

Explainable AI and explainable multiobjective optimization

Giovanni Misitano¹

`giovanni.a.misitano@jyu.fi`
`giovanni.misitano.xyz`

¹University of Jyväskylä (Finland), The Multiobjective Optimization Group

Presented remotely during the CCLS seminar on November, 14, 2022.



JYVÄSKYLÄN YLIOPISTO
UNIVERSITY OF JYVÄSKYLÄ



Overview

- 1 Motivation
- 2 Background
 - Multiobjective optimization
 - Scalarization
 - Reference point based interactive methods
 - Explainable Artificial Intelligence
 - Shapley values
- 3 Explainable interactive multiobjective optimization
- 4 Tests and results
- 5 What next?
- 6 Conclusions
- 7 Appendices

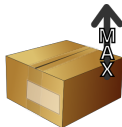
- 1 Motivation
- 2 Background
- 3 Explainable interactive multiobjective optimization
- 4 Tests and results
- 5 What next?
- 6 Conclusions
- 7 Appendices

Multiobjective optimization

- Real-life problems often consist of multiple conflicting objectives.
- These problems have many compromise, non-comparable solutions with various trade-offs.
- A domain expert, known as the decision maker, is needed to find the *best* solution.
- The decision maker can provide preferences, which are used to find the best solution.



Income



Inventory



Efficiency



Time



Pollution

Multiobjective optimization

- Multiobjective optimization methods support the decision maker in finding the best solution.
- The solution is then used in real-life decision-making.
- Often decision makers lack support in providing preferences.
- Can the decision maker *trust* the solution found? Can the solution be *justified* in any way?



Explainability and multiobjective optimization

- Could we make multiobjective optimization methods explainable?
- **Idea:** borrow existing techniques from explainable artificial intelligence (XAI).
- We will explore a new paradigm: explainable (interactive) multiobjective optimization.

Explainable multiobjective optimization

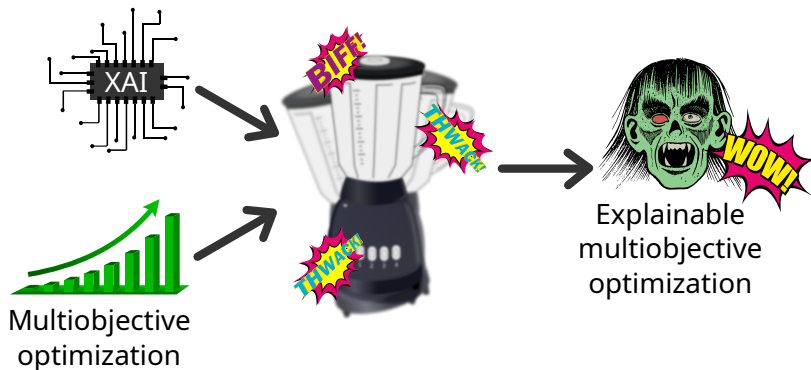


Figure: The main theme of this presentation and the main theme of my PhD.

Background

- 1 Motivation
- 2 Background
 - Multiobjective optimization
 - Scalarization
 - Reference point based interactive methods
 - Explainable Artificial Intelligence
 - Shapley values
- 3 Explainable interactive multiobjective optimization
- 4 Tests and results
- 5 What next?
- 6 Conclusions
- 7 Appendices

Multiobjective optimization problems

- A multiobjective optimization problem has many conflicting objectives, which are to be optimized simultaneously¹.

Multiobjective optimization problem

A multiobjective optimization problem can be defined as

$$\min F(\mathbf{x}) = (f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_k(\mathbf{x})), \quad (1)$$

where $f_1 \dots f_i$, $i \in [1, k]$ are objective functions and \mathbf{x} is a decision variable vector. The vectors \mathbf{x} can be subject to both **box-constraints** and **function constraints**. Feasible \mathbf{x} belong to the *feasible variable space* S or $\mathbf{x} \in S$.

¹Kaisa Miettinen. *Nonlinear multiobjective optimization*. Boston: Kluwer Academic Publishers, 1999.

Box-constraints

$$x_i^{\text{low}} \leq x_i \leq x_i^{\text{high}}, x_i \in \mathbf{x} \quad (2)$$

Function constraints

$$\begin{aligned} g(\mathbf{x}) - \delta_g &> 0 \\ h(\mathbf{x}) - \delta_h &= 0 \\ \delta_g, \delta_h &\in \mathbb{R} \end{aligned} \quad (3)$$

- In (2) x_i^{low} and x_i^{high} are the lower and higher limits for the i th element in \mathbf{x} , respectively.
- In (3) δ_g and δ_h are scalar values which should be exceeded or be exactly matched by $g(\mathbf{x})$ and $h(\mathbf{x})$, respectively.

More definitions

Objective vector

An objective vector \mathbf{z} is the image of the solution $\mathbf{x} \in S$ such that $F(\mathbf{x}) = \mathbf{z}$. The set of objective vectors Z consists of all the images \mathbf{z} .

Pareto optimality

A solution $\mathbf{x}^* \in S$ is said to be Pareto optimal if, and only if, there does not exist any other solution $\mathbf{x} \in S$ such that $f_i(\mathbf{x}) \leq f_i(\mathbf{x}^*) \forall i \in [1, k]$ and $f_i(\mathbf{x}) < f_i(\mathbf{x}^*)$ for at least some $i \in [1, k]$.

Pareto front

The Pareto front Z^{Pareto} consists of the images of all the Pareto optimal solutions. The set of Pareto optimal solutions is the Pareto optimal solution set.

More definitions

Ideal and nadir points

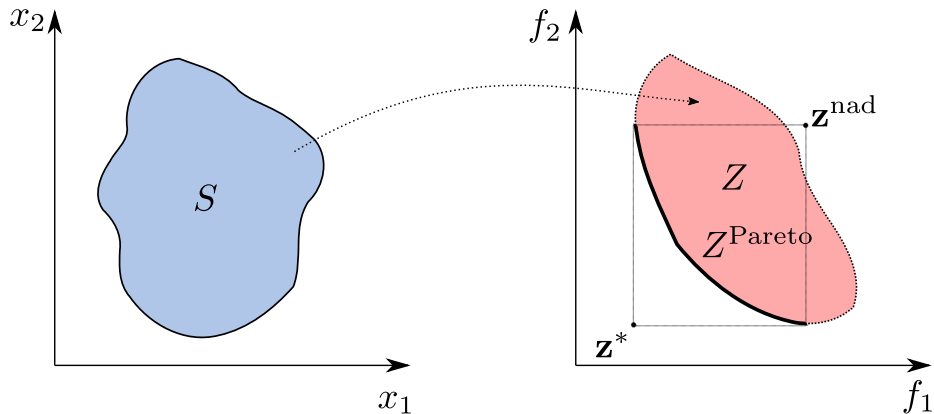
The ideal \mathbf{z}^* and nadir \mathbf{z}^{nad} points represent the best (lowest) and worst (highest) values of the objective function values on the Pareto front, respectively.

Reference point

A reference point $\bar{\mathbf{z}}$ is a vector of aspiration levels $\bar{z}_i, i = 1 \dots k$. The reference point can be provided by a decision maker, in which case, the reference point represents the decision maker's preferences.

Important concepts graphically

$$\min_{\mathbf{x} \in S} \{f_1(\mathbf{x}), f_2(\mathbf{x})\}$$



- Multiobjective optimization problems can be scalarized using a scalarizing function $s : \mathbb{R}^k \rightarrow \mathbb{R}$.

Scalarized problem

$$\begin{aligned} \min \quad & s(\mathbf{F}(\mathbf{x}); \mathbf{p}) \\ \text{subject to} \quad & \mathbf{x} \in S, \end{aligned} \tag{4}$$

where \mathbf{p} is a set of additional parameters given to the scalarizing function .

- Scalarizing functions usually have some desirable properties, such as guaranteeing (weak) Pareto optimality of the solution found.

- Scalarizing function used in STOM²:

STOM

$$\text{STOM}(\mathbf{F}; \bar{\mathbf{z}}, \mathbf{z}^{**}) = \min_{\mathbf{x} \in S} \max_{i=1, \dots, k} \left[\frac{f_i(\mathbf{x}) - z_i^{**}}{\bar{z}_i - z_i^{**}} \right] + \rho \sum_{i=1}^k \frac{f_i(\mathbf{x})}{\bar{z}_i - z_i^{**}}, \quad (5)$$

where $\mathbf{z}^{**} = (z_1^* - \delta, z_2^* - \delta, \dots, z_k^* - \delta)$ is an utopian point with $\delta \in \mathbb{R}^+$, and $\rho \in \mathbb{R}^+$.

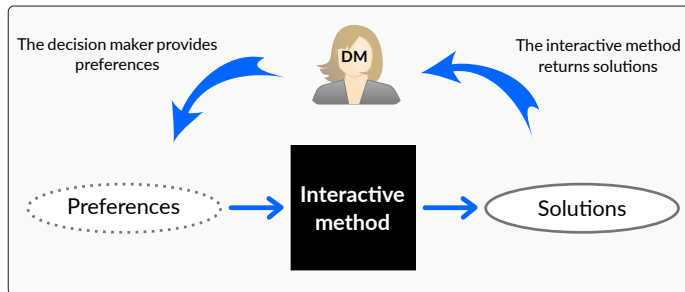
- A reference point $\bar{\mathbf{z}}$ can be incorporated in scalarizing functions.
- More examples of scalarizing functions in³.

²Hirota Nakayama. "Aspiration Level Approach to Interactive Multi-Objective Programming and Its Applications". In: *Advances in Multicriteria Analysis*. Ed. by Panos M. Pardalos, Yannis Siskos, and Constantin Zopounidis. Boston, MA: Springer, 1995, pp. 147–174. ISBN: 978-1-4757-2383-0. DOI: 10.1007/978-1-4757-2383-0_10.

³Kaisa Miettinen and Marko M. Mäkelä. "On scalarizing functions in multiobjective optimization". In: *OR Spectrum* 24.2 (2002), pp. 193–213. DOI: 10.1007/s00291-001-0092-9.

Interactive method

- A decision maker (DM) iteratively provides preference information as a reference point.
 - New solution(s) are computed for the problem after each iteration.
- We focus on reference point based interactive methods.



Explainable Artificial Intelligence

- The field of explainable artificial intelligence (XAI)⁴ focuses on the study and development of artificial intelligence that is capable of functioning in a way understandable by humans.
- Clear focus on machine learning methods, especially deep neural networks and deep learning in general.
- Usually model interpretability (by humans) and predictive power correlate negatively. I.e., most powerful machine learning models are also black-boxes.
- Roughly two main approaches to explainability: using interpretable models and model agnostic approaches.

⁴David Gunning et al. "XAI—Explainable artificial intelligence". In: *Science Robotics* 4.37 (2019). DOI: 10.1126/scirobotics.aay7120.

Why explainability?

- If we use artificial intelligence (AI) for decision-making, we cannot blindly trust any model.
- How to tell if a model works correctly? How to justify decisions made based on these models?
- Explainability aims to uncover these issues by shedding light on the black-box.
- There is societal pressure in EU as well to consider explainability in AI: GDPR recital 71⁵ (right to explanation⁶).

⁵<https://www.privacy-regulation.eu/en/r71.htm>

⁶Bryce Goodman and Seth Flaxman. "European Union regulations on algorithmic decision-making and a "right to explanation"". In: *AI magazine* 38.3 (2017), pp. 50–57.

Motivating examples

If the field of XAI is new to you and you work with AI/ML, I would highly suggest the following reads:

- Essay published in *Nature* by Cynthia Rudin on why we should stop explaining black-box models and use interpretable models instead.⁷
- Examples on how usage of AI can do more harm than good in society.⁸
- A very handy reference to start using interpretable AI.⁹
- A more traditional text book on XAI.¹⁰

⁷Cynthia Rudin. “Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead”. In: *Nature Machine Intelligence* 1.5 (2019), pp. 206–215.

⁸Hans de Bruijn, Martijn Warnier, and Marijn Janssen. “The perils and pitfalls of explainable AI: Strategies for explaining algorithmic decision-making”. In: *Government Information Quarterly* 39.2 (2022), p. 101666.

⁹Christoph Molnar. *Interpretable Machine Learning. A Guide for Making Black Box Models Explainable*. 2nd ed. 2022. URL: <https://christophm.github.io/interpretable-ml-book>.

¹⁰Uday Kamath and John Liu. *Explainable Artificial Intelligence: An Introduction to Interpretable Machine Learning*. Springer, 2021.

Shapley values and SHAP

- Shapley values¹¹ is a game-theoretical concept.
- Shapley values are a way to quantify the contribution of each player to the payoff in an n -player game.



¹¹Lloyd S Shapley. 17. A value for n -person games. Santa Monica: RAND Corporation, 1951.

Shapley values and SHAP

- Shapley values have been used in the field of XAI to explain black-box machine learning models.
- Because of the nature how Shapley values are computed (remove player from game, compute partial payoff) makes them hard to be used with arbitrary machine learning models.
- Instead, we may rely on SHAP values¹², particularly kernel SHAP, which are computationally less expensive to compute than pure Shapley values.

¹²Scott M Lundberg and Su-In Lee. "A Unified Approach to Interpreting Model Predictions". In: *Advances in Neural Information Processing Systems 30*. Ed. by I. Guyon et al. California: Curran Associates, Inc., 2017, pp. 4765–4774.

So?



"I am very confused Giovanni, what are you getting at?"

Explainable interactive multiobjective optimization

- 1 Motivation
- 2 Background
- 3 Explainable interactive multiobjective optimization**
- 4 Tests and results
- 5 What next?
- 6 Conclusions
- 7 Appendices

Could we somehow utilize SHAP values to probe interactive multiobjective optimization methods and get insight on how the preferences provided (the reference point) has affected the computed solution(s)?

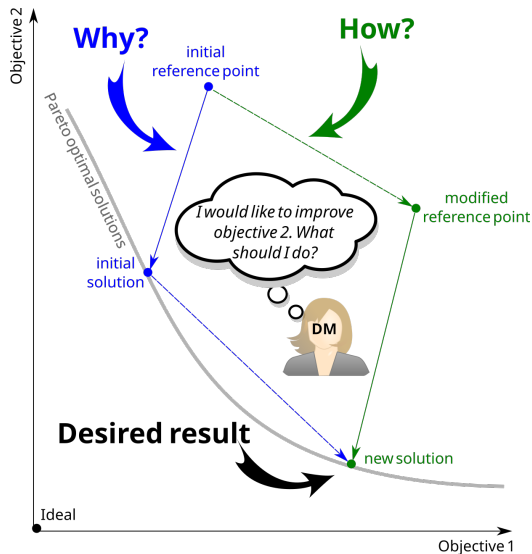
YES!

- A new method, R-XIMO, that can provide *explanations* on how reference points have affected computed solutions.
- From the explanations, we can derive *suggestions*.
- Suggestions to support the decision maker in providing new preferences in the next iteration.
- R-XIMO can be used with any multiobjective optimization method that takes as its input a reference point and computes a solution.

Explainable interactive multiobjective optimization

- A decision maker can express a wish to improve some objective function in a solution.
- A suggestion on how to modify the current reference point to achieve the desired improvement is provided.
- E.g., the decision maker wishes to improve objective 1, R-XIMO suggest that the decision maker should improve objective 1 in the reference point and impair objective 3.
 - We know this by computing SHAP values.

Explainable interactive multiobjective optimization



Explainable interactive multiobjective optimization

Decision maker: I would like to improve the first objective.

Example explanation:

Objective 1 was most improved in the solution by the second component and most impaired by the third component in the reference point.

Example suggestion:

Try improving the first¹³ component and impairing the third component in the reference point.

¹³We always improve the component that matches the objective the decision maker wishes to improve.

Tests and results

- 1 Motivation
- 2 Background
- 3 Explainable interactive multiobjective optimization
- 4 Tests and results**
- 5 What next?
- 6 Conclusions
- 7 Appendices

- We utilized three reference point based interactive methods that consisted of minimizing different scalarizing functions (5).
- We tested R-XIMO with two real-life multiobjective optimization problems (3 and 5 objectives).
 - Check what happens if we follow the suggestion provided by R-XIMO fully, partly, or not at all, and see if we were successful in improving the desired objective and how much it improved.
 - We did the above many times and got statistical data.

Results of numerical experiment

Some key finding based on the numerical tests:

- It is best to follow the suggestion provided by R-XIMO.
- Even only partly following the suggestion had some value.
- R-XIMO seems to work just as well for different scalarizing functions.


- Piloted the suggestions and explanations generated by R-XIMO with a human decision maker.
- Problem in Finnish forest management with three objectives.
- The decision maker was a domain expert in the field of forest management.
- The decision maker solved the problem twice.

- The suggestions were found to be useful by our human decision maker in the case study.
- The decision maker thought that R-XIMO supported them in reaching a satisfying solution in less iterations than without.
- However, the actual explanations were too complicated and the decision maker did not want to read them.

Autonomous Agents and Multi-Agent Systems (2022) 36:43
<https://doi.org/10.1007/s10458-022-09577-3>



Towards explainable interactive multiobjective optimization: R-XIMO

Giovanni Misitano¹  · Bekir Afsar¹ · Giomara Lárraga¹ · Kaisa Miettinen¹

Accepted: 6 July 2022 / Published online: 13 August 2022
© The Author(s) 2022

Abstract

In interactive multiobjective optimization methods, the preferences of a decision maker are incorporated in a solution process to find solutions of interest for problems with multiple conflicting objectives. Since multiple solutions exist for these problems with various trade-offs, preferences are crucial to identify the best solution(s). However, it is not necessarily clear to the decision maker how the preferences lead to particular solutions and,

¹⁴Giovanni Misitano et al. “Towards explainable interactive multiobjective optimization: R-XIMO”. In: *Autonomous Agents and Multi-Agent Systems* 36.2 (2022), pp. 1–43.

What next?

- 1 Motivation
- 2 Background
- 3 Explainable interactive multiobjective optimization
- 4 Tests and results
- 5 What next?**
- 6 Conclusions
- 7 Appendices

What next?

- With our work, we have taken an important step towards a new paradigm in (interactive) multiobjective optimization: **Explainable Interactive Multiobjective optimization** or **XIMO**.
- Took inspiration from what has been done in the field of explainable artificial intelligence
- Novel approaches tailored especially to multiobjective optimization probably needed.

Explainability in evolutionary multiobjective optimization

- Next paper brewing: exploring explainability in an evolutionary multiobjective optimization setting.
- Evolutionary multiobjective optimization is based on heuristics. Lots of opacity and unknowns about how solutions are generated. Maybe we could *explain* some of it?
- Tons of solutions are generated during an evolutionary multiobjective optimization method. This means there is a lot of data available by default! Maybe we could *learn* something from it?
- $1 + 1 = 3 \dots$

- Explainability in multiobjective optimization is still very much in its incubation stage, but cracks in the shell have appeared outside our current work as well, especially in the context of evolutionary multiobjective optimization¹⁵¹⁶¹⁷¹⁸¹⁹.

¹⁵Jinkun Wang et al. "Diversified recommendation incorporating item content information based on MOEA/D". In: *2016 49th Hawaii international conference on system sciences (HICSS)*. IEEE. 2016, pp. 688–696. DOI: [10.1109/HICSS.2016.91](https://doi.org/10.1109/HICSS.2016.91).

¹⁶Roykrong Sukkerd, Reid Simmons, and David Garlan. "Toward explainable multi-objective probabilistic planning". In: *2018 IEEE/ACM 4th International Workshop on Software Engineering for Smart Cyber-Physical Systems (SEsCPS)*. IEEE. 2018, pp. 19–25. DOI: [10.1145/3196478.3196488](https://doi.org/10.1145/3196478.3196488).

¹⁷Huixin Zhan and Yongcan Cao. "Relationship Explainable Multi-objective Optimization Via Vector Value Function Based Reinforcement Learning". In: *arXiv preprint arXiv:1910.01919* (2019).

¹⁸Giovanni Misitano. "Interactively Learning the Preferences of a Decision Maker in Multi-objective Optimization Utilizing Belief-rules". In: *2020 IEEE Symposium Series on Computational Intelligence (SSCI)*. IEEE. 2020, pp. 133–140. DOI: [10.1109/SSCI47803.2020.9308316](https://doi.org/10.1109/SSCI47803.2020.9308316).

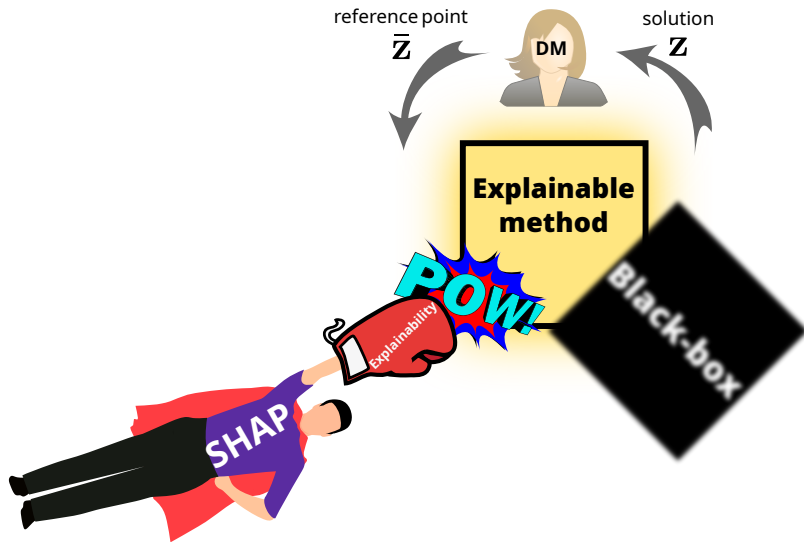
¹⁹Salvatore Corrente et al. "Explainable Interactive Evolutionary Multiobjective Optimization". In: *Available at SSRN 3792994* (2021).

Conclusions

- 1 Motivation
- 2 Background
- 3 Explainable interactive multiobjective optimization
- 4 Tests and results
- 5 What next?
- 6 Conclusions**
- 7 Appendices

- Explainability is an exciting and important concept to be studied in the context of and applied to multiobjective optimization.
- Makes the life of decision makers easier.
- Very much an unexplored area in the field of multiobjective optimization.
- New and wild ideas are needed!

Conclusions



Appendices

- 1 Motivation
- 2 Background
- 3 Explainable interactive multiobjective optimization
- 4 Tests and results
- 5 What next?
- 6 Conclusions
- 7 Appendices**

Acknowledgements

This work has been supported by the Academy of Finland (grant numbers 311877 and 322221) and the Vilho, Yrjö and Kalle Väisälä Foundation of the Finnish Academy of Science and Letters. This work is a part of the thematic research area Decision Analytics Utilizing Causal Models and Multiobjective Optimization (DEMO, jyu.fi/demo) at the University of Jyväskylä.

- DESDEO framework¹
- DESDEO website²
- Multiobjective Optimization (research) Group³
- Follow me on LinkedIn⁴
- If Twitter does not burn down, then you can find me there as well: @misitano_g

¹<https://desdeo.it.jyu.fi>

²G. Misitano et al. "DESDEO: The Modular and Open Source Framework for Interactive Multiobjective Optimization". In: *IEEE Access* 9 (2021), pp. 148277–148295. DOI: 10.1109/ACCESS.2021.3123825

³<http://www.mit.jyu.fi/optgroup/>

⁴<https://linkedin.com/in/misitano>

- [1] Kaisa Miettinen. *Nonlinear multiobjective optimization*. Boston: Kluwer Academic Publishers, 1999.
- [2] Hirotaka Nakayama. “Aspiration Level Approach to Interactive Multi-Objective Programming and Its Applications”. In: *Advances in Multicriteria Analysis*. Ed. by Panos M. Pardalos, Yannis Siskos, and Constantin Zopounidis. Boston, MA: Springer, 1995, pp. 147–174. ISBN: 978-1-4757-2383-0. DOI: 10.1007/978-1-4757-2383-0_10.
- [3] Kaisa Miettinen and Marko M. Mäkelä. “On scalarizing functions in multiobjective optimization”. In: *OR Spectrum* 24.2 (2002), pp. 193–213. DOI: 10.1007/s00291-001-0092-9.
- [4] David Gunning et al. “XAI–Explainable artificial intelligence”. In: *Science Robotics* 4.37 (2019). DOI: 10.1126/scirobotics.aay7120.

- [5] Bryce Goodman and Seth Flaxman. “European Union regulations on algorithmic decision-making and a “right to explanation””. In: *AI magazine* 38.3 (2017), pp. 50–57.
- [6] Cynthia Rudin. “Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead”. In: *Nature Machine Intelligence* 1.5 (2019), pp. 206–215.
- [7] Hans de Bruijn, Martijn Warnier, and Marijn Janssen. “The perils and pitfalls of explainable AI: Strategies for explaining algorithmic decision-making”. In: *Government Information Quarterly* 39.2 (2022), p. 101666.
- [8] Christoph Molnar. *Interpretable Machine Learning. A Guide for Making Black Box Models Explainable*. 2nd ed. 2022. URL: <https://christophm.github.io/interpretable-ml-book>.
- [9] Uday Kamath and John Liu. *Explainable Artificial Intelligence: An Introduction to Interpretable Machine Learning*. Springer, 2021.

- [10] Lloyd S Shapley. 17. *A value for n -person games*. Santa Monica: RAND Corporation, 1951.
- [11] Scott M Lundberg and Su-In Lee. “A Unified Approach to Interpreting Model Predictions”. In: *Advances in Neural Information Processing Systems 30*. Ed. by I. Guyon et al. California: Curran Associates, Inc., 2017, pp. 4765–4774.
- [12] Giovanni Misitano et al. “Towards explainable interactive multiobjective optimization: R-XIMO”. In: *Autonomous Agents and Multi-Agent Systems 36.2* (2022), pp. 1–43.
- [13] Jinkun Wang et al. “Diversified recommendation incorporating item content information based on MOEA/D”. In: *2016 49th Hawaii international conference on system sciences (HICSS)*. IEEE. 2016, pp. 688–696. DOI: 10.1109/HICSS.2016.91.

- [14] Roykrong Sukkerd, Reid Simmons, and David Garlan. “Toward explainable multi-objective probabilistic planning”. In: *2018 IEEE/ACM 4th International Workshop on Software Engineering for Smart Cyber-Physical Systems (SEsCPS)*. IEEE. 2018, pp. 19–25. DOI: 10.1145/3196478.3196488.
- [15] Huixin Zhan and Yongcan Cao. “Relationship Explainable Multi-objective Optimization Via Vector Value Function Based Reinforcement Learning”. In: *arXiv preprint arXiv:1910.01919* (2019).
- [16] Giovanni Misitano. “Interactively Learning the Preferences of a Decision Maker in Multi-objective Optimization Utilizing Belief-rules”. In: *2020 IEEE Symposium Series on Computational Intelligence (SSCI)*. IEEE. 2020, pp. 133–140. DOI: 10.1109/SSCI47803.2020.9308316.
- [17] Salvatore Corrente et al. “Explainable Interactive Evolutionary Multiobjective Optimization”. In: *Available at SSRN 3792994* (2021).

- [18] G. Misitano et al. “DESDEO: The Modular and Open Source Framework for Interactive Multiobjective Optimization”. In: *IEEE Access* 9 (2021), pp. 148277–148295. DOI: [10.1109/ACCESS.2021.3123825](https://doi.org/10.1109/ACCESS.2021.3123825).