



Robust Ordinal Regression for Subsets Comparisons with Interactions

Hugo Gilbert, Mohamed Ouaguenouni, Meltem Ozturk, Olivier Spanjaard

► To cite this version:

Hugo Gilbert, Mohamed Ouaguenouni, Meltem Ozturk, Olivier Spanjaard. Robust Ordinal Regression for Subsets Comparisons with Interactions. 2023. hal-04177872

HAL Id: hal-04177872

<https://hal.science/hal-04177872v1>

Preprint submitted on 6 Aug 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Robust Ordinal Regression for Subsets Comparisons with Interactions[★]

Hugo Gilbert^b, Mohamed Ouaguenouni^{a,*}, Meltem Öztürk^{b,**}, Olivier Spanjaard^a

^a*Sorbonne Université, CNRS, LIP6, Paris, F-75005, France*

^b*Université Paris Dauphine, PSL Research University, CNRS, LAMSADE, Paris, F-75016, France*

Abstract

This paper is dedicated to a robust ordinal method for learning the preferences of a decision maker between subsets. The decision model, derived from Fishburn and LaValle [20] and whose parameters we learn, is general enough to be compatible with any strict weak order on subsets, thanks to the consideration of possible interactions between elements. Moreover, we accept not to predict some preferences if the available preference data are not compatible with a reliable prediction. A predicted preference is considered reliable if all the simplest models (Occam’s razor) explaining the preference data agree on it. Following the robust ordinal regression methodology, our predictions are based on an uncertainty set encompassing the possible values of the model parameters. We define a robust ordinal dominance relation between subsets and we design a procedure to determine whether this dominance relation holds. Numerical tests are provided on synthetic and real-world data to evaluate the richness and reliability of the preference predictions made.

Keywords: robust ordinal regression, preference elicitation, positive and negative interactions, subsets comparisons

[★]This paper is a revised and extended version of a workshop paper at MPREF 2022, and an extended abstract at AAMAS 2023:

H. Gilbert, M. Ouaguenouni, M. Öztürk, O. Spanjaard, *Cautious Learning of Multiattribute Preferences*, 13th Workshop MPREF, Jul 2022, Vienna, Austria.

H. Gilbert, M. Ouaguenouni, M. Öztürk, O. Spanjaard, *Robust Ordinal Regression for Collaborative Preference Learning with Opinion Synergies*, AAMAS 2023, pp. 2439-2441.

^{*}Corresponding author

^{**}A significant part of the work presented here has been carried out while Meltem Öztürk was on delegation at LIP6.

1. Introduction

Preference elicitation (or preference learning) is an important step in setting up a recommender system for decision making. In this preference elicitation setting, our focus is on determining the parameters of a decision model that accurately captures the pairwise preferences of a Decision Maker (DM) over subsets, by comparing subsets of elements. The preferences are depicted using a highly adaptable model whose versatility stems from its ability to incorporate positive or negative synergies between elements [24]. Moreover, we provide an ordinally robust approach, in the sense that the preferences we infer do not rely on arbitrarily specified parameter values, but on the set of all parameter values that are compatible with the observed preferences. Importantly, another distinctive feature of our approach is its ability to learn the parameter set itself (not only the *values* of parameters).

The preference model we consider can be used in different contexts, depending on the nature of the subsets we are comparing. The subsets are represented by binary vectors, showing the presence or absence of an element in the subset. The elements of a subset can be for example:

- individuals (in the comparison of coalitions, teams, etc.),
- binary attributes (in the comparison of multiattribute alternatives),
- objects (in the comparison of subsets in a subset choice problem), etc.

For illustration, a toy example of such an elicitation context could be a coffee shop trying to determine its customers' favorite frozen yogurt flavor combination by offering them to test a small number of flavor combinations rather than having them taste each combination.

Objective of the paper. Our objective is to design a preference elicitation procedure that complies with the two following principles.

First, the sophistication of the learned preference model should be able to fit any level of complexity of the stated preferences. For this purpose, we use a utility function f general enough to represent any order \succ of preference, i.e., for any strict weak ordering \succ on a set \mathcal{A} of alternatives (i.e., subsets) there exists f such that, for any pair $\{A, B\} \subseteq \mathcal{A}$, $f(A) > f(B)$ iff $A \succ B$. Note that we also aim to make the model as simple as possible, in the sense that the parameter set remains as concise as possible (*sparse* model).

Second, the predicted pairwise preferences should not depend on the partly arbitrary choice of precise numerical values for the parameters of the model but solely on the stated preferences. Hence, we design an *ordinally robust* elicitation procedure that maintains an isomorphism between the collected preferential data and the learned model (in the same spirit as ordinal

measurement for problem solving [4]) by using a polyhedron of possible values for the parameters, reflecting the uncertainty about them. As a consequence, when predicting an unknown pairwise preference between two alternatives A and B , apart from the predictions “ A is preferred to B ” and “ B is preferred to A ”, it is possible that the model does not make a prediction due to a lack of sufficiently rich preferential data (the absence of prediction is preferred to a wrong prediction, although a compromise must obviously be made between the reliability of the prediction and the predictive power of the learned model).

Elicitation setting. The input of our elicitation procedure is a learning set consisting of pairwise comparisons of various alternatives. More precisely, we consider an offline elicitation setting (passive learning) where we assume that a dataset of comparison examples is available, from which the parameters of the preference model are (partially) specified. This is a separate framework from the online elicitation setting (active learning) where we would incrementally select pairwise preference queries to enrich the learning set. The output of the elicitation procedure consists of pairwise comparisons that were not present in the learning set, which we call (preference) *predictions* hereafter. Note that, in some cases, the model may choose not to provide a prediction. The elicitation procedure thus results in a strict partial order on the alternatives.

Organization of the paper. After an overview of the related work (Section 2), we present the θ -additive utility model (Section 3), as well as the robust ordinal dominance relation inferred from it, based on the knowledge of a collection of preference examples. We then show how to determine whether a subset dominates another subset given the known pairwise preferences of the DM (Section 4), which enables to make preference predictions. The paper ends with numerical tests on synthetic and real-world preference data, and comparison with other preference learning methods (Section 5).

2. Related work

Preference elicitation (see e.g. Dias et al. [16]) and preference learning (see e.g. Fürnkranz and Hüllermeier [21], Corrente et al. [15]) have been studied for a long time in operations research and artificial intelligence. This is a prerequisite in many applications across a wide range of fields, such as recommender systems, banking, financial management, chemistry, energy resources, health, investments, and industrial location [2]. Several issues can be tackled in preference elicitation, among which:

1. to handle a set of alternatives of combinatorial nature: an incremental preference elicitation is then often adopted, where comparison examples are interactively generated with the DM, in order to determine a necessary “optimal” alternative [e.g., 5, 10, 37];
2. to cope with preferences that cannot be represented by an additive utility function: for instance, the elicitation of generalized utility functions has been considered in the literature [11], but also the elicitation of several other involved decision models [6, 34];
3. to deal with “incorrect” preference examples: Bayesian approaches have been considered in this matter [9, 26], but also possibilistic approaches [1].

We focus here on the second challenge, by studying the elicitation of a set function taking into account positive and negative interactions between elements. Furthermore, preference elicitation problems differ in their purpose: some aim to produce a recommendation, others a set of recommendations, and still others pairwise comparisons. We will, in our case, produce a set of pairwise comparisons.

The Choquet integral is the most studied decision model for taking into account positive and negative interactions between criteria in multicriteria decision making [23]. It turns out that a Choquet integral defined on binary vectors representing subsets can be viewed as a set function. Note that a Choquet integral is parameterized by a capacity v on the criteria set N , i.e., a set function on N that is *monotone* ($A \subseteq B \Rightarrow v(A) \leq v(B)$) and *normalized* ($v(N) = 1$). As will become clear in the remainder of the paper, we do not impose such constraints in the model we consider. There are some recent works dealing with the elicitation of the parameters of a Choquet-related aggregation function: Bresson et al. [12] use a perceptron approach to learn the parameters of a 2-additive hierarchical Choquet integral, while Herin et al. [28] propose an algorithm to learn sparse Möbius representations from preference examples, without a prior k -additivity assumption. For a broad literature review about learning the parameters of a Choquet integral, the reader may refer to the article by Grabisch, Kojadinovic, and Meyer [24]. Let us also mention the work by Marichal and Roubens [32], which use a polyhedron to characterize the set of parameters that are compatible with a training set of examples. The idea of defining a polyhedron of uncertainty on the parameters of a utility function goes back at least to the work of Charnetski and Soland [13]. Their model state that $A \succ B$ if the proportion of parameters that give a better value for A than for B among those that are compatible with the stated preferences is greater than the proportion of

parameters that give a better value for B than for A . This principle was also adapted to the case of a Choquet integral by Angilella, Corrente and Greco [3]. In the sequel, we will use a similar polyhedron.

More precisely, we elicit a partial specification of a set function, namely the components of the parameter set and the set of parameter values, which yields an ordinal dominance relation between subsets. As already mentioned, we do not assume interactions with the DM but only the knowledge of a “static” training set of examples of pairwise preferences in order to predict pairwise comparisons between alternatives.

Predicting a comparison between alternatives can be framed as a binary classification problem by considering, as a training set, a set of triples (A, B, c) , where A and B are two alternatives and $c = 1$ if $A \succ B$, and $c = 0$ otherwise. In this setting, many approaches have been proposed, going from perceptrons [18] to Gaussian processes [14] or Support Vector Machines (SVM) [17].

An important feature of our elicitation procedure is that it may lead to not making predictions for some pairwise comparisons if the available preferential information is not conclusive enough. Other classification models also have such a possibility to not predict a class for some examples, either because of an ambiguity in the class to predict (ambiguity rejection) or because the example is too far from the examples that are in the learning set (novelty rejection). This type of approaches are generally used in safety-sensitive domains, e.g. to predict a disease in medical applications [29]. For a complete review of learning with reject option, we refer the reader to the survey made by Hendrickx et al. [27].

The two closest works to ours are those by Domshlak and Joachims [17] and by Bigot et al. [7]. Similarly to our approach, Domshlak and Joachims consider a function that could represent any weak order on the alternatives. More precisely, they consider a multiattribute utility function that is a sum of 4^n subutilities over subsets of attribute values, where n is the number of attributes. The subutility values are then learned using an efficient SVM approach based on the *kernel trick* [see e.g., 35]. Bigot et al. study the use of generalised additively independent decompositions of utility functions [19, 22]. They give a PAC-learner that is polynomial time if a constant bound is known on the degree of the function, where the *degree* is the size of the greatest subset of attributes in the decomposition. Yet, both works do not fit the robust ordinal learning framework we consider in this article.

3. From the θ -additive model to robust ordinal dominance

Given a set $\mathcal{F} = \{a_1, a_2, \dots, a_n\}$ of elements, we aim to reason on the preferences of the DM on a set \mathcal{A} of subsets $A \subseteq \mathcal{F}$, representing alternatives. The characteristic vector \vec{A} of a subset A is the n -dimensional binary vector whose i^{th} component is 1 if $a_i \in A$, and 0 otherwise. For instance, the characteristic vector of $A = \{a_1, a_2, a_4\}$ is $\vec{A} = (1, 1, 0, 1)$ if $\mathcal{F} = \{a_1, a_2, a_3, a_4\}$. In the following, we may use one or the other notation for describing a subset. Here are some examples of alternatives represented by subsets:

- If \mathcal{F} is a set of reference users expressing opinions on cultural products (e.g., movies), a cultural product may be represented by the subset A of reference users in \mathcal{F} that have a positive opinion on it, i.e., $a_i \in A$ if reference user a_i has a positive opinion on it, otherwise $a_i \notin A$.
- If \mathcal{F} is the set of players in a squad, a team lineup may be represented by the subset A of players that compound it.
- If \mathcal{F} is a set of binary features of technological products (e.g., smart-phones), a technological product may be represented by a subset A of features, i.e., $a_i \in A$ if the product has feature a_i , otherwise $a_i \notin A$.

We assume for simplicity that there are no two distinct alternatives corresponding to the same subset $A \subseteq \mathcal{F}$, which implies in particular that $2^{|\mathcal{F}|} \geq |\mathcal{A}|$. We infer strict pairwise preferences from strict preferences given by a DM on some subset of alternatives in \mathcal{A} , and we use this training set of pairwise preferences on alternatives (each viewed as a subset) to elicit the parameters of a utility function f defined on \mathcal{A} . The role of the utility function f is to represent the (unknown) strict weak order on \mathcal{A} corresponding to the DM's preferences, with $A \succ B$ iff $f(A) > f(B)$ and $A \sim B$ iff $f(A) = f(B)$.

We do not perform a full elicitation of the parameters of f , but we consider an uncertainty set of parameters values consistent with the known preferences of the DM, as in robust ordinal regression. If $f(A) > f(B)$ for all parameters values in this uncertainty set, then A is predicted to be strictly preferred to B . Actually, we do not only learn the parameters values, but also the components of the parameter set themselves, as we explain below.

3.1. The θ -additive model

Before coming to the proposed θ -additive model, we first recall the standard additive utility model, and its extension, the k -additive utility model.

The additive and k -additive utility models. As the DM's preferences over \mathcal{A} are modeled as a strict weak order, there exists a real-valued function

f such that $\forall A, B \in \mathcal{A}, f(A) > f(B) \Leftrightarrow A \succ B$. Many models assume that f can be represented in a compact way using some sort of additivity property. The simplest and most used one is the additive model [19]. This model makes the strong assumption that we can find a parameter value $v(a) \in \mathbb{R}$ for each element $a \in \mathcal{F}$ such that for all $A \in \mathcal{A}$, the utility of A is $f(A) = \sum_{a \in A} v(a)$. This assumption is strong because it implies that there is no interaction between the elements. A weaker assumption is that of k -additivity where we suppose the existence of a parameter $v(S) \in \mathbb{R}$ for each $S \in [\mathcal{F}]^k$, where $[\mathcal{F}]^k = \{S \subseteq \mathcal{F} : 1 \leq |S| \leq k\}$. Hence, in the k -additive model, for all $A \in \mathcal{A}$, $f(A) = \sum_{S \in [\mathcal{F}]^k} I_A(S) v_S$, where $I_A(S) = 1$ if $S \subseteq A$ and 0 otherwise, and v_S is an abbreviation for $v(S)$. Obviously, the 1-additive model amounts to the additive model. Taking k strictly greater than 1 makes it possible to account for (positive or negative) synergies between subsets of k or less elements. For example, the 2-additive model makes it possible to account for binary synergies. The utility of the alternative $A = (1, 1, 0, 1)$ with the 2-additive model is $f(A) = v(\{a_1\}) + v(\{a_2\}) + v(\{a_4\}) + v(\{a_1, a_2\}) + v(\{a_1, a_4\}) + v(\{a_2, a_4\})$. If there is a positive synergy between a_1 and a_2 then $f(\{a_1, a_2\}) > v(\{a_1\}) + v(\{a_2\})$ holds because $f(\{a_1, a_2\}) = v(\{a_1\}) + v(\{a_2\}) + v(\{a_1, a_2\})$. Note incidentally that $f(\{a_1, a_2\}) \neq v(\{a_1, a_2\})$. The n -additive model is general enough to represent *any* strict weak order on \mathcal{A} because it can represent any real-valued set function $f : 2^{\mathcal{F}} \rightarrow \mathbb{R}$ [25], provided that $f(\emptyset) = 0$. However, it requires to specify $2^n - 1$ parameters. We therefore restrict our attention to additive models requiring fewer parameters.

The θ -additive model. Given a set $\theta \subseteq 2^{\mathcal{F}}$, and a set function $v : \theta \rightarrow \mathbb{R}$, we assume that f is of the form $f(A) = \sum_{S \in \theta} I_A(S) v_S$, where v_S stands again for $v(S)$. We call this the θ -additive model. For this model, we may also use the notation $f_{\theta, v}(A)$ instead of $f(A)$. The 1-additive (resp. k -additive) model is the special case in which $\theta = [\mathcal{F}]^1$ (resp. $\theta = [\mathcal{F}]^k$).

Example 1. Let $\mathcal{F} = \{a_1, a_2, a_3, a_4\}$ be a set of 4 elements, $\mathcal{A} = \{0, 1\}^4$ and the DM's preferences be the strict weak order \succsim given by :

$$\begin{array}{llll}
\{a_2, a_3, a_4\} \succ & \{a_1, a_3, a_4\} \succ & \{a_1, a_2, a_4\} \succ & \{a_3, a_4\} \\
\succ & \{a_2, a_4\} \succ & \{a_2, a_3\} \succ & \{a_1, a_4\} \succ & \{a_1, a_3\} \\
\succ & \{a_1, a_2\} \succ & \{a_4\} \succ & \{a_3\} \succ & \{a_2\} \\
\succ & \{a_1\} \succ & A = \{a_1, a_2, a_3, a_4\} \sim & \emptyset \succ & B = \{a_1, a_2, a_3\}.
\end{array}$$

These preferences can be explained by a clear negative synergy when a_1 , a_2 , and a_3 are chosen together (in A and B). Interestingly, instead of using

a complete 3-additive model, which would require the definition of 14 parameters, this strict weak order can be obtained by using a θ -additive model with $\theta = \{\{a_1\}, \{a_2\}, \{a_3\}, \{a_4\}, \{a_1, a_2, a_3\}\}$ and $v_{\{a_1\}} = 1$, $v_{\{a_2\}} = 2$, $v_{\{a_3\}} = 3$, $v_{\{a_4\}} = 4$, $v_{\{a_1, a_2, a_3\}} = -10$. This allows us to benefit from the expressiveness offered by 3-additivity while restricting the number of parameters.

3.2. The θ -ordinal dominance relation

In our elicitation setting, we assume that we have only access to a partial set R of strict pairwise preferences provided by the DM. This set may contain only a few comparisons. Our aim is to use these observed preferences to infer other strict pairwise preferences on the set of alternatives. We formalize R as a set of pairs $(A, B) \in \mathcal{A}^2$ such that $(A, B) \in R \Leftrightarrow A \succ B$.

Moreover, given θ , the set of value functions on θ that are compatible with the preferences observed in R is denoted by V_θ^R :

$$V_\theta^R = \{v : \theta \rightarrow \mathbb{R} \mid \forall (A, B) \in R, f_{\theta, v}(A) > f_{\theta, v}(B)\}.$$

Note that, for a given θ , this set V_θ^R can be either empty or composed of an infinity of possible value functions on θ . Notably, if this set is empty then the preferences of the user cannot be represented by a θ -additive function. We denote by Θ^R the set $\{\theta \mid V_\theta^R \neq \emptyset\}$, i.e., the θ 's such that the preferences in R are consistent with a θ -additive function.

Unfortunately, given $\theta \in \Theta^R$ such that $V_\theta^R \neq \emptyset$, a pair $\{v, v'\}$ of value functions in V_θ^R may lead to infer opposite preferences, as illustrated below.

Example 2. Let $\mathcal{F} = \{a_1, a_2, a_3, a_4\}$. Let us assume that, contrary to Example 1, we now only observe preferences on the singletons $\{a_1\}, \{a_2\}, \{a_3\}, \{a_4\}$:

$\{a_4\} \succ \{a_3\} \succ \{a_2\} \succ \{a_1\}$, or equivalently:

$$R = \{(\{a_4\}, \{a_3\}), (\{a_4\}, \{a_2\}), (\{a_4\}, \{a_1\}), (\{a_3\}, \{a_2\}), (\{a_3\}, \{a_1\}), (\{a_2\}, \{a_1\})\}.$$

The two additive functions v and v' defined by:

$$v(\{a_1\}) = 1, v(\{a_2\}) = 2, v(\{a_3\}) = 3, v(\{a_4\}) = 5$$

$$\text{and } v'(\{a_1\}) = 1, v'(\{a_2\}) = 3, v'(\{a_3\}) = 4, v'(\{a_4\}) = 5$$

are both in V_θ^R , but we infer $\{a_1, a_4\} \succ \{a_2, a_3\}$ from v while we infer $\{a_2, a_3\} \succ \{a_1, a_4\}$ from v' .

This example shows that, given R , choosing a specific function $v \in V_\theta^R$ can lead to infer preferences that are only related to this arbitrary choice [4]. Our aim is to infer preferences for pairs outside R in a reliable way by eliminating such arbitrary choices. In this purpose, we turn to a robust ordinal regression approach based on the observed preferences in R .

Fishburn and Lavalley [20] showed how one can obtain an *ordinal dominance relation* from a partially specified 2-additive numerical model. We now explain how their idea can be extended to a θ -additive model.

Definition 1. Let \mathcal{F} be a set of elements, $\mathcal{A} \subseteq 2^{\mathcal{F}}$ a set of subsets and R a set of pairs $(A, B) \in \mathcal{A}^2$ where $(A, B) \in R \Leftrightarrow A \succ B$. Given $\theta \in \Theta^R$, the θ -ordinal dominance relation, denoted by \succ_{θ}^R , is defined for $A, B \in \mathcal{A}$ by:

$$A \succ_{\theta}^R B \Leftrightarrow \forall v \in V_{\theta}^R, f_{\theta,v}(A) > f_{\theta,v}(B).$$

The θ -ordinal dominance relation is independent from the choice of a specific $v \in V_{\theta}^R$. Naturally, $(A, B) \in R \Rightarrow A \succ_{\theta}^R B$. Nevertheless, note that the binary relation \succ_{θ}^R is obviously partial, and we define the incomparability relation \sim_{θ}^R as:

$$A \sim_{\theta}^R B \Leftrightarrow \exists v, v' \in V_{\theta}^R, f_{\theta,v}(A) \geq f_{\theta,v}(B) \text{ and } f_{\theta,v'}(B) \geq f_{\theta,v'}(A).$$

If $A \succ_{\theta}^R B$ then we can predict, based on R and for a θ -additive model, that A is strictly preferred to B . If $A \sim_{\theta}^R B$ then no prediction is made

We conclude this section by mentioning some properties of \succ_{θ}^R :

- Unlike \succ , the relation \succ_{θ}^R is not a strict weak order: it is asymmetric but it may not be complete nor negatively-transitive. The absence of preference prediction may occur in two situations that are not equivalent: either A and B belong to the same incomparability class of the (unknown) strict weak order \succ on \mathcal{A} , i.e., $A \sim B$, or there is not enough preferential information in R to conclude that $A \succ B$ or $B \succ A$.
- Since the ordinal dominance relation depends on the preference set R and on the model θ , the relation \succ_{θ}^R evolves when θ or R are restricted or extended. In particular, if $\theta' \subseteq \theta$ then any prediction that is yielded using ordinal dominance with the model θ is also yielded using ordinal dominance with the model θ' ; thus, if $V_{\theta}^R \neq \emptyset$ and $V_{\theta'}^R \neq \emptyset$, then θ' appears as more appealing from a preference learning standpoint since it allows more predictions to be made. Furthermore, one could prefer θ' over θ because of the philosophical principle of parsimony [e.g. 8].

A more formal and detailed description of the properties of \succ_{θ}^R can be found in the supplementary material (Appendix A).

3.3. The robust ordinal dominance relation

Note that the ordinal dominance relation is dependent on the choice of a specific set $\theta \in \Theta^R$. However, as shown in the following example, there may be several θ 's in Θ^R .

Example 3. Assume that R consists of all pairwise preferences resulting from \succ in Example 1. Setting $\theta = \{\{a_1\}, \{a_2\}, \{a_3\}, \{a_4\}\}$ yields then $V_{\theta}^R = \emptyset$.

In contrast, setting $\theta_1 = \{\{a_1\}, \{a_2\}, \{a_3\}, \{a_4\}, \{a_1, a_2, a_3\}\}$ yields $V_{\theta_1}^R \neq \emptyset$. Actually, there are many other sets θ compatible with the preferences in R : it can be shown¹ that $\Theta^R = \{\theta : \theta_1 \subseteq \theta\}$ for this example.

The question that naturally arises is whether we could find two different models $\theta_1, \theta_2 \in \Theta^R$ that are both compatible with the observed preferences in R and such that $A \succ_{\theta_1}^R B$ and $B \succ_{\theta_2}^R A$ for a pair of alternatives $(A, B) \in \mathcal{A}^2$. Unfortunately, this situation may indeed happen:

Example 4. Let $R = \{(\{a_1\}, \{a_2\})\}$, $\theta_1 = \{\{a_1\}\}$ and $\theta_2 = \{\{a_2\}\}$. Note that both θ_1 and θ_2 belong to Θ^R . If we consider $\theta_1 = \{\{a_1\}\}$, the set $V_{\theta_1}^R$ is compounded of value functions v defined on θ_1 such that $v(\{a_1\}) > 0$. Hence, for all $v \in V_{\theta_1}^R$ we have $f_{\theta_1, v}(\{a_1, a_2\}) = v(\{a_1\}) > 0 = f_{\theta_1, v}(\emptyset)$ and thus $\{a_1, a_2\} \succ_{\theta_1}^R \emptyset$. Conversely, if we consider $\theta_2 = \{\{a_2\}\}$, the set $V_{\theta_2}^R$ is compounded of value functions v defined on θ_2 such that $v(\{a_2\}) < 0$. This yields $f_{\theta_2, v}(\{a_1, a_2\}) = v(\{a_2\}) < 0$ for each $v \in V_{\theta_2}^R$ and thus $\emptyset \succ_{\theta_2}^R \{a_1, a_2\}$.

In what follows, we will define a more robust variant of the ordinal dominance relation. This variant will take into account the plurality of models compatible with the observed preferences.

Note that there always exists a θ able to represent R (at worst, $\theta = 2^{\mathcal{F}}$) and that if a θ -additive model is compatible with R , then any θ' -additive model with $\theta \subseteq \theta'$ is also compatible with R . For this reason, the number of sets θ compatible with the observed preferences may be very large.

For this reason, we start by restricting the set of models to take into account. In this purpose, we need a binary relation \sqsubseteq on Θ^R , such that $\theta \sqsubseteq \theta'$ if θ is considered simpler than θ' . Our idea is to only consider sets θ that are minimal according to such a binary relation, i.e., θ such that $\nexists \theta' \in \Theta^R$ for which $\theta' \sqsubseteq \theta$. This is motivated by the philosophical principle of parsimony that the simpler of two explanations is to be preferred (Occam's razor [8]). Different possible definitions for \sqsubseteq will be discussed upon in the following subsection.

We call \sqsubseteq -simplest θ of Θ^R the parameter sets $\theta \in \Theta^R$ which are minimal w.r.t. \sqsubseteq , and we denote by Θ_{\sqsubseteq}^R their set. Based on Θ_{\sqsubseteq}^R , we extend the ordinal dominance relation to define the \sqsubseteq -robust ordinal dominance relation.

Definition 2. Let \mathcal{F} be a set of elements, $\mathcal{A} \subseteq 2^{\mathcal{F}}$ a set of subsets and R a set of pairs $(A, B) \in \mathcal{A}^2$ where $(A, B) \in R \Leftrightarrow A \succ B$. The \sqsubseteq -robust ordinal

¹It has been computer tested by brute force enumeration.

dominance relation, denoted by \succ_{\sqsubseteq}^R , is defined, for $A, B \in \mathcal{A}$, as follows:

$$\begin{aligned} A \succ_{\sqsubseteq}^R B &\iff \forall \theta \in \Theta_{\sqsubseteq}^R, A \succ_{\theta}^R B, \\ &\iff \forall \theta \in \Theta_{\sqsubseteq}^R, \forall v \in V_{\theta}^R, f_{\theta,v}(A) > f_{\theta,v}(B). \end{aligned}$$

In other words, A \sqsubseteq -robustly ordinally dominates B if A θ -ordinally dominates B according to all θ in Θ_{\sqsubseteq}^R , i.e., all the \sqsubseteq -simplest θ 's of Θ^R .

3.4. Different definitions for \sqsubseteq

We say that a relation \sqsubseteq is *based on* a function ξ when $\theta \sqsubseteq \theta'$ if and only if $\xi(\theta) \leq \xi(\theta')$. Several aspects can be taken into account to define ξ :

- A first idea is to favor parameter sets θ that minimize the complexity of synergies between the attributes. To measure this complexity, we use the *degree* of θ , namely $\deg(\theta) = \max\{|S| : S \in \theta\}$ (i.e., the greatest cardinality of a subset of interacting attributes). This leads to the binary relation \sqsubseteq_{\deg} based on \deg , i.e., $\theta_1 \sqsubseteq_{\deg} \theta_2 \iff \deg(\theta_1) \leq \deg(\theta_2)$.
- A second idea is to favor parameter sets θ having the *sparsest* possible representation [38], i.e., those which minimize $\text{card}(\theta) = |\theta|$. This choice yields the binary relation $\sqsubseteq_{\text{card}}$, which is the relation based on the function card , i.e., $\theta_1 \sqsubseteq_{\text{card}} \theta_2 \iff \text{card}(\theta_1) \leq \text{card}(\theta_2)$.
- Alternatively, we define a binary relation combining the ideas of \sqsubseteq_{\deg} and $\sqsubseteq_{\text{card}}$ by considering both the number and the size of elements in a parameter set θ . In this purpose, we define \sqsubseteq_{ws} , the relation based on the function $\text{ws}(\theta) = \sum_{S \in \theta} |S|$, i.e., $\theta_1 \sqsubseteq_{\text{ws}} \theta_2 \iff \text{ws}(\theta_1) \leq \text{ws}(\theta_2)$.
- Lastly, we define the binary relation \sqsubseteq_{lex} , defined by using lexicographically the binary relations \sqsubseteq_{\deg} , $\sqsubseteq_{\text{card}}$, and \sqsubseteq_{ws} , in this order. This relation could be seen as based on the function lex where $\text{lex}(\theta) = n4^n \deg(\theta) + n2^n \text{card}(\theta) + \text{ws}(\theta)$.

Example 5. Let $R = \{(\{a_1, a_2\}, \{a_3, a_4\}), (\{a_1, a_2\}, \{a_1, a_3\})\}$. It is easy to see that $V_{\theta}^R \neq \emptyset$ for $\theta = \{\{a_1, a_2\}\}$, which corresponds to a model of degree 2. However, we may prefer being consistent with a model of degree 1, even if there are more elements in it: $\theta' = \{\{a_1\}, \{a_2\}\}$ or $\theta'' = \{\{a_1\}, \{a_3\}\}$ or $\theta''' = \{\{a_2\}\}$. In this example, the minimal parameter set θ among $\theta', \theta'', \theta'''$ w.r.t. relation \sqsubseteq_{\deg} (resp. $\sqsubseteq_{\text{card}}$, \sqsubseteq_{ws} , \sqsubseteq_{lex}) is $\{\theta', \theta'', \theta'''\}$ (resp. $\{\theta'''\}$ in the three cases).

4. Preference prediction by using robust ordinal dominance

Given a set R of pairwise preferences and a binary relation \sqsubseteq on Θ^R , the preference learning method we propose consists in predicting that a subset

A is preferred to B if $A \succ_{\sqsubseteq}^R B$, i.e., A is preferred to B for all simplest models $\theta \in \Theta^R$ and value functions $v \in V_\theta^R$. The purpose of this section is to detail the procedure for determining whether $A \succ_{\sqsubseteq}^R B$. It is organized as follows:

- We show that determining if $A \succ_\theta^R B$ is polytime in $|R|$ and $|\theta|$, while determining if $A \succ_{\sqsubseteq}^R B$ amounts to testing whether $\Theta_{\sqsubseteq}^R \cap \Theta_{B \succsim A}^R = \emptyset$, where $\Theta_{B \succsim A}^R = \{\theta \in \Theta^R : B \succ_\theta^R A \text{ or } B \sim_\theta^R A\}$ (Subsection 4.1).
- As determining an explicit representation of Θ_{\sqsubseteq}^R is likely to be cumbersome (as the size of Θ_{\sqsubseteq}^R can be very large), we turn to an implicit representation based on the values $\mathbf{deg}(\theta)$, $\mathbf{card}(\theta)$, $\mathbf{ws}(\theta)$ for $\theta \in \Theta_{\sqsubseteq}^R$. We thus study the computational complexity of determining $\mathbf{deg}(\theta)$ (resp. $\mathbf{card}(\theta)$, $\mathbf{ws}(\theta)$, $\mathbf{lex}(\theta)$) for $\theta \in \Theta_{\sqsubseteq}^R$ and $\sqsubseteq = \sqsubseteq_{\mathbf{deg}}$ (resp. $\sqsubseteq = \sqsubseteq_{\mathbf{card}}$, $\sqsubseteq = \sqsubseteq_{\mathbf{ws}}$, $\sqsubseteq = \sqsubseteq_{\mathbf{lex}}$), showing that the former problem can be solved in polynomial time, while the others are NP-hard (Subsection 4.2).
- The implicit representation of Θ_{\sqsubseteq}^R is based on the following idea: if we know that $\theta_0 \in \Theta_{\sqsubseteq}^R$, then $\theta \in \Theta_{\sqsubseteq}^R \Leftrightarrow (\mathbf{deg}(\theta), \mathbf{card}(\theta), \mathbf{ws}(\theta)) = (\mathbf{deg}(\theta_0), \mathbf{card}(\theta_0), \mathbf{ws}(\theta_0))$. It is thus enough to determine a single model $\theta_0 \in \Theta_{\sqsubseteq}^R$ to be able to determine whether a model belongs to Θ_{\sqsubseteq}^R . This is why we propose a Mixed Integer Program (MIP) to compute a model $\theta \in \Theta_{\sqsubseteq}^R$, derived from a linear program for determining whether a model θ belongs to Θ^R (Subsection 4.3).
- We derive from it another MIP to compute a model in $\Theta_{\sqsubseteq}^R \cap \Theta_{B \succsim A}^R$, concluding $A \not\succ_{\sqsubseteq}^R B$ if it exists, $A \succ_{\sqsubseteq}^R B$ otherwise (Subsection 4.4).

4.1. Determining whether $A \succ_\theta^R B$ and whether $A \succ_{\sqsubseteq}^R B$

We first show that, unsurprisingly, linear programming provides an operational tool for determining whether $A \succ_\theta^R B$. Viewing a value function on θ as a vector $v = (v_S)_{S \in \theta}$ where $v_S = v(S)$, the set V_θ^R corresponds to the polyhedron defined by the following linear constraints in the $|\theta|$ -dimensional parameter space (where each parameter v_S corresponds to a dimension)²:

$$\forall (X, Y) \in R, \sum_{S \in \theta} I_X(S) v_S - \sum_{S \in \theta} I_Y(S) v_S \geq 1.$$

For a given set R of strict pairwise preferences and a model $\theta \in \Theta^R$, checking whether $A \succ_\theta^R B$ can be evaluated in polynomial time in $|R|$ and in $|\theta|$ by solving the following linear program $\mathcal{P}_{A \succ_\theta^R B}$, where there is one

²The right hand side of the constraint is here set to 1, but it could be set to any strictly positive constant as utilities v_S are always compatible with R to within a positive multiplicative factor.

variable $v_S \in \mathbb{R}$ for each pair $S \in \theta$:

$$(\mathcal{P}_{A \succ_{\theta}^R B}) \left\{ \begin{array}{ll} \min \sum_{S \in \theta} I_A(S) v_S - \sum_{S \in \theta} I_B(S) v_S \\ \sum_{S \in \theta} (I_X(S) - I_Y(S)) v_S \geq 1 & \forall (X, Y) \in R \setminus \{(A, B)\}, \\ v_S \in \mathbb{R} & \forall S \in \theta. \end{array} \right.$$

We have that $A \succ_{\theta}^R B$ if and only if the optimal value of $\mathcal{P}_{A \succ_{\theta}^R B}$ is strictly positive, as it implies that $\sum_{S \in \theta} I_A(S) v_S > \sum_{S \in \theta} I_B(S) v_S$ for all $v \in V_{\theta}^R$.

In contrast with this positive complexity result for ordinal dominance, determining whether $A \succ_{\sqsubseteq}^R B$ by direct use of the definition of robust ordinal dominance would require a high computational burden. We overcome this difficulty by reducing this problem to testing whether $\Theta_{\sqsubseteq}^R \cap \Theta_{B \succsim A}^R$ is empty.

To achieve this reduction, let us study the relationships between $\Theta_{A \succsim B}^R$, $\Theta_{B \succsim A}^R$ and Θ_{\sqsubseteq}^R . For visual support, the reader may refer to Figure 1. We recall that we denote by $\Theta_{B \succsim A}^R$ the set $\{\theta \in \Theta^R : B \succ_{\theta}^R A \text{ or } B \sim_{\theta}^R A\}$. As one of the relations $A \succ_{\theta} B$ or $B \succ_{\theta} A$ or $A \sim_{\theta} B$ holds for any $\theta \in \Theta^R$, we have that $\Theta^R = \Theta_{A \succsim B}^R \cup \Theta_{B \succsim A}^R$. Consequently, $\Theta_{\sqsubseteq}^R \subseteq \Theta_{A \succsim B}^R \cup \Theta_{B \succsim A}^R$ because $\Theta_{\sqsubseteq}^R \subseteq \Theta^R$. Furthermore, $\Theta_{A \succsim B}^R \cap \Theta_{B \succsim A}^R = \{\theta \in \Theta^R : A \sim_{\theta} B\} \neq \emptyset$ as soon as there exists $\theta \in \Theta^R$ for which $A \sim_{\theta} B$.

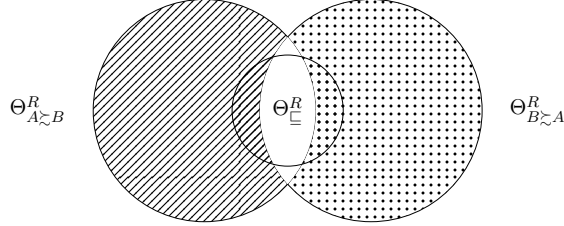


Figure 1: $(\Theta_{\sqsubseteq}^R \cap \Theta_{B \succsim A}^R = \emptyset \Leftrightarrow A \succ_{\sqsubseteq}^R B)$ and $(\Theta_{\sqsubseteq}^R \cap \Theta_{A \succsim B}^R = \emptyset \Leftrightarrow B \succ_{\sqsubseteq}^R A)$.

To evaluate whether a robust ordinal dominance relation holds between two subsets A and B , we examine if one of the following conditions holds:

- (i) $\Theta_{\sqsubseteq}^R \cap \Theta_{B \succsim A}^R = \emptyset$,
- (ii) $\Theta_{\sqsubseteq}^R \cap \Theta_{A \succsim B}^R = \emptyset$.

We have indeed the following result:

Proposition 1. *For any $A, B \subseteq \mathcal{F}$, we have $A \succ_{\sqsubseteq}^R B \Leftrightarrow \Theta_{\sqsubseteq}^R \cap \Theta_{B \succsim A}^R = \emptyset$.*

Proof. It follows from the following sequence of equivalences:

$$\begin{aligned} A \succ_{\sqsubseteq}^R B &\Leftrightarrow \forall \theta \in \Theta_{\sqsubseteq}^R, A \succ_{\theta}^R B \Leftrightarrow \forall \theta \in \Theta_{\sqsubseteq}^R, B \not\succ_{\theta}^R A \text{ and } A \not\prec_{\theta}^R B \\ &\Leftrightarrow \Theta_{\sqsubseteq}^R \cap \Theta_{B \succsim A}^R = \emptyset. \end{aligned} \quad \square$$

Symmetrically, we have obviously that $B \succ_{\Theta}^R A \Leftrightarrow \Theta_{\sqsubseteq}^R \cap \Theta_{A \succsim B}^R = \emptyset$. To test whether $\Theta_{\sqsubseteq}^R \cap \Theta_{B \succsim A}^R = \emptyset$, the mathematical programming approach we propose applies to cases where relation \sqsubseteq is based on a function ξ . The approach starts by computing a *single* model $\theta \in \Theta^R$ minimizing $\xi(\theta)$, which is enough for determining the value $\xi(\theta)$ of *any* $\theta \in \Theta_{\sqsubseteq}^R$, as they all share the same optimal value $\xi(\theta)$. We now study the complexity of computing such an optimal θ in Θ^R . More precisely, we study the complexity of the following decision problem MIN- θ - ξ , for $\xi \in \{\text{card}, \text{ws}, \text{deg}, \text{lex}\}$ (as is well-known, the optimization problem is at least as hard as its decision variant):

MIN- θ - ξ
INPUT: A set \mathcal{A} of alternatives, a set $R = \{(A, B), A, B \in \mathcal{A}\}$ of strict pairwise preferences, an integer $\tau \in \mathbb{Z}^+$.
QUESTION: Does there exist $\theta \in \Theta^R$ such that $\xi(\theta) \leq \tau$?

4.2. Computational complexity of MIN- θ - ξ for $\xi \in \{\text{card}, \text{ws}, \text{lex}, \text{deg}\}$

We show here that MIN- θ - ξ is NP-hard for $\sqsubseteq \in \{\text{ws}, \text{card}, \text{lex}\}$, while it can be solved in polynomial time for $\sqsubseteq = \text{deg}$.

Theorem 1. *MIN- θ -card and MIN- θ -ws are NP-complete.*

Proof. The membership of MIN- θ -card to NP follows from the fact that $\min_{\theta} \text{card}(\theta) \leq 2|R|$ and checking that $\theta \in \Theta^R$ can be done in polynomial time in $|R|$ and $|\theta|$. Indeed, the parameter set $\theta = \{A \in \mathcal{A} : (A, \cdot) \in R \text{ or } (\cdot, A) \in R\}$ obviously belongs to Θ^R , and $|\theta| \leq 2|R|$. The proof that MIN- θ -ws belongs to NP is similar, based on the fact that $\min_{\theta} \text{ws}(\theta) \leq 2|R| \times n$.

To prove the NP-hardness, we use a reduction from Hitting Set:

HITTING SET
INPUT: Given a set of n elements: $\mathcal{X} = \{x_i\}_{1 \leq i \leq n}$, a family of m sets $\mathcal{S} = \{S_i : S_i \subseteq \mathcal{X}, 1 \leq i \leq m\}$, and an integer $\tau \in \mathbb{Z}^+$.
QUESTION: Does there exist $\mathcal{X}' \subseteq \mathcal{X}$ such that $\forall S_i \in \mathcal{S}, S_i \cap \mathcal{X}' \neq \emptyset$ and $ \mathcal{X}' \leq \tau$?

Given an instance $(\mathcal{X}, \mathcal{S}, \tau)$ of the Hitting Set problem, we define the following instance (\mathcal{A}, R, τ') of MIN- θ -card (resp. MIN- θ -ws).

We let $\mathcal{A} = \mathcal{S} \cup \{\emptyset\}$, $\tau' = \tau$, and consider the following set of preferences:

$$R = \{(S, \emptyset) : S \in \mathcal{S}\}.$$

Now we show that $(\mathcal{X}, \mathcal{S}, \tau)$ is a yes-instance of Hitting Set iff (\mathcal{A}, R, τ') is a yes-instance of MIN- θ -**card** (resp. MIN- θ -**ws**). Note that a set θ belongs to Θ^R if and only if it satisfies the following condition:

$$\forall (S, \emptyset) \in R, \exists T \in \theta \text{ such that } T \subseteq S. \quad (\text{C})$$

Indeed, each preferences in R can then be satisfied by assigning positive values to parameters entailed by the elements of θ . Moreover, note that if a set θ satisfies C and $\exists T \in \theta$ such that $|T| > 1$, then the set θ' obtained from θ by replacing T by any singleton $\{x\} \subset T$ also satisfies C. Hence, within the sets satisfying C and minimizing **card**, there exists a set θ' compounded only of singletons, minimizing both **card** and **ws** (because $\text{card}(\theta) = \text{ws}(\theta)$ if θ is compounded only of singletons). By taking $\mathcal{X}' = \{x : \{x\} \in \theta'\}$, we obtain a hitting set of size $|\mathcal{X}'| \leq \tau$. This yields the following conclusion: there exists a hitting set of size $s \leq \tau$ if and only if there exists a set θ satisfying C such that $\text{card}(\theta) = s$ (resp. $\text{ws}(\theta) = s$). This argument completes the proof. \square

The following result is a direct consequence of the previous one:

Corollary 1. *MIN- θ -lex is NP-hard.*

Proof. Given an instance (\mathcal{A}, R, τ) of the MIN- θ -**card** problem, we could solve for each degree $d \in \{0, 1, \dots, |\mathcal{F}|\}$ an instance (\mathcal{A}, R, τ') of the MIN- θ -**lex** problem where $\tau' = dn4^n + (\tau + 1)n2^n$. \square

In contrast, we show a polynomial-time complexity result for MIN- θ -**deg**, by resorting to the *kernel trick*, widely used in machine learning [see e.g., 35]. Given a vector space \mathcal{X} of dimension $n_{\mathcal{X}}$ and a transformation function $\varphi : \mathcal{X} \rightarrow \mathcal{Y}$, where the dimension $n_{\mathcal{Y}}$ of vector space \mathcal{Y} is exponential in $n_{\mathcal{X}}$, the kernel trick consists in computing the scalar products $\langle \varphi(x), \varphi(y) \rangle$ of $x, y \in \mathcal{X}$ in polynomial time in $n_{\mathcal{X}}$, by using a kernel function $K(x, y)$ that returns the value $\langle \varphi(x), \varphi(y) \rangle$ without requiring to explicit $\varphi(x)$ and $\varphi(y)$. In our setting, \mathcal{X} is the set of characteristic vectors of subsets A of \mathcal{F} , and \mathcal{Y} the set of “augmented” characteristic vectors containing additional dimensions corresponding to binary values $I_A(S)$ for $S \in [\mathcal{F}]^\tau$ (more details in the proof). The complexity result is formulated as follows:

Theorem 2. *MIN- θ -deg can be solved in polynomial time in $|R|$ and n .*

Proof. Let (\mathcal{A}, R, τ) be an instance of MIN- θ -**deg**. We wish to determine if preferences in R can be represented by a θ -additive model with $\theta = [\mathcal{F}]^\tau$. For notational convenience, we set $\theta^{(\tau)} = [\mathcal{F}]^\tau$ and $n_\tau = |\theta^{(\tau)}| = \sum_{i=1}^\tau \binom{n}{i}$. We associate to $\theta^{(\tau)}$ the vector $\overrightarrow{\theta^{(\tau)}} = (S_1, \dots, S_{n_\tau})$, where subsets $S = \{a_{i_1}, \dots, a_{i_k}\}$ ($i_1 < \dots < i_k$) are indexed in lexicographic order of vectors $(|S|, i_1, \dots, i_k)$.

For instance, if $\mathcal{F} = \{a_1, a_2, a_3\}$ and $\theta = \theta^{(3)}$ then $\vec{\theta} = (\{a_1\}, \{a_2\}, \{a_3\}, \{a_1, a_2\}, \{a_1, a_3\}, \{a_2, a_3\}, \{a_1, a_2, a_3\})$. Additionally, for a value function $v : \theta \rightarrow \mathbb{R}$, we denote by $\vec{v} = (v_{S_1}, \dots, v_{S_{n_\tau}})$ the vector of values associated to the elements of $\vec{\theta}^{(\tau)}$ ordered in the same fashion. Finally, given $A \in \mathcal{A}$, we denote by \vec{A}_τ the binary vector $\vec{A}_\tau = (I_A(S_1), \dots, I_A(S_{n_\tau}))$ where $I_A(S_i)$ is the indicator function of $S_i \in \theta^{(\tau)}$.

Problem MIN- θ -deg evaluates if the following proposition holds:

$$\exists \vec{v} \in \mathbb{R}^{n_\tau} \text{ s.t. } \forall (A, B) \in R; \vec{A}_\tau \vec{v}^T > \vec{B}_\tau \vec{v}^T.$$

A value vector \vec{v} of minimum norm can be determined by solving the following convex quadratic program:

$$\begin{aligned} \min_{\vec{v} \in \mathbb{R}^{n_\tau}} \quad & \frac{1}{2} \vec{v} \vec{v}^T \\ \text{s.t.} \quad & \vec{A}_\tau \vec{v}^T \geq \vec{B}_\tau \vec{v}^T + 1 \quad \forall (A, B) \in R \end{aligned}$$

Using the same trick as Domshlak and Joachims [17], instead of solving this program whose number n_τ of variables is not polynomial in the size of our instance of MIN- θ -deg (because τ is an input variable and not a constant), we consider its Wolfe dual defined by:

$$\begin{aligned} \max_{\alpha \in \mathbb{R}^{|R|}} \quad & \sum_{(A,B) \in R} \alpha_{(A,B)} - \frac{1}{2} \sum_{(A,B) \in R} \sum_{(C,D) \in R} \alpha_{(A,B)} \alpha_{(C,D)} (\vec{A}_\tau - \vec{B}_\tau)(\vec{C}_\tau - \vec{D}_\tau)^T \\ \text{s.t.} \quad & \alpha \geq 0 \end{aligned}$$

By defining the kernel function $K^{(\tau)}(A, B) = \vec{A}_\tau \vec{B}_\tau^T$, the previous program can be written as:

$$\begin{aligned} \max_{\alpha \in (\mathbb{R}^+)^{|R|}} \quad & \sum_{(A,B) \in R} \alpha_{(A,B)} - \frac{1}{2} \sum_{(A,B) \in R} \sum_{(C,D) \in R} \alpha_{(A,B)} \alpha_{(C,D)} \\ & (K^{(\tau)}(A, C) - K^{(\tau)}(A, D) - K^{(\tau)}(B, C) + K^{(\tau)}(B, D)) \end{aligned}$$

which can be solved in polynomial time in $|R|$ and n provided that $K^{(\tau)}(X, Y)$ can be evaluated in polynomial time in n without expliciting X and Y .

Indeed, since the reformulation yields a convex quadratic program of polynomial size in the input data, the problem can then be solved in polynomial time (by polynomial time solvability of convex quadratic programming [30, 31]). We now prove that $K^{(\tau)}(X, Y)$ can be efficiently computed without expliciting X and Y . Let k be the size of the intersection between X and Y , i.e., $k = |X \cap Y|$. Note that $K^{(\tau)}(X, Y)$ counts the number of parameters of $\theta^{(\tau)}$ that are subsets of both X and Y . We conclude by noting that the number of such elements corresponds to $\sum_{i=1}^{\tau} \binom{k}{i}$, i.e., the number ($< 2^n$) of

non-empty subsets of size less than or equal to τ in $X \cap Y$. \square

Remark 1. Note that Tehrani et al. [36] and Herin et al. [28] have proposed kernel functions $K(x, y)$ that return the scalar product $\langle \varphi(x), \varphi(y) \rangle$ of augmented vectors $\varphi(x), \varphi(y)$ used to obtain an additive expression $\langle m, \varphi(x) \rangle$ of a discrete Choquet integral $C(x)$, where m is the vector of Möbius masses obtained from the capacity used in $C(x)$. It turns out that there is a close link between $f_{\theta, v}$ and a Choquet integral $C(x)$ expressed as $\langle m, \varphi(x) \rangle$ (note however that we do not impose the constraints on the $v(S)$ values ensuring the monotonicity of the capacity, or the normalization constraint $\sum_S v(S) = 1$). However, their kernel functions do not use the same calculations as ours: we take advantage of the particular case we study, where all components of x take binary values, to compute the kernel function in $O(n)$ instead of $O(n^2)$.

Algorithm 1 takes as input a set R of strict pairwise preferences and computes $\min\{\deg(\theta) : \theta \in \Theta^R\}$ by solving a sequence of convex quadratic programs establishing whether there exists $\theta \in \Theta^R$ such that $\xi(\theta) = \tau$ (which holds if the optimal value of the program is bounded). The variable τ is gradually incremented from 1. At each iteration, the objective function parameters are updated by using the kernel trick, which makes the procedure polynomial-time in $|R|$ and n .

Algorithm 1 Compute $\min\{\deg(\theta) : \theta \in \Theta^R\}$

Input: set R of strict pairwise preferences

Output: $\min\{\deg(\theta) : \theta \in \Theta^R\}$

$\tau \leftarrow 1$

for $(A, B) \in R$ **do**

for $(C, D) \in R$ **do**

\triangleright Initialization of dictionary Q

$Q[A, B, C, D] \leftarrow |A \cap C| - |A \cap D| - |B \cap C| + |B \cap D|$

while $\max_{\alpha \geq 0} \sum_{(A, B) \in R} \alpha_{(A, B)} - \frac{1}{2} \sum_{(A, B) \in R} \sum_{(C, D) \in R} \alpha_{(A, B)} \alpha_{(C, D)} Q[A, B, C, D]$ is unbounded **do**

\triangleright the $\alpha_{(X, Y)}$'s are the variables of the convex quadratic program

$\triangleright \alpha \geq 0$ means that $\alpha_{(X, Y)} \geq 0$ for all $(X, Y) \in R$

$\triangleright Q$ contains the coefficients of the objective function, updated at each iteration

$\tau \leftarrow \tau + 1$

for $(A, B) \in R$ **do**

for $(C, D) \in R$ **do**

$Q[A, B, C, D] \leftarrow Q[A, B, C, D] + \binom{|A \cap C|}{\tau} - \binom{|A \cap D|}{\tau} - \binom{|B \cap C|}{\tau} + \binom{|B \cap D|}{\tau}$

return τ

4.3. Computing $(\deg(\theta), \text{card}(\theta), \text{ws}(\theta))$ for $\theta \in \Theta_{\sqsubseteq_{\text{lex}}}^R$

As all models $\theta \in \Theta_{\sqsubseteq_{\text{lex}}}^R$ share the same vector $(\deg(\theta), \text{card}(\theta), \text{ws}(\theta))$, it is enough to compute a single model $\theta \in \Theta_{\sqsubseteq_{\text{lex}}}^R$ to deduce this vector, which

will be required to determine whether $A \succ_{\sqsubseteq_{\text{lex}}}^R B$. The negative complexity result (Corollary 1) regarding the computation of a model $\theta \in \Theta_{\sqsubseteq_{\text{lex}}}^R$ does not prevent us from proposing an exact solution method that will prove efficient in practice. For this purpose, we first present a Linear Program (LP) allowing us to determine in polynomial time in $|R|$ and $|\theta|$ whether $\theta \in \Theta^R$, given a model θ and a set R of strict pairwise preferences. From this LP, we will then develop a MIP formulation for computing $\theta \in \Theta_{\sqsubseteq_{\text{lex}}}^R$.

For a given set R of strict pairwise preferences and a given model θ , checking whether $\theta \in \Theta^R$ can be evaluated in polynomial time in $|R|$ and in $|\theta|$ by solving the following linear program \mathcal{P}_θ , where there is one variable $e_{(A,B)} \geq 0$ for each pair (A, B) in R :

$$(\mathcal{P}_\theta^R) \begin{cases} \min \sum_{(A,B) \in R} e_{(A,B)} \\ \sum_{S \in \theta} (I_A(S) - I_B(S))v_S + e_{(A,B)} \geq 1 & \forall (A, B) \in R, \\ e_{(A,B)} \geq 0 & \forall (A, B) \in R, \\ v_S \in \mathbb{R} & \forall S \in \theta. \end{cases}$$

We have that $\theta \in \Theta^R$ if and only if the optimal value of \mathcal{P}_θ^R is 0, because in this case we can find values for variables v_S that respect all the preferences in R without the help of the additional slack variables $e_{(A,B)}$.

We now show how to derive, from \mathcal{P}_θ^R , a MIP formulation for computing a model $\theta \in \Theta_{\sqsubseteq_{\text{lex}}}^R$. For this, we first compute $\deg(R) = \min\{\deg(\theta) : \theta \in \Theta^R\}$, by using Algorithm 1. We then add a binary variable b_S for each $S \in [\mathcal{F}]^{\deg(\theta)}$, as well as big-M constraints to ensure that $b_S = 1$ iff $S \in \theta$ (i.e., S belongs to the model $\theta \in \Theta_{\sqsubseteq_{\text{lex}}}^R$). Determining a model $\theta \in \Theta_{\sqsubseteq_{\text{lex}}}^R$ can be done by solving the following lexicographic optimization problem:

$$(\mathcal{P}_{\sqsubseteq_{\text{lex}}}^R) \begin{cases} \min \text{lex} \sum_{S \in [\mathcal{F}]^{\deg(R)}} b_S, \sum_{S \in [\mathcal{F}]^{\deg(R)}} b_S |S| \\ \sum_{S \in [\mathcal{F}]^{\deg(R)}} (I_A(S) - I_B(S))v_S \geq 1 & \forall (A, B) \in R, \\ -b_S M \leq v_S \leq b_S M & \forall S \in [\mathcal{F}]^{\deg(R)}, \\ b_S \in \{0, 1\} & \forall S \in [\mathcal{F}]^{\deg(R)}. \end{cases} \quad (1)$$

where $M = (2 \sum_{i=1}^{\deg(R)} \binom{n}{i} + |R|) \times (|R|)^{2|R|+2}$, so that if the values v_S can be set to satisfy constraints 1, then there exist such values in the interval $[-M, M]$ (see [33]). Every feasible instantiation of variables v_S, b_S in $\mathcal{P}_{\sqsubseteq_{\text{lex}}}^R$ corresponds to an element $\theta \in \Theta^R$, namely $\theta = \{S \in [\mathcal{F}]^{\deg(R)} : b_S = 1\}$. Lexicographic optimization amounts to determine, among feasible instantiations of v_S, b_S that minimize the first objective $\sum_{S \in [\mathcal{F}]^{\deg(R)}} b_S$, one that

minimizes the second objective $\sum_{S \in [\mathcal{F}]^{\deg(R)}} b_S |S|$. It is well-known that this can be achieved as follows, using a mixed integer programming solver:

- first, we solve the MIP \mathcal{P}_1 obtained by replacing the lexicographic objective function in $\mathcal{P}_{\sqsubseteq_{\text{lex}}}^R$ by $\min \sum_{S \in [\mathcal{F}]^{\deg(R)}} b_S$;
- denoting by opt_1 the optimal value of \mathcal{P}_1 , we then solve the MIP \mathcal{P}_2 where the objective function in $\mathcal{P}_{\sqsubseteq_{\text{lex}}}^R$ is replaced by $\min \sum_{S \in [\mathcal{F}]^{\deg(R)}} b_S |S|$, under the additional constraint $\sum_{S \in [\mathcal{F}]^{\deg(R)}} b_S \leq \text{opt}_1$.

As every feasible instantiation corresponds to a model θ of minimal degree $\deg(\theta)$ (i.e., $\deg(\theta) = \deg(R)$), we thus obtain a model $\theta \in \Theta_{\sqsubseteq_{\text{lex}}}^R$, from which we deduce $(\deg(\theta), \text{card}(\theta), \text{ws}(\theta))$ for $\theta \in \Theta_{\sqsubseteq_{\text{lex}}}^R$. In the following, we denote by $(\deg_{\text{lex}}, \text{card}_{\text{lex}}, \text{ws}_{\text{lex}})$ the vector $(\deg(\theta), \text{card}(\theta), \text{ws}(\theta))$ for $\theta \in \Theta_{\sqsubseteq_{\text{lex}}}^R$.

4.4. Determining whether $A \succ_{\sqsubseteq_{\text{lex}}}^R B$

Determining whether $A \succ_{\sqsubseteq_{\text{lex}}}^R B$ amounts to solve:

$$(\mathcal{P}_{A \succ_{\sqsubseteq_{\text{lex}}}^R B}) \left\{ \begin{array}{ll} \min \sum_{S \in [\mathcal{F}]^{\deg(R)}} b_S |S| & \\ \sum_{S \in [\mathcal{F}]^{\deg(R)}} b_S \leq \text{card}_{\text{lex}}, & (2) \\ \sum_{S \in [\mathcal{F}]^{\deg(R)}} (I_B(S) - I_A(S)) v_S \geq 0, & (3) \\ \sum_{S \in [\mathcal{F}]^{\deg(R)}} (I_X(S) - I_Y(S)) v_S \geq 1 \quad \forall (X, Y) \in R, & (4) \\ -b_S M \leq v_S \leq b_S M & \forall S \in [\mathcal{F}]^{\deg(R)}, \\ b_S \in \{0, 1\} & \forall S \in [\mathcal{F}]^{\deg(R)}. \end{array} \right.$$

A feasible solution of $\mathcal{P}_{A \succ_{\sqsubseteq_{\text{lex}}}^R B}$ yields a model θ satisfying $\deg(\theta) = \deg_{\text{lex}}$ (variables b_S are only defined for $S \in [\mathcal{F}]^{\deg(R)}$) and $\text{card}(\theta) = \text{card}_{\text{lex}}$ (by constraint 2 on the value of $\text{card}(\theta)$). Furthermore, constraint 3 ensures that $\theta \in \Theta_{B \succsim A}^R$, while constraint 4 ensures that $\theta \in \Theta^R$. If the optimal value of $\mathcal{P}_{A \succ_{\sqsubseteq_{\text{lex}}}^R B}$ is ws_{lex} , then the corresponding model θ belongs to $\theta \in \Theta_{\sqsubseteq_{\text{lex}}}^R$ (because then $(\deg(\theta), \text{card}(\theta), \text{ws}(\theta)) = (\deg_{\text{lex}}, \text{card}_{\text{lex}}, \text{ws}_{\text{lex}})$), and thus there exists $\theta \in \Theta_{\sqsubseteq_{\text{lex}}}^R \cap \Theta_{B \succsim A}^R$. Consequently:

- if the optimal value is strictly greater than ws_{lex} , or the polyhedron is empty, then $\Theta_{\sqsubseteq_{\text{lex}}}^R \cap \Theta_{B \succsim A}^R = \emptyset$ and hence $A \not\succ_{\sqsubseteq_{\text{lex}}}^R B$ (by Proposition 1);
- if the optimal value of $\mathcal{P}_{A \succ_{\sqsubseteq_{\text{lex}}}^R B}$ is ws_{lex} , then $A \not\succ_{\sqsubseteq_{\text{lex}}}^R B$.

5. Numerical tests

We call hereafter ORD the learning approach consisting in computing $(\text{deg}(R), \text{card}(R), \text{ws}(R))$ and using $\succ_{\sqsubseteq_{1\text{ex}}}^R$ for preference prediction. Numerical tests were carried out on Google Colab³, with the aim of comparing ORD with state of the art approaches in two different settings:

- A first set of experiments were carried out on *synthetic data*, i.e., obtained by simulating a user. They aimed at evaluating our approach in an ideal setting where a θ -additive model perfectly fits the preferences.
- A second set of experiments were carried out on real-world data for content-based filtering methods (more precisely, movies described by binary attributes). These tests aimed at evaluating how our approach deals with partially described alternatives (i.e., with possible “collisions” if two distinct alternatives share the same description), compared to other state of the art approaches.

In both sets of experiments, we start with a learning set of preferences. Based on this learning set, pairwise preference predictions are then requested on random pairs of alternatives (pairs not in the learning set). As said earlier, the model may not make a prediction if it is not robust enough given the available preference data (i.e., if there is no robust ordinal dominance).

5.1. The synthetic and real-world datasets

The dataset consists of ratings assigned by a user (DM) on a set \mathcal{A} of N alternatives. Given a set $\mathcal{F} = \{a_1, \dots, a_n\}$ of binary features, a learning set $\mathcal{A}_{\text{train}}$ consists of $k \leq N$ ratings of alternatives in \mathcal{A} , where each alternative A_i ($i = 1, \dots, N$) is described by a binary vector $\vec{A}_i = (A_i^1, \dots, A_i^n)$, with $A_i^j = 1$ if $a_j \in A_i$, and $A_i^j = 0$ otherwise. The user rating of A_i is denoted by r_i . The set of known strict preferences is $R = \{(A_i, A_j) \in \mathcal{A}_{\text{train}}^2 : r_i > r_j\}$.

The *real-world data* consist of ratings of movies by users picked up from the IMDb dataset⁴. This is a dataset of movie reviews that contains over 50k reviews. Each movie A_i is described by a set of binary features A_i^j , and the ratings r_i are integer values ranging from 1 to 10. The experiments were conducted with a dataset of 50 users (randomly sampled) who each rated at least $k=100$ movies. Each movie is described using a subset of $n=8$ binary features (corresponding to the main genres of the movie, e.g., “adventure”, “animation”, “children”, “comedy”, “fantasy”, etc.).

³two virtual CPU at 2.2GHz, 13GB RAM.

⁴www.kaggle.com/datasets/gauravduttakiit/imdb-recommendation-engine.

The *synthetic data* are generated in two steps: first a θ -additive function $f_{\theta,v}$ is randomly sampled, then a rating function is inferred from $f_{\theta,v}$. The procedure is precisely detailed in the following two paragraphs.

Sampling a θ -additive function $f_{\theta,v}$. For sampling a function $f_{\theta,v}$, we first sample a set θ and then we sample parameters v_S for $S \in \theta$. More precisely, the generation of θ is achieved as follows. First, θ is initialised as the set of singletons $\{a_1\}, \{a_2\}, \dots, \{a_n\}$, then we add $\lfloor \alpha \times (2^n - n) \rfloor$ subsets of attributes, where the coefficient $\alpha \in [0, 1]$ makes it possible to control the model's complexity: for $\alpha = 0$, only the singletons are in θ , which yields the simple additive utility model, and for $\alpha = 1$, all subsets of attributes are present, which yields the most general utility model. Each subset S is sampled according to a parameter $p \in (0, 1]$:

1. Initialize S as a singleton by uniformly sampling in \mathcal{F} .
2. Uniformly sample another attribute in \mathcal{F} and add it to S .
3. Exit this process if $S = \mathcal{F}$.
4. Exit this process with a probability p otherwise go to 2.

The expected size of sets S we sample is $\mathbb{E}[|S|] = 2 + (1 - p - (1 - p)^{n-1})/p$. Once θ is set, we sample the parameters v_S for each $S \in \theta$ with a normal distribution $\mathcal{N}(0, \sigma)$. The sampling of $f_{\theta,v}$ thus depends on three parameters p , α and σ . In the tests, p varies in $[0.1, 0.9]$, α in $[0.1, 0.5]$, and $\sigma = 100$.

From $f_{\theta,v}$ to a rating function. A function $r : \mathcal{A} \rightarrow \{1, \dots, t\}$ simulates the ratings of the user (of which only a subset of examples $r(A_i) = r_i$, for $i \in \{1, \dots, k\}$, is known to the model). The definition of r from $f_{\theta,v}$ depends on a parameter t defining the domain $\{1, \dots, t\}$ of possible ratings. The range of scores $f_{\theta,v}(A) = \sum_{S \in \theta} v_S I_A(S)$ of alternatives A is partitioned into t equally-sized intervals $(v_{k-1}, v_k]$ between the min score $v_0 = \min_{A \in \mathcal{A}} f_{\theta,v}(A)$ and the max score $v_t = \max_{A \in \mathcal{A}} f_{\theta,v}(A)$. The function r is then:

$$r(A) = \min\{1 \leq k \leq t : f_{\theta,v}(A) \leq v_k\}.$$

Put another way, the rating of A corresponds to the index k of the interval $(v_{k-1}, v_k]$ in which $f_{\theta,v}(A)$ lies. In general, the wider the domain of possible ratings, the fewer incomparabilities (alternatives with the same rating).

5.2. Baseline models

We briefly describe here the baseline models to which ORD is compared. Throughout the subsection, we have $\theta = [\mathcal{F}]^{\deg(R)}$ and each alternative A is described by an *augmented* binary vector $\vec{A} = (I_A(S_1), \dots, I_A(S_{|\theta|}))$, where $S_1, \dots, S_{|\theta|}$ are the subsets of \mathcal{F} of size less than or equal to $\deg(R)$.

Linear Regression (LR). We consider the θ -additive model, and we use linear regression to determine the value function \hat{v} such that $f_{\theta, \hat{v}}$ best approximates the utility function f , by minimizing $\sum_{i=1}^k (\vec{A}_i \vec{v}^T - \text{normalized}(r_i))^2$, where $\text{normalized}(r_i) = \frac{r_i - \min_i r_i}{\max_i r_i - \min_i r_i}$ (note that $\vec{A}_i \vec{v}^T = f_{\theta, v}(A_i)$). Put another way, we use the least squares method⁵ with ratings normalized in $[0, 1]$. We predict $A \succ B$ if $f_{\theta, \hat{v}}(\vec{A}) > f_{\theta, \hat{v}}(\vec{B})$.

Support Vector Machine (SVM). This baseline model is inspired by the approach proposed by Domshlak and Joachims [17]. An SVM approach is a supervised learning method for binary classification: each example in the dataset is labeled by 0 or 1; an SVM is learned from the dataset⁶, from which labels are inferred for new examples. In our setting, each preference $A \succ B$ in R yields two examples: a $(|\theta|+1)$ -dimensional vector $(\vec{A} - \vec{B}, 1)$ and another vector $(\vec{B} - \vec{A}, 0)$. That is, the third component of $(\vec{A} - \vec{B}, c)$ is $c = 1$ if A is preferred to B , and $c = 0$ if it is not. For predicting the preference between two alternatives A and B , we infer the labels of $(\vec{A} - \vec{B})$ and $(\vec{B} - \vec{A})$ by using the SVM. If the label of $(\vec{A} - \vec{B})$ is 1 (resp. 0) and that of $(\vec{B} - \vec{A})$ is 0 (resp. 1), then we predict $A \succ B$ (resp. $B \succ A$).

K-Nearest Neighbours (KNN). The distance-based models are widely used in the context of recommender systems. The distance-based model we consider is implemented as follows. The predicted rating of an alternative A is obtained by making a weighted sum $\sum_{i=1}^K w_i r_i$ of the ratings r_1, \dots, r_K of its K nearest neighbours A_1, \dots, A_K in the learning set⁷, with each weight w_i proportional to the Euclidean distance of the neighbour \vec{A}_i to \vec{A} . The value of K was set to $K = 5$ in our experiments, after preliminary tests showing this was the value yielding the best results for the dataset considered here. For predicting the preference between two alternatives A and B , we compute the predicted ratings of them, and predict the preference accordingly.

5.3. Experimental setup

In all experiments, the dataset is a set \mathcal{A} of N alternatives, described by a set \mathcal{F} of n binary features, and an associated rating vector r (integer values). The rating $r(A)$ of each alternative $A \in \mathcal{A}$ is known. To compare the performances of the different learning methods, we extract a subset \mathcal{A}_{train} of

⁵Precisely the `LinearRegression` function from the `scikit-learn` python library.

⁶We use the `SVC` function from the `scikit-learn` python library.

⁷We use the `KNeighborsClassifier` function from the `scikit-learn` python library.

k alternatives from \mathcal{A} , on which the models are trained. The alternatives in \mathcal{A}_{train} are chosen uniformly at random. We then randomly sample 100 pairs $\{A, B\}$ in \mathcal{A} such that $A \notin \mathcal{A}_{train}$ or $B \notin \mathcal{A}_{train}$ (possibly neither A nor B belongs to \mathcal{A}_{train}), and we compare the predicted pairwise preference with the actual preference: $A \succ B$ if $r(A) > r(B)$, $B \succ A$ if $r(B) > r(A)$, $A \sim B$ (incomparability) if $r(A) = r(B)$. The extraction of a subset \mathcal{A}_{train} from \mathcal{A} , the training of each model and the (100) pairwise preference predictions are performed 10 times, and the prediction performances are averaged over the 10 runs. We detail below the parameters that are used for the experiments on synthetic data and for the experiments on real-world data.

Synthetic data. The experiments on synthetic data were conducted with $|\mathcal{F}| = 8$ binary features, which yields a set \mathcal{A} of $2^{|\mathcal{F}|} = 256$ alternatives, a scale of $t = 12$ possible ratings, and the set of parameters $(\alpha, p, \sigma) = (0.1, 0.9, 100)$ for the generation of $f_{\theta, v}$. This set of parameters yields functions $f_{\theta, v}$ that are usually up to 4-additive, with an average $|\theta|$ equal to 12. This setting is not really restrictive as, given the number of strict pairwise preferences in R that are considered in our experiments (i.e., $|R| \leq \binom{|\mathcal{A}_{train}|}{2}$), it is unlikely that R cannot be represented by using a function $f_{\theta, v}$ of degree up to 4. The size of \mathcal{A}_{train} indeed varies between 12 and 29, from which between $|R| = \binom{12}{2} = 66$ and $\binom{29}{2} = 400$ pairwise preferences can be inferred.

Real-world data. For each of the 50 users that have rated at least 100 movies, a dataset \mathcal{A} including between 45 and 100 alternatives is first extracted. A training set \mathcal{A}_{train} is then extracted from \mathcal{A} , with $|\mathcal{A}_{train}|$ corresponding to 90% of $|\mathcal{A}|$ (which is common practice in machine learning, in particular for performing 10-fold cross-validation). The size of \mathcal{A}_{train} thus varies from 5 to 10, from which between $|R| = \binom{5}{2} = 10$ and $\binom{10}{2} = 45$ pairwise preferences can be inferred.

5.4. Evaluation metrics

We outline here the specific metrics that will be used to evaluate the ORD approach and compare it to other methods. To define our metrics we consider the 9 cases that can occur in the confusion matrix defined below.

Confusion Matrix. For a given pair of alternatives $(A, B) \in \mathcal{A}^2$ each model could either infer (predicted output) that A is *better* than B ($A \succ B$), or that A is *worse* than B ($B \succ A$) or it could return that the relation between A and B is *unknown*. Then, as outlined earlier, by comparing $r(A)$ and $r(B)$, we can have (real outputs) that A is indeed better than B if $r(A) > r(B)$ or that A is worse than B if $r(B) > r(A)$ or that the relation between them is

unknown if they share the same rating (incomparability). Our metrics are based on the confusion matrix defined in Table 1, where the rows symbolizes the predicted outputs and the columns the real outputs.

Predicted/Real	(B)etter	(W)orst	(U)nknown
(B)etter	BB	BW	BU
(W)orst	WB	WW	WU
(U)nknown	UB	UW	UU

Table 1: Confusion matrix.

Precision. The precision is defined as the ratio between the number of correct predictions among all the predictions that *were* made.

$$P = \frac{BB + WW}{BB + WW + BW + WB + BU + WU}.$$

Recall. The recall is defined as the ratio between the number of correct predictions among all the predictions that *could be* made.

$$R = \frac{BB + WW}{BB + WW + BW + WB + UB + UW}.$$

The precision metric penalizes the models making unreliable predictions, while the recall metric penalizes the models that avoid making predictions.

F1-score. F1-score is a metric that combines precision and recall to provide a balanced evaluation of a model’s performance. It is obtained by computing the harmonic mean of precision and recall:

$$F = 2 \frac{P \times R}{P + R}.$$

As the F1-score captures both precision and recall, it is an ideal metric for evaluating the robustness and accuracy of the studied models. Hence, we strongly rely on it when presenting our results.

Prediction Correctness. This metric is similar to precision, except that it does not take into account predictions that cannot be evaluated for lack of preferential information to check whether they are correct or incorrect.

$$PC = \frac{BB + WW}{BB + WW + BW + WB}.$$

Prediction Rate. This metric does not take into account the correctness of the predictions, it simply evaluates the rate at which the model produces predictions:

$$PR = 1 - \frac{UB + UW + UU}{M},$$

where M represents all the cases of Table 1 ($BB + WW + BW + WB + BU + WU + UB + UW + UU$).

5.4.1. Results on synthetic data

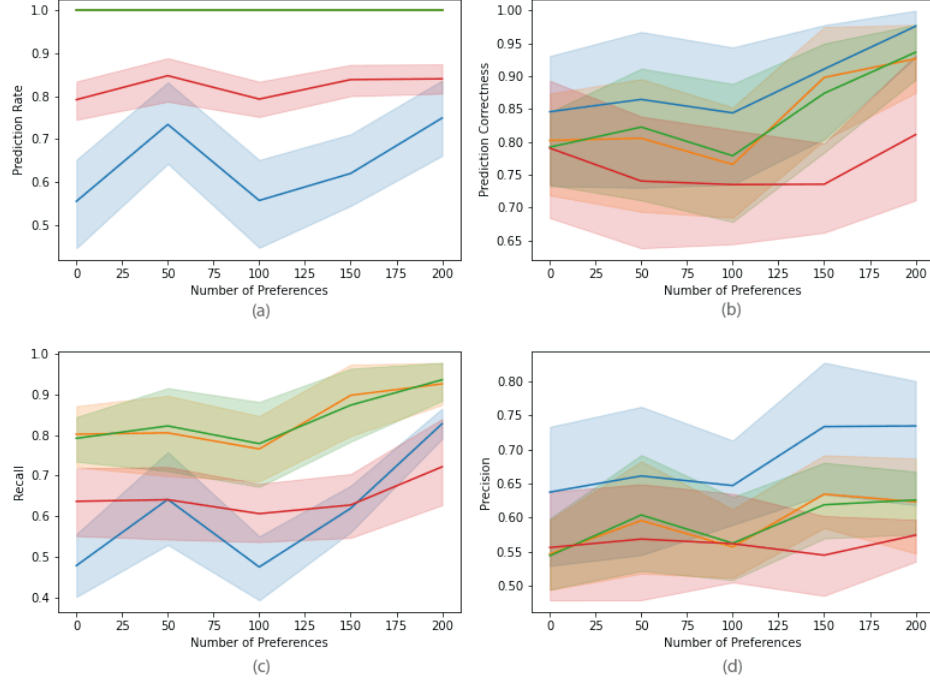


Figure 2: Precision, Recall, Prediction Rate and Prediction Correctness according to the number $|R|$ of preferences for models ORD (blue), KNN (red), SVM (orange), LR (green).

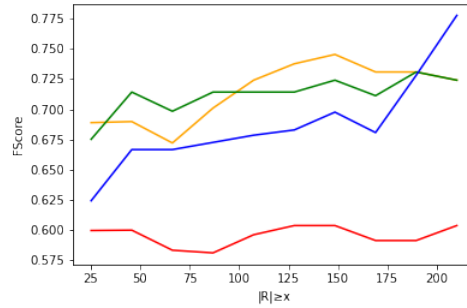


Figure 3: Average F1-score according to the threshold x on the number of preferences in R , for models ORD (blue), KNN (red), SVM (orange), LR (green).

The results on synthetic data are presented in Figures 2 and 3, where the x-axis gives the number of preferences in R (inferred from the ratings of the alternatives in \mathcal{A}_{train}) and the curves show the mean and 95% confidence

interval. The curves show how the different metrics evolve with $|R|$.

Figure 2 shows that each approach produces a different compromise between the number of predictions and their quality. The LR and SVM approaches have, by design, a prediction rate of 1 (the orange line is covered by the green one in the figure) but with predictions that are always less accurate than the predictions made by ORD. Since the KNN approach averages the rates of the K nearest neighbours of the instance to predict, it may occur that two alternatives obtain the exact same score (e.g., if the K nearest neighbours are the same for both alternatives) and thus that no strict preference prediction is made. The prediction rate of KNN is 0.8 on average, but the curve of prediction correctness (and thus the curves of recall and precision) shows that it does not improve the accuracy of the predictions compared to the other methods, quite the contrary.

The ORD model, in contrast, outperforms the other models in terms of precision, as illustrated by the average prediction correctness that is almost always above 0.85. However, since the recall metric penalizes the models that do not make enough predictions, the performances of ORD are below the average performance of the other models in terms of recall. This behavior is, in a sense, intrinsic to an approach that prioritizes the robustness of predictions. Nevertheless, as can be seen in Figure 2c, the recall significantly improves with $|R|$. There are a few irregularities in the curve of the prediction rate for ORD, due to the fact that $\deg(R)$ grows in steps with $|R|$, and this degree impacts $|\Theta^R|$ and thus the number of predictions made (the ordinal dominance relationship becoming more stringent).

While the interest of a compromise between quantity and quality of predictions inherently varies depending on the specific context of an application, the F1-score is a commonly adopted metric to navigate these trade-offs. Figure 3 shows the average F1-score on learning instances where $|R| \geq x$, in function of x . We observe that the average F1-scores of ORD, LR and SVM are close for $x = 50$ (i.e., R include at least 50 preferences). Notably, as x grows so that R encompasses at least 170 preferences, the ORD approach demonstrates a significant performance advantage over LR and SVM.

The curves in Figure 4 gives the average running times of ORD (in seconds, averaged over 20 instances) according to the number n of features (for $300 \leq |R| \leq 400$) and the number $|R|$ of known strict pairwise preferences (for $n = 8$). The orange curve gives the average running time for *one* pairwise preference prediction; this is the most time-consuming phase: note indeed that learning $(\deg(R), \text{card}(R), \text{ws}(R))$ is only performed once for each R , while 100 preference predictions are made for each R in our tests.

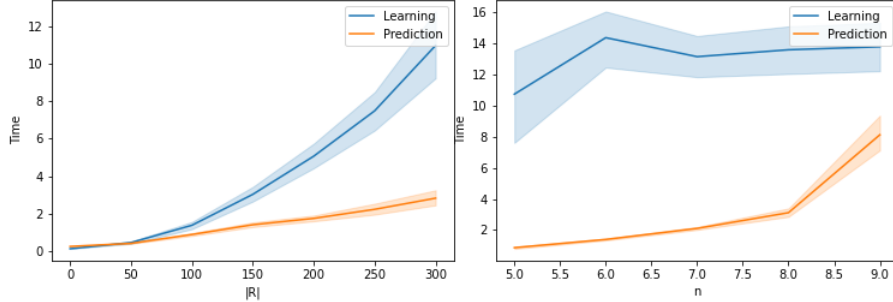


Figure 4: Running times of ORD (in seconds).

5.4.2. Results on real-world data

The results obtained on real-world data from IMDb are summarized in Table 2. Compared to the synthetic data, the precision rates of KNN, LR and SVM significantly decrease, while the precision rate of ORD is holding up better. The recall of ORD remains lower than the recall of LR and that of SVM, but this is overcompensated by the reduced precision performance gap between ORD and LR/SVM. This allows ORD to achieve a better compromise between precision and recall, thus yielding a better F1-Score.

Table 2: Model performances averaged on all the users

Model	Prediction Rate	Precision	Recall	F1-Score
KNN	0.82	0.48	0.65	0.59
LR	1	0.55	0.90	0.69
ORD	0.60	0.76	0.83	0.81
SVM	1	0.55	0.92	0.70

6. Conclusion

We have presented here a robust ordinal method for subsets comparisons with interactions. The model we use is not restrictive, in the sense that any strict weak order on subsets can be represented. The learning method achieves a trade-off between the number of predicted preferences and the accuracy of the predictions, by relying on a robust ordinal dominance relation between subsets.

Several research directions are worth investigating, among which the adaptation of the approach to an active learning setting where one interactively determines a sequence of queries to minimize the cognitive burden for the decision maker, or a better consideration of potential “errors” in the preferences used as a learning set.

Acknowledgements

We acknowledge the support of the French Agence Nationale de la Recherche (ANR), under grant ANR20-CE23-0018 (project THEMIS).

References

- [1] Adam, L. and Destercke, S. (2021). Possibilistic preference elicitation by minimax regret. In *Uncertainty in Artificial Intelligence*, pages 718–727.
- [2] Andreopoulou, Z., Koliouka, C., and Zopounidis, C. (2017). *Multicriteria and Clustering*. Springer.
- [3] Angilella, S., Corrente, S., and Greco, S. (2015). Stochastic multiobjective acceptability analysis for the Choquet integral preference model and the scale construction problem. *European J. of Operational Research*, 240(1):172–182.
- [4] Bartee, E. M. (1971). Problem solving with ordinal measurement. *Management Science*, 17(10):B–622.
- [5] Benabbou, N., Leroy, C., Lust, T., and Perny, P. (2021). Combining preference elicitation with local search and greedy search for matroid optimization. In *Proc. AAAI 2021*, pages 12233–12240. AAAI Press.
- [6] Benabbou, N. and Perny, P. (2015). Combining Preference Elicitation and Search in Multiobjective State-Space Graphs. In *The 24th International Joint Conference on AI (IJCAI’15)*, pages 297–303.
- [7] Bigot, D., Fargier, H., Mengin, J., and Zanuttini, B. (2012). Using and learning gai-decompositions for representing ordinal rankings. In *ECAI’2012 workshop on Preference Learning (PL 2012)*, pages 5–10.
- [8] Blumer, A., Ehrenfeucht, A., Haussler, D., and Warmuth, M. K. (1987). Occam’s razor. *Information Processing Letters*, 24(6):377–380.
- [9] Bourdache, N., Perny, P., and Spanjaard, O. (2019). Incremental elicitation of rank-dependent aggregation functions based on bayesian linear regression. In *Proceedings of IJCAI-19*, pages 2023–2029.
- [10] Boutilier, C., Patrascu, R., Poupart, P., and Schuurmans, D. (2006). Constraint-based optimization and utility elicitation using the minimax decision criterion. *Artificial Intelligence*, 170(8):686–713.
- [11] Brazianus, D. and Boutilier, C. (2007). Minimax regret based elicitation of generalized additive utilities. In *Proceedings of UAI*, pages 25–32.

- [12] Bresson, R., Cohen, J., Hüllermeier, E., Labreuche, C., and Sebag, M. (2020). Learning 2-additive hierarchical choquet integrals with non-monotonic utilities. In *DA2PL 2020*.
- [13] Charnetski, J. R. and Soland, R. M. (1978). Multiple-attribute decision making with partial information: the comparative hypervolume criterion. *Naval Research Logistics Quarterly*, 25(2):279–288.
- [14] Chu, W. and Ghahramani, Z. (2005). Preference learning with gaussian processes. In *Proceedings of ICML-05*, pages 137–144.
- [15] Corrente, S., Greco, S., Kadziński, M., and Słowiński, R. (2013). Robust ordinal regression in preference learning and ranking. *Machine Learning*, 93(2):381–422.
- [16] Dias, L. C., Morton, A., and Quigley, J., editors (2018). *Elicitation : The Science and Art of Structuring Judgement*. Springer.
- [17] Domshlak, C. and Joachims, T. (2005). Unstructuring user preferences: efficient non-parametric utility revelation. In *Proceedings of the Twenty-First Conference on Uncertainty in Artificial Intelligence*, pages 169–177.
- [18] Dragone, P., Teso, S., and Passerini, A. (2017). Constructive preference elicitation over hybrid combinatorial spaces. *CoRR*, abs/1711.07875.
- [19] Fishburn, P. C. (1970). *Utility theory for decision making*. Wiley.
- [20] Fishburn, P. C. and Lavalley, I. H. (1996). Binary interactions and subset choice. *European J. of Operational Research*, 92:182–192.
- [21] Fürnkranz, J. and Hüllermeier, E. (2003). Pairwise preference learning and ranking. In *Proceedings of ECML*, pages 145–156. Springer.
- [22] Gonzales, C. and Perny, P. (2005). GAI networks for decision making under certainty. In *Multidisciplinary IJCAI-05 Workshop on Advances in Preference Handling*, pages 100–105, Edinburgh, United Kingdom.
- [23] Grabisch, M. (1996). The application of fuzzy integrals in multicriteria decision making. *European J. of Operational Research*, 89(3):445–456.
- [24] Grabisch, M., Kojadinovic, I., and Meyer, P. (2008). A review of methods for capacity identification in Choquet integral based multi-attribute utility theory: Applications of the Kappalab R package. *European J. of Operational Research*, 186(2):766–785.
- [25] Grabisch, M., Marichal, J.-L., and Roubens, M. (2000). Equivalent representations of set functions. *Mathematics of OR*, 25(2):157–178.

- [26] Guo, S. and Sanner, S. (2010). Multiattribute Bayesian Preference Elicitation with Pairwise Comparison Queries. In *International Symposium on Neural Networks*, pages 396–403. Springer.
- [27] Hendrickx, K., Perini, L., Van der Plas, D., Meert, W., and Davis, J. (2021). Machine learning with a reject option: A survey. *arXiv preprint arXiv:2107.11277*.
- [28] Herin, M., Perny, P., and Sokolovska, N. (2023). Learning preference models with sparse interactions of criteria. In *Proc. of IJCAI 2023*.
- [29] Kompa, B., Snoek, J., and Beam, A. L. (2021). Second opinion needed: communicating uncertainty in medical machine learning. *NPJ Digital Medicine*, 4(1).
- [30] Kozlov, M. K., Tarasov, S. P., and Khachiyan, L. G. (1979). Polynomial solvability of convex quadratic programming. In *Doklady Akademii Nauk*, volume 248(5), pages 1049–1051. Russian Academy of Sciences.
- [31] Kozlov, M. K., Tarasov, S. P., and Khachiyan, L. G. (1980). The polynomial solvability of convex quadratic programming. *USSR Computational Mathematics and Mathematical Physics*, 20(5):223–228.
- [32] Marichal, J.-L. and Roubens, M. (2000). Determination of weights of interacting criteria from a reference set. *European J. of Operational Research*, 124(3):641–650.
- [33] Papadimitriou, C. H. (1981). On the complexity of integer programming. *Journal of the ACM (JACM)*, 28(4):765–768.
- [34] Schmeidler, D. (1986). Integral representation without additivity. *Proceedings of the American Mathematical Society*, 97(2):255–261.
- [35] Scholkopf, B. and Smola, A. J. (2018). *Learning with kernels: Support Vector Machines, regularization, optimization, and beyond*. MIT press.
- [36] Tehrani, A. F., Strickert, M., and Hüllermeier, E. (2014). The Choquet kernel for monotone data. In *Proc. of ESANN 2014*, pages 337–342.
- [37] Wang, T. and Boutilier, C. (2003). Incremental utility elicitation with the minimax regret decision criterion. In *IJCAI*, volume 3, pages 309–316.
- [38] Zhang, Z., Xu, Y., Yang, J., Li, X., and Zhang, D. (2015). A survey of sparse representation: algorithms and applications. *IEEE access*, 3:490–530.

Appendix A. Properties of the θ -ordinal dominance relation

Proposition 2. *The following properties hold for \succ_θ^R :*

- (i) \succ_θ^R is asymmetric.
- (ii) \succ_θ^R may not be complete.
- (iii) \succ_θ^R is not necessarily negatively-transitive.

Proof. (i) $A \succ_\theta^R B \Rightarrow \forall v \in V_\theta^R, f_{\theta,v}(A) > f_{\theta,v}(B)$. Thus there is no function $v' \in V_\theta^R$ such that $f_{\theta,v'}(A) < f_{\theta,v'}(B)$.

(ii) As shown in Example 2, we may have $v, v' \in V_\theta^R$ such that $f_{\theta,v}(A) > f_{\theta,v}(B)$ and $f_{\theta,v'}(B) > f_{\theta,v'}(A)$. We have then neither $A \succ_\theta^R B$ nor $B \succ_\theta^R A$, and thus \succ_θ^R may not be complete.

(iii) Let $\mathcal{F} = \{a_1, a_2, a_3\}$, $R = \{(\{a_1\}, \{a_3\})\}$, $\theta = \{\{a_1\}, \{a_2\}, \{a_3\}\}$ and v, v' two value functions defined as follows:

$$\begin{aligned} v(\{a_1\}) &= 2, v(\{a_2\}) = 3, v(\{a_3\}) = 1, \\ v'(\{a_1\}) &= 3, v'(\{a_2\}) = 1, v'(\{a_3\}) = 2. \end{aligned}$$

We have that $v, v' \in V_\theta^R$ as $f_{\theta,v}(\{a_1\}) > f_{\theta,v}(\{a_3\})$ and $f_{\theta,v'}(\{a_1\}) > f_{\theta,v'}(\{a_3\})$. It follows from $f_{\theta,v}(\{a_2\}) > f_{\theta,v}(\{a_1\})$ that $\neg(\{a_1\} \succ_\theta^R \{a_2\})$.

It follows from $f_{\theta,v'}(\{a_3\}) > f_{\theta,v'}(\{a_2\})$ that $\neg(\{a_2\} \succ_\theta^R \{a_3\})$.

Yet $\{a_1\} \succ_\theta^R \{a_3\}$ by definition of R . □

Proposition 3. *Given a set R of strict pairwise comparisons, and $\theta \in \Theta^R$, if $R' \subseteq R$ then: (i) $\theta \in \Theta^{R'}$; (ii) $A \succ_\theta^{R'} B \Rightarrow A \succ_\theta^R B$; (iii) $A \succ_\theta^R B \Rightarrow \neg(B \succ_\theta^{R'} A)$.*

Proof. (i) If all the preferences in R can be represented by a θ -additive function, then so can the preferences in R' as R' is compounded of a subset of the preferences in R .

(ii) If the preferences in R' imply that A should be necessarily strictly preferred to B , then R will imply the same conclusion as $\Theta^R \subseteq \Theta^{R'}$ (because R contains all the preference constraints in R' , along with additional constraints).

(iii) The contrapositive is proved as follows: $B \succ_\theta^{R'} A \Rightarrow B \succ_\theta^R A$ by (ii), and $B \succ_\theta^R A \Rightarrow \neg(A \succ_\theta^R B)$ because strict preferences are asymmetrical. □

Proposition 4. *Let $\theta, \theta' \in \Theta^R$. If $\theta' \subseteq \theta$, then the following assertions hold:*

- (i) $A \succ_\theta^R B \Rightarrow A \succ_{\theta'}^R B$; (ii) $A \sim_{\theta'}^R B \Rightarrow A \sim_\theta^R B$; (iii) $A \succ_{\theta'}^R B \Rightarrow \neg(B \succ_\theta^R A)$.

Proof. (i) is true because if $f_{\theta,v}(A) > f_{\theta,v}(B)$ for all $v \in V_\theta^R$, then we should also have $f_{\theta',v}(A) > f_{\theta',v}(B)$ for all $v \in V_{\theta'}^R$. Indeed, each element of $V_{\theta'}^R$ can

be seen as a value function in V_θ^R in which the parameters v_S are set to 0 for $S \in \theta \setminus \theta'$.

(ii) follows by a similar argument as for (i).

(iii) The contrapositive is proved as follows: $B \succ_\theta^R A \Rightarrow B \succ_{\theta'}^R A$ by (i), and $B \succ_{\theta'}^R A \Rightarrow \neg(A \succ_{\theta'}^R B)$ because strict preferences are asymmetrical. \square