: We use the ALLBUS data and account for three dimensions of poverty: *income* (numeric), *health* (ordinal, 6 levels) and *education* (ordinal, 8 levels)

Example:



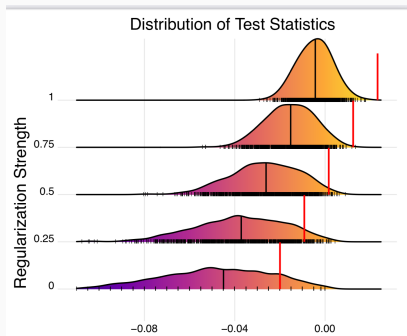
$$\langle \star, \heartsuit \rangle \in \mathbb{R}_2^*$$

$$\begin{aligned} 2500 - 1000 &\geq 2000 - 800 \\ \text{good} &\geq \text{okay} \geq \text{bad} \geq \text{very bad} \\ \text{M.Sc.} &\geq \text{B.Sc.} \geq \text{High School} \geq \text{elementary} \end{aligned}$$

Multidimensional Poverty Analysis, cont.

For the ALLBUS data, we focus on a subsample with $n = m = 100$ men and women each. Again, we operationalize poverty by the variables *income* (numeric), *health* (ordinal, 6 levels) and *education* (ordinal, 8 levels)

Test results:

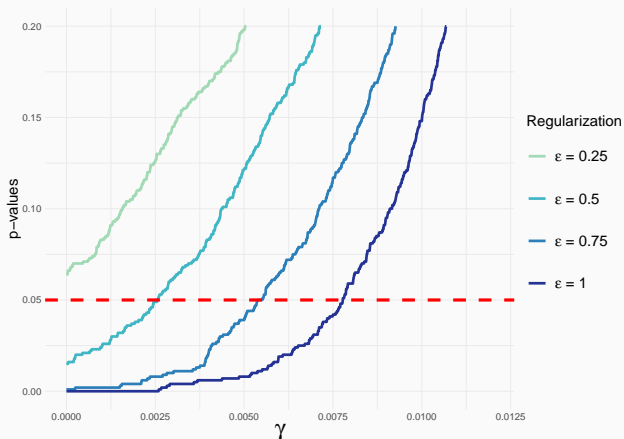


Results: All tests significant for $\alpha = 0.05$.

Reversed test: No evidence for incomparability: All reversed p -values ≥ 0.95 .

Multidimensional Poverty Analysis, cont.

Results of the robustified test:



Application II

Statistical Multicriteria Comparison of Classifiers

Question of interest: How to utilize our decision-theoretical approach for comparing classifiers under **multiplicity** of **quality criteria** and **data sets**?

Setup: Let

- \mathcal{D} denote the set of all relevant **data sets**,
- \mathcal{C} denote the finite set of all relevant **classifiers**,
- $(\phi_i : \mathcal{C} \times \mathcal{D} \rightarrow [0, 1])_{i \in \{1, \dots, r\}}$ denote a family of **quality criteria**,
- $\phi := (\phi_1, \dots, \phi_r) : \mathcal{D} \times \mathcal{C} \rightarrow [0, 1]^r$ be a **multidimensional quality criterion**.

Statistical Multicriteria Comparison of Classifiers

Question of interest: How to utilize our decision-theoretical approach for comparing classifiers under **multiplicity** of **quality criteria** and **data sets**?

Setup: Let

- \mathcal{D} denote the set of all relevant data sets,
- \mathcal{C} denote the finite set of all relevant classifiers,
- $(\phi_i : \mathcal{C} \times \mathcal{D} \rightarrow [0, 1])_{i \in \{1, \dots, r\}}$ denote a family of quality criteria,
- $\phi := (\phi_1, \dots, \phi_r) : \mathcal{D} \times \mathcal{C} \rightarrow [0, 1]^r$ be a multidimensional quality criterion.

Assumptions:

- For $0 \leq z \leq r$, the criteria ϕ_1, \dots, ϕ_z are of cardinal scale (differences may be interpreted)
- The remaining criteria are purely ordinal (differences are meaningless apart from sign).

Statistical Multicriteria Comparison of Classifiers

Three **levels of problems** when comparing classifiers w.r.t. multiple quality criteria on multiple data sets simultaneously.

classifier \ data sets	D_1	...	D_s
C_1	$\begin{pmatrix} \phi_1(C_1, D_1) \\ \vdots \\ \phi_n(C_1, D_1) \end{pmatrix}$...	$\begin{pmatrix} \phi_1(C_1, D_s) \\ \vdots \\ \phi_n(C_1, D_s) \end{pmatrix}$
\vdots	\vdots	\vdots	\vdots
C_q	$\begin{pmatrix} \phi_1(C_q, D_1) \\ \vdots \\ \phi_n(C_q, D_1) \end{pmatrix}$...	$\begin{pmatrix} \phi_1(C_q, D_s) \\ \vdots \\ \phi_n(C_q, D_s) \end{pmatrix}$

Statistical Multicriteria Comparison of Classifiers

Three **levels of problems** when comparing classifiers w.r.t. multiple quality criteria on multiple data sets simultaneously.

		data sets		
		D_1	...	D_s
classifier	C_1	$\begin{pmatrix} 0.8 \\ \vdots \\ 0.7 \end{pmatrix}$...	$\begin{pmatrix} \phi_1(C_1, D_s) \\ \vdots \\ \phi_n(C_1, D_s) \end{pmatrix}$
	\vdots	\vdots	\vdots	\vdots
	C_q	$\begin{pmatrix} 0.7 \\ \vdots \\ 0.8 \end{pmatrix}$...	$\begin{pmatrix} \phi_1(C_q, D_s) \\ \vdots \\ \phi_n(C_q, D_s) \end{pmatrix}$

Level 1: On a fixed data set D it may hold

$$\phi_1(C_1, D) > \phi_1(C_2, D) \wedge \phi_2(C_1, D) < \phi_2(C_2, D).$$

Statistical Multicriteria Comparison of Classifiers

Three **levels of problems** when comparing classifiers w.r.t. multiple quality criteria on multiple data sets simultaneously.

		data sets		
		D_1	...	D_s
classifier	C_1	$\begin{pmatrix} 0.8 \\ \vdots \\ 0.8 \end{pmatrix}$...	$\begin{pmatrix} 0.6 \\ \vdots \\ \phi_n(C_1, D_s) \end{pmatrix}$
	\vdots	\vdots	\vdots	\vdots
	C_q	$\begin{pmatrix} 0.7 \\ \vdots \\ 0.7 \end{pmatrix}$...	$\begin{pmatrix} 0.9 \\ \vdots \\ \phi_n(C_q, D_s) \end{pmatrix}$

Level 2: Even if, for all $i \in \{1, \dots, n\}$, we have

$$\phi_i(C_1, D_1) > \phi_i(C_2, D_1)$$

there may exist some $i_0 \in \{1, \dots, n\}$ such that

$$\phi_{i_0}(C_1, D_2) < \phi_{i_0}(C_2, D_2).$$

Statistical Multicriteria Comparison of Classifiers

Three **levels of problems** when comparing classifiers w.r.t. multiple quality criteria on multiple data sets simultaneously.

		data sets		
		D_1	\dots	D_S
classifier	C_1	$\begin{pmatrix} 0.8 \\ \vdots \\ 0.8 \end{pmatrix}$	\dots	$\begin{pmatrix} 0.8 \\ \vdots \\ 0.8 \end{pmatrix}$
	\vdots	\vdots	\vdots	\vdots
	C_q	$\begin{pmatrix} 0.7 \\ \vdots \\ 0.7 \end{pmatrix}$	\dots	$\begin{pmatrix} 0.7 \\ \vdots \\ 0.7 \end{pmatrix}$

Level 3: Even if a decision can be made for a sample (D_1, \dots, D_S) of data sets,

Statistical Multicriteria Comparison of Classifiers

Three **levels of problems** when comparing classifiers w.r.t. multiple quality criteria on multiple data sets simultaneously.

		data sets		
		D_1^*	...	D_s^*
classifier	C_1	$\begin{pmatrix} 0.7 \\ \vdots \\ 0.9 \end{pmatrix}$...	$\begin{pmatrix} 0.75 \\ \vdots \\ 0.4 \end{pmatrix}$
	\vdots	\vdots	\vdots	\vdots
	C_q	$\begin{pmatrix} 0.85 \\ \vdots \\ 0.67 \end{pmatrix}$...	$\begin{pmatrix} 0.33 \\ \vdots \\ 0.98 \end{pmatrix}$

Level 3: Even if a decision can be made for a sample (D_1, \dots, D_s) of data sets, no clear decision might be possible for a different sample (D_1^*, \dots, D_s^*) .

Transferring GSD to Classifier Comparison

Idea: Embed the range $\Phi(\mathcal{C} \times \mathcal{D})$ of Φ in the following preference system $\mathcal{P} = [\mathbb{R}^r, R_1^*, R_2^*]$ from before.

Then:

- To transfer the GSD-relation, interpret the data sets in \mathcal{D} as realizations of a random variable $T : \Omega \rightarrow \mathcal{D}$ on some probability space $(\Omega, \mathcal{S}, \pi)$.
- Associate each $C \in \mathcal{C}$ with the variable $\Phi_C := \Phi(C, T(\cdot))$ on Ω and compare classifiers by comparing the associated random variables by precise GSD.

Formally:

GSD for Classifier Comparison

Let \mathcal{P}_Φ be the preference system obtained by restricting \mathcal{P} to $\Phi(\mathcal{C} \times \mathcal{D})$.

Further, let \mathcal{C} be such that $\{\Phi_C : C \in \mathcal{C}\} \subseteq \mathcal{F}_{(\mathcal{P}_\Phi, \pi)}$.

For $C, C' \in \mathcal{C}$, say that C **dominates** C' , abbreviated $C \succsim C'$, whenever

$$(\Phi_C, \Phi_{C'}) \in R_{(\mathcal{P}_\Phi, \pi)}.$$

Theoretical and Empirical GSD-Front

We associate the following two sets to the relation \succsim :

The GSD-Front

Let \mathcal{C} be such that $\{\Phi_C : C \in \mathcal{C}\} \subseteq \mathcal{F}_{(\mathcal{P}_\Phi, \pi)}$ and T_1, \dots, T_s be *i.i.d.* copies of T .

i) The **GSD-front** is the set

$$\text{gsd}(\mathcal{C}) := \{C \in \mathcal{C} : \nexists C' \in \mathcal{C} \text{ s.t. } C' \succ C\},$$

where \succ denotes the strict part of \succsim .

ii) Let $\rho \in [0, 1]$. The ρ -**empirical GSD-front** is the (random) subset of \mathcal{C} defined by

$$\text{egsd}_s^\rho(\mathcal{C}) = \left\{ C : \nexists C' \in \mathcal{C} \text{ s.t. } \begin{array}{l} d_{(\Phi_{C'}, \Phi_C)} \geq -\rho \\ d_{(\Phi_C, \Phi_{C'})} < 0 \end{array} \right\}.$$

Consistent Estimability of the GSD-Front

The following theorem on the consistent estimability of the GSD-front holds:

Estimating the GSD-Front

Denote by \mathcal{I}_Φ the set of all sets $\{a : u(a) \geq c\}$, where $c \in [0, 1]$ and $u \in \mathcal{U}_{\mathcal{P}_\Phi}$.

Assume that \succsim is antisymmetric.

If the VC-dimension² of \mathcal{I}_Φ is finite and $\rho : \mathbb{N} \rightarrow [0, 1]$ converges to 0 with at most $\Theta(1/\sqrt[4]{s})$, then $(\text{egsd}_s^{\rho(s)}(\mathcal{C}))_{s \in \mathbb{N}}$ is consistent, i.e.,

$$\pi \left(\left\{ \omega \in \Omega : \lim_{s \rightarrow \infty} \text{egsd}_s^{\rho(s)}(\mathcal{C}) = \text{gsd}(\mathcal{C}) \right\} \right) = 1,$$

where set convergence is defined via the trivial metric.

²The VC-dimension of a family of sets \mathcal{S} is the largest possible cardinality of a set A , such that $2^A = \{A \cap S : S \in \mathcal{S}\}$, i.e., A can be shattered by \mathcal{S} .

Consistent Tests for the GSD-Front

Goal: Compare the (multivariate, mixed-scaled) quality of a newly developed classifier C with a set \mathcal{C} of state-of-the-art classifiers.

Consistent Tests for the GSD-Front

Goal: Compare the (multivariate, mixed-scaled) quality of a newly developed classifier C with a set \mathcal{C} of state-of-the-art classifiers.

How to proceed? Develop a statistical test for the pair

$$H_0 : C \notin \text{gsd}(\mathcal{C}) \text{ vs. } H_1 : C \in \text{gsd}(\mathcal{C})$$

Consistent Tests for the GSD-Front

Goal: Compare the (multivariate, mixed-scaled) quality of a newly developed classifier C with a set \mathcal{C} of state-of-the-art classifiers.

How to proceed? Develop a statistical test for the pair

$$H_0 : C \notin \text{gsd}(\mathcal{C}) \text{ vs. } H_1 : C \in \text{gsd}(\mathcal{C})$$

How exactly? Note that H_0 can be rewritten as:

$$H_0 : \exists C' \in \mathcal{C} \setminus \{C\} : C' \succsim C.$$

Thus, H_0 is false iff the hypothesis $H_0^{C'} : C' \succsim C$ is false for every $C' \in \mathcal{C} \setminus \{C\}$.

Consistent Tests for the GSD-Front

Goal: Compare the (multivariate, mixed-scaled) quality of a newly developed classifier C with a set \mathcal{C} of state-of-the-art classifiers.

How to proceed? Develop a statistical test for the pair

$$H_0 : C \notin \text{gsd}(\mathcal{C}) \text{ vs. } H_1 : C \in \text{gsd}(\mathcal{C})$$

How exactly? Note that H_0 can be rewritten as:

$$H_0 : \exists C' \in \mathcal{C} \setminus \{C\} : C' \succsim C.$$

Thus, H_0 is false iff the hypothesis $H_0^{C'} : C' \succsim C$ is false for every $C' \in \mathcal{C} \setminus \{C\}$.

Good news:

- The pairs $(H_0^{C'}, \neg H_0^{C'})$ can be tested using the test from Application I.
- Thus, $(H_0, \neg H_0)$ can (essentially) be tested by running these tests multiple times, while rejecting H_0 if all $H_0^{C'}$ are rejected.
- This even allows to construct **consistent** tests.

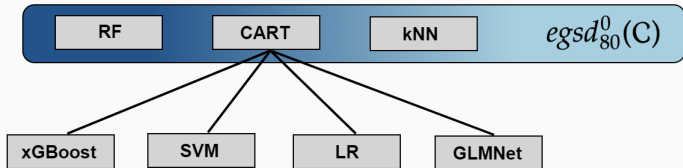
OpenML Benchmarking Experiments: Setup

- We use 80 binary classification datasets from the Open Multimedia Library (OpenML) [Van Rijn et al., 2013].
- We compare the performance of *Support Vector Machine* (SVM) with
 - *Random Forest* (RF),
 - *Decision Tree* (CART),
 - *Logistic Regression* (LR),
 - *Generalized Linear Model with Elastic net* (GLMNet),
 - *Extreme Gradient Boosting* (xGBoost), and
 - *k-Nearest Neighbors* (kNN).
- Comparison is based on the multivariate metric Φ composed of
 - *predictive accuracy*,
 - *computation time on the test data*, and
 - *computation time on the training data*.

Since computation time strongly depends on the computing environment used, we treat the time-related metrics as purely ordinal.

OpenML Benchmarking Experiments: Empirical GSD-Front

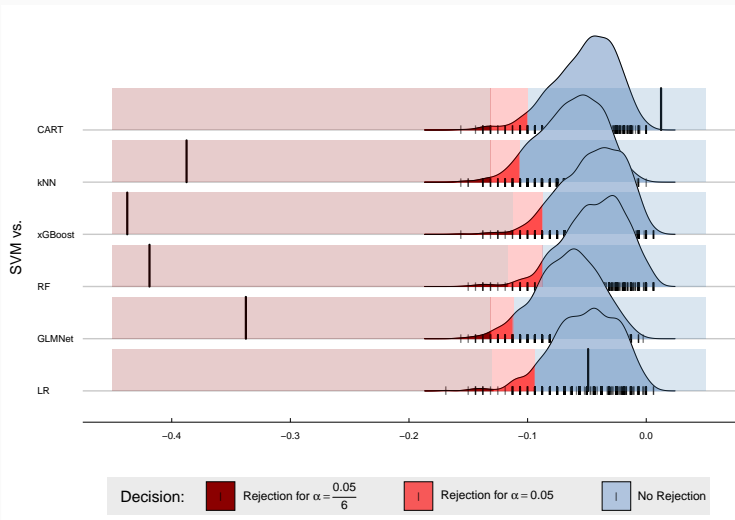
The Hasse graph of the empirical GSD relation:



The blue shaded region symbolizes the 0-empirical GSD-front.

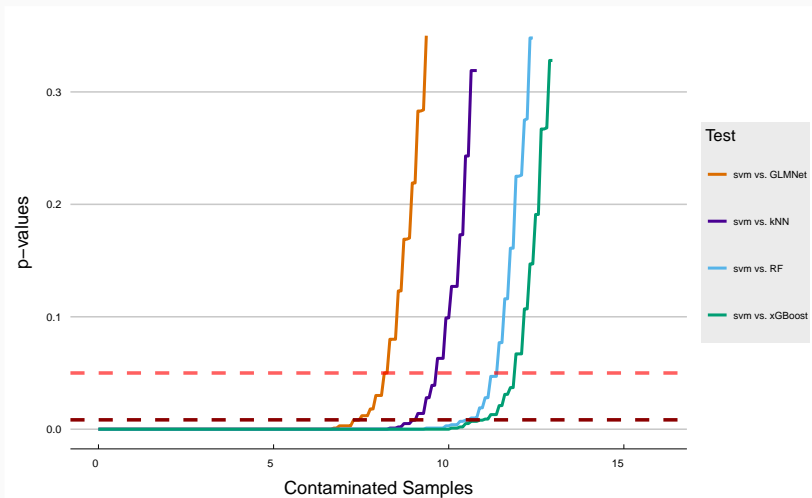
OpenML Benchmarking Experiments: Tests for GSD-Front

Results of the GSD-front test:



OpenML Benchmarking Experiments: Robustness

Robustness of test decision under contamination of the benchmark suite:



Summary and Outlook

Summary:

- Presented a framework for decision making under weakly structured information
- Demonstrated two applications of this framework in problems of robust statistics and machine learning

What is next?

- Exploit other problems/fields where a decision-theoretic perspective might be fruitful

Thank you very much for your attention!



Bacharach, M. (1975).

Group decisions in the face of differences of opinion.

Management Science, 22:182–191.



Bewley, T. F. (2002).

Knightian decision theory. part i.

Decisions in economics and finance, 25:79–110.



Bradley, S. (2019).

Aggregating belief models.

In *Proceedings of ISIPTA 2019*, Proceedings of Machine Learning Research.

-  Etner, J., Jeleva, M., and Tallon, J.-M. (2012).
Decision theory under ambiguity.
Journal of Economic Surveys, 26(2):234–270.
-  Jansen, C., Blocher, H., Augustin, T., and Schollmeyer, G. (2022).
Information efficient learning of complexly structured preferences: Elicitation procedures and their application to decision making under uncertainty.
International Journal of Approximate Reasoning, 144:69–91.
-  Jansen, C., Schollmeyer, G., and Augustin, T. (2018).
Concepts for decision making under severe uncertainty with partial ordinal and partial cardinal preferences.
International Journal of Approximate Reasoning, 98:112–131.



Jansen, C., Schollmeyer, G., Blocher, H., Rodemann, J., and Augustin, T. (2023).

Robust statistical comparison of random variables with locally varying scale of measurement.

In *Uncertainty in Artificial Intelligence (UAI)*. PMLR.

To appear.



Kikuti, D., Cozman, F., and Filho, R. (2011).

Sequential decision making with partially ordered preferences.

Artificial Intelligence, 175:1346 – 1365.



Krantz, D., Luce, R., Suppes, P., and Tversky, A. (1971).

Foundations of Measurement. Volume I: Additive and Polynomial Representations.

Academic Press, San Diego and London.



Levi, I. (1974).

On indeterminate probabilities.

The Journal of Philosophy, 71:391–418.



Mosler, K. and Scarsini, M. (1991).

Some theory of stochastic dominance.

In Mosler, K. and Scarsini, M., editors, *Stochastic Orders and Decision under Risk*, pages 203–212. Institute of Mathematical Statistics, Hayward, CA.



Nau, R. (2006).

The shape of incomplete preferences.

Annals of Statistics, 34:2430–2448.



Savage, L. (1954).

The Foundations of Statistics.

Wiley.



Seidenfeld, T., Kadane, J., and Schervish, M. (1995).

A representation of partially ordered preferences.

Annals of Statistics, 23:2168–2217.



Sen, A. (1985).

Commodities and Capabilities.

Elsevier.



Troffaes, M. (2007).

Decision making under uncertainty using imprecise probabilities.

International Journal of Approximate Reasoning, 45:17–29.



Van Rijn, J. N., Bischl, B., Torgo, L., Gao, B., Umaashankar, V., Fischer, S., Winter, P., Wiswedel, B., Berthold, M. R., and Vanschoren, J. (2013).

Openml: A collaborative science platform.

In Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2013, Prague, Czech Republic, September 23–27, 2013, Proceedings, Part III 13, pages 645–649. Springer.



von Neumann, J., Morgenstern, O., Kuhn, H., and Rubinstein, A. (1944).

Theory of Games and Economic Behavior (60th Anniversary Commemorative Edition).

Princeton University Press.



Walley, P. (1991).

Statistical Reasoning with Imprecise Probabilities.

Chapman and Hall, London.