

Alma Mater Studiorum Università di Bologna  
Archivio istituzionale della ricerca

A psychometric modeling approach to fuzzy rating data

This is the final peer-reviewed author's accepted manuscript (postprint) of the following publication:

*Published Version:*

Antonio Calcagni, Niccolò Cao, Enrico Rubaltelli, Luigi Lombardi (2022). A psychometric modeling approach to fuzzy rating data. FUZZY SETS AND SYSTEMS, 447, 76-99 [10.1016/j.fss.2022.01.008].

*Availability:*

This version is available at: <https://hdl.handle.net/11585/956033> since: 2024-02-07

*Published:*

DOI: <http://doi.org/10.1016/j.fss.2022.01.008>

*Terms of use:*

Some rights reserved. The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website.

This item was downloaded from IRIS Università di Bologna (<https://cris.unibo.it/>).  
When citing, please refer to the published version.

(Article begins on next page)

# A psychometric modeling approach to fuzzy rating data

Antonio Calcagni<sup>1\*</sup>, Niccolò Cao<sup>1</sup>, Enrico Rubaltelli<sup>1</sup>, Luigi Lombardi<sup>2</sup>

<sup>1</sup>*University of Padova*, <sup>2</sup>*University of Trento*

\* E-mail: antonio.calcagni@unipd.it

## Abstract

Modeling fuzziness and imprecision in human rating data is a crucial problem in many research areas, including applied statistics, behavioral, social, and health sciences. Because of the interplay between cognitive, affective, and contextual factors, the process of answering survey questions is a complex task, which can barely be captured by standard (crisp) rating responses. Fuzzy rating scales have progressively been adopted to overcome some of the limitations of standard rating scales, including their inability to disentangle decision uncertainty from individual responses. The aim of this article is to provide a novel fuzzy scaling procedure which uses Item Response Theory trees (IRTrees) as a psychometric model for the stage-wise latent response process. In so doing, fuzziness of rating data is modeled using the overall rater's pattern of responses instead of being computed using a single-item based approach. This offers a consistent system for interpreting fuzziness in terms of individual-based decision uncertainty. A simulation study and two empirical applications are adopted to assess the characteristics of the proposed model and provide converging results about its effectiveness in modeling fuzziness and imprecision in rating data.

Keywords: fuzzy rating data, fuzzy rating scale, item response model, fuzzy numbers, decision uncertainty

## 1 Introduction

Rating scales are the most common tools for collecting data involving the assessment of interests, motivations, attitudes, personality traits, and a wide variety of health-related and sociodemographic

constructs. A typical use of rating scales is in self-report questionnaires and social surveys where a set of questions (items) are presented individually and respondents are asked to indicate the extent of their agreement/disagreement on a scale with multiple response categories. Overall, rating scales are effective, reliable, and easy to use instruments [1]. However, it is widely recognized that they are not immune to problems such as response biases [2, 3], faking behaviors [4, 5], violation of rating rules [6, 7], and cultural or cognitive differences in the use of response categories [8]. In addition, rating scales do not allow for an in-depth inquire into respondents’ rating process [9]. As many studies have shown, the process of answering multiple choice questions is a complex task since it involves both individual-dependent cognitive and affective factors as well as individual-independent contextual factors (e.g., [10, 11, 12]). For instance, when a respondent is presented with an item like “I am satisfied with my current work”, which is rated on a five-point scale from “strongly disagree” to “strongly agree”, he or she first retrieves long-term memory information about events, attitudes, beliefs about his or her job. The retrieved events may activate affective components which influence positively or negatively the opinion formation (for example, a recent promotion may enhance the chance for answering the item positively). Then, cognitive and affective information are integrated to activate the decision making stage, which includes the answer editing step where a set of candidate answers is pruned to produce the final response [13]. As a result of conflicting demands from the latter stages, some levels of decision uncertainty can impact the final rating choice. Consequently, final responses on questionnaires only reflect a portion of the entire response process.

There have been numerous attempts to make rating scales more sensitive to detecting components of response process such as decision uncertainty. Generally, there are three types of solutions to this problem, a first one involving the use of additional measures like response times or latencies along with standard multiple choice questions [14, 15, 16], a second one using extended item response theory (IRT) models on standard rating data [17, 18, 19], and a third one involving the use of alternative rating instruments such as those based on tracing methodology [20] and fuzzy rating scales [21, 22, 23]. Since the seminal work of [23], the latter has become popular only in recent years. Typically, there are two ways to define a fuzzy rating scale, one involving fuzzy conversion systems and the other involving a direct fuzzy rating system. In the first case, a fuzzy conversion system is used to transform standard rating responses into fuzzy numbers (e.g., see [24]). In the second case, a tailor-made rating interface is instead adopted in order to map fuzzy numbers to a rating process by means of implicit [21] or explicit [22] procedures. Despite their differences, both the approaches aim at modeling decision uncertainty or its counterpart, fuzziness and imprecision of rating data, as emerging from multiple choice rating tasks. The role of fuzziness in rating and psychometric data has been highlighted by several researchers working at the interface between statistics and applied mathematics (e.g., [25, 26, 27, 28, 29]).

In this paper, we contribute to this research stream by proposing a novel method which places fuzzy rating scales in the context of Item Response Theory trees (IRTrees) models [17]. The aim is to provide an approach to fuzzy rating data that incorporates a stage-wise cognitive formalization

of the process that respondents use to answer survey questions. IRTrees are a novel class of item response models aiming at representing the internal decision stages behind final rating outcomes. By adopting a sequence of linear or nested binary trees, they allow for disentangling the result of the rating process (e.g., the choice of the category “strongly disagree” on a common Likert-type scale) and the sequential steps needed by raters to reach their final outcomes. In this manner they provide an elegant way to mine information from rating data, which can be used to model fuzziness and imprecision encapsulated into rating data.

Although the proposed method does not conflict with existing standard methods for fuzzy rating, there are some differences that should be highlighted along with some advantages as well as disadvantages. First, the novel approach is grounded on a psychometric formalization of the rating process and uses a statistical method (IRTree) to model the observed rating data in advance. This results in a different characterization of the rating fuzziness, which does no longer represent the actual degree of confidence a rater has in providing his/her rating response. Rather, it represents the conflicting demands provoked by the decision stage which precedes the expression of the final rating response. In this sense, fuzzy-IRTree based fuzzy sets are computed as a function of the IRTree parameters instead of being derived from the data directly. Second, the new method does not require specialized computerized interfaces through which measuring fuzzy sets, with the consequence that it can be widely used with standard rating scale formats. Finally, fuzzy-IRTree avoids the use of direct rating methods which might be potentially affected by cognitive biases regarding the direct estimation of numerical quantities. However, as for any statistical model, a potential limitation of the proposed method is that it requires a sufficiently large sample size and number of items in order to get reliable results. Similarly, as it is based on IRTrees, the psychometric model formalizing the rating response process should be chosen in advance. In this case, it might be advisable for data analysts to refer to existing scientific literature or to use a statistically-oriented procedure to find the best IRTree model given the sample data (e.g., this can be done by means of AIC based model comparison).

The reminder of this article is organized as follows. Section 2 offers a review of the major literature about fuzzy rating scales. Section 3 describes our method for modeling imprecision and uncertainty in rating data using IRTrees. Section 4 reports the results of a simulation study designed to validate our proposal whereas Section 5 describes two applications using empirical case studies. Finally, Section 6 concludes the article by providing final remarks and suggestion for future research. All the materials like algorithms and datasets used throughout the paper are available to download at <https://github.com/antcalcagni/firtree/>.

## 2 Currently used methods in fuzzy rating

Fuzzy rating scales aim at quantifying fuzziness and imprecision of human subjective responses. Typically, there are two approaches known in the literature to construct a fuzzy rating instrument,

namely fuzzy direct or indirect scales and fuzzy conversion scales.

In *fuzzy direct rating*, a computerized rating scale is adopted and raters are asked to draw their responses using fuzzy sets according to their perceived uncertainty [23, 30]. This method usually require a two-step response process. First, raters draw an interval or a point on a pseudo-continuous graphical scale which represent the set of admissible responses compatible with their assessment of the item being rated. Then, they are asked to express the degree of confidence by drawing another interval about their previous interval or point-wise responses. Finally, the two information are combined to form triangular or trapezoidal fuzzy responses. An overview of direct fuzzy rating is described in [31]. By contrast, the *fuzzy indirect rating* uses implicit subjective information to quantify the fuzziness of rating data. On this research line, for instance, [21] adopted a system which includes biometric measures of cognitive response process (e.g., response time, computer-mouse trajectories) in the construction of fuzzy responses. Despite their differences, both the approaches have successfully been adopted to measure psychological constructs [32], to evaluate students' perceptions and feelings [33], to measure gendered beliefs [34], to inspect experience of perplexity [35], to evaluate the quality of linguistic descriptions [36], to explore physicals' perception of mental patients [37], to evaluate service quality [38] as well as the quality of products [39].

Unlike direct or indirect fuzzy rating, fuzzy conversion scales adopt stochastic or deterministic procedures (e.g., fuzzy systems) to convert crisp rating data - usually collected by means of traditional rating tools (e.g., Likert-type scales) - into fuzzy sets with the aim of obtaining an improvement of the scaling procedure. To this end, a number of conversion systems have been proposed, which are mainly based on expert-knowledge, empirical-based or indirect methods [40]. Among them, *expert-knowledge* conversion systems use a-priori information to derive fuzzy categories through which crisp data are fuzzified. For instance, [24] proposed an improved Likert-type scale based on a deterministic Mamdadi fuzzy system which includes fuzzification and defuzzification steps. On this line, [41] compared Likert-type scale and three fuzzy conversion scales based on triangular, trapezoidal, and Gaussian fuzzy numbers, respectively. This type of fuzzy scaling has been widely applied, for instance, in measuring user experience [42], workers' motivation [43], teachers' beliefs about mathematics [44], students' perceptions about learning through a computer algebra system [45], motivation, attention and anxiety [46], job satisfaction [47], tourists' satisfaction [48, 49, 50], in evaluating healthcare services [51], educational services [52, 53, 54], and to develop methodologies for service quality analysis [55, 56, 57, 58]. Instead, *empirical-based* fuzzy conversion methods transform crisp responses into fuzzy data using the information gathered directly from the empirical sample of responses. For example, [59] developed a fuzzy system in which fuzzy categories are built based on empirical distribution of Likert-type responses. Similarly, [60] developed a fuzzy system to measure xenophobia through pollster method and frequency-based fuzzy set assignment. Still, [61, 62] and [63] proposed to generate fuzzy categories via Dombi-intersection of sigmoid-shaped functions based on the most likely, worst and best values assigned by raters. In a similar way, [64] derived fuzzy numbers using histograms of Likert-type responses and ideal histograms-based distances for modeling

response-bias. Finally, *indirect methods* to fuzzy conversion scales use hybrid systems through which fuzzy data are obtained by means of statistical models which are adapted on empirical crisp data first. For instance, [65] proposed an innovative method where CUB models are used as a back-end tool for quantifying fuzziness of rating responses. Similarly, [40] used ordinal regression in order to generate well-founded fuzzy response categories. [66] proposed a statistically-oriented procedure by means of which fuzzy sets are computed using non-parametric spline methods. On the same line, [67] and [68] used an Item Response Theory model (i.e., Partial Credit Model) to convert linguistic response categories in fuzzy numbers by means of the estimated IRT parameters.

There have been several attempts to compare fuzzy rating and conversion scales with respect to more traditional rating methods. To this end, comparisons have been made based on hypothesis testing about means [69, 70], descriptive summary measures [30], ratings accordance criterion in empirical and simulated context [71, 72], scale reliability [73, 74]. Other research used validated questionnaires to study the differences between traditional and fuzzy rating. For example, [75] used the WHOQOL-BREF questionnaire to compare standard Likert-type scale, fuzzy direct scale, and two fuzzy conversion scales. In a similar way, [76] proposed and compared four fuzzy version of the pain intensity scales, namely fuzzy visual analogue scale, fuzzy numerical rating scale, fuzzy qualitative pain scale, and fuzzy face pain scale.

### 3 An IRTree-based model for fuzzy rating

In this section we illustrate our approach to fuzzy rating scales, which is based upon the use of IRT trees as computational models of the response process [17]. In particular, we adopt a two-stage modeling strategy where IRTrees are first fit on rating data and then their estimated parameters are mapped to parametric fuzzy numbers [68]. In so doing, a psychometric model is used to model response data for each rater and item combination, which is in turn used as a building block for representing final ratings in terms of fuzzy numbers.

#### 3.1 IRT models

Item Response Theory (IRT) models represent a class of statistical models which are used to formalize the measurement process underlying self-reported responses in questionnaires, tests, and surveys. Being at the intersection of psychometrics and statistics, they offer a way to formalize the underlying process a rater  $i$  responds to a given item  $j$  [77]. Although there are a number of IRT models available nowadays (for an extensive review, see [78]), they all revolve around the assumption that the probability  $\mathbb{P}(Y_{ij} = y; \boldsymbol{\theta})$  of responding to an item  $j$  for a given rater  $i$  is a function of at least two parameters, namely the quality of the item  $\alpha_j$  (e.g., difficulty, informativeness) and the characteristics of the rater  $\eta_i$  (e.g., latent personality trait, response style). In formulae, we have

$$\mathbb{P}(Y_{ij}=y; \boldsymbol{\theta}) = g(\alpha_i, \eta_j)$$

with  $g(\cdot)$  being a twice differentiable link function (e.g., logistic, probit, generalized logistic). Depending on the complexity of the psychometric model being used, the basic IRT formulation can be generalized to include many other information such as covariates (e.g., gender, age), additional rater’s information (e.g., careless responding, lucky guessing), and questionnaire structure (e.g., latent dimensions connecting items among them). Because of their characteristics, IRT models are quite closed to Generalized Linear Mixed-Effect Models (GLMMs), another class of linear statistical models widely used in applied statistical analyses [79]. Indeed, parameters of IRT models are conventionally estimated using methods typically adopted by GLMMs such as marginal maximum likelihood, expectation-maximization, and pairwise maximum likelihood. The simplest IRT model is the well-known Rasch model (also called, 1-PL IRT model) which formalizes the probability of responding to a dichotomous item  $Y_{ij} \in \{0, 1\}$  as follows:

$$\mathbb{P}(Y_{ij} = 1; \boldsymbol{\theta}) = \frac{1}{1 + \exp(\eta_i - \alpha_j)}$$

where  $\eta_i \sim \mathcal{N}(0, \sigma_\eta^2)$  and  $\alpha_j \in \mathbb{R}$  with the constraint  $\sum_{j=1}^J \alpha_j = 0$ . The model formalizes the intuition that for a fixed item  $j$ , the probability of a correct response  $\mathbb{P}(Y_{ij} = 1)$  increases with the rater’s ability  $\eta$  and, conversely, for a fixed rater  $i$  the probability of a right response decreases as the item difficulty  $\alpha$  increases. Given a set of  $n$  responses to  $J$  items, the 1PL parameters can be estimated via maximum-likelihood theory and their estimates  $\hat{\boldsymbol{\alpha}}$  and  $\hat{\boldsymbol{\eta}}$  can be used for further analyses, including assessing tests or questionnaires to measure a particular ability or trait, estimating raters information about response styles, evaluating difficulties of items.

### 3.2 IRTrees and response process

IRT trees are conditional linear models that represent final rating responses in terms of binary trees. They formalize the response process as a sequence of conditional stages going through the tree to end nodes. Intermediate nodes are defined such that they represent specific cognitive components of the rating process whereas end nodes represent the possible outcomes of the decision process. Figure 1a depicts the simplest IRT tree model for three-point rating scales (0: “perhaps”; 2: “yes”; 1: “no”). It contains two intermediate nodes, one representing the first stage of the response process  $Z_1$  (e.g., answering with uncertainty vs. answering with certainty) with a single outcome (e.g.,  $Y = 0$ : “perhaps”) and the other representing the second decision stage  $Z_2$  (e.g., answering with certainty) with two possible outcomes (e.g.,  $Y = 2$ : “yes” vs.  $Y = 1$ : “no”). For instance, the probability of uncertain responses (i.e.,  $Y = 0$ : “perhaps”) is simply given by the probability to activate the first stage of the decision process, i.e.  $\mathbb{P}(Y = 0) = \mathbb{P}(Z_1; \boldsymbol{\theta}_1)$ . By contrast, the probability of a negative response (i.e.,  $Y = 1$ : “no”) is computed as  $\mathbb{P}(Y = 1) = (\mathbb{P}(Z_1; \boldsymbol{\theta}_1)(1 - \mathbb{P}(Z_2; \boldsymbol{\theta}_2)))$ . The simplest case described by Figure 1a is paradigmatic of the cognitive modeling underlying IRT trees [80, 81]. These models assume the rater’s response process to be stage-wise: raters would first decide whether or not provide their responses ( $Z_1$ ) and, then, decide on the direction and strength of their answers

( $Z_2$ ). The latent random variables  $Z_1$  and  $Z_2$  govern the two sub-processes of the rater's response. Similarly, Figure 1b generalizes a two-stage decision tree for the common five-point rating scale (e.g., from 1: “strongly disagree”; to 5: “strongly agree”). It contains three decision nodes, one for the uncertain response category (i.e.,  $Y = 3$ : “neither agree, nor disagree”), a second one for the levels of disagreement (i.e.,  $Y = 1$ : “strongly disagree”,  $Y = 2$ : “disagree”), and the last node for the levels of agreements (i.e.,  $Y = 4$ : “strongly agree”,  $Y = 5$ : “disagree”). Probabilities for each response are computed as before. Figures 1c-1d represents two cases of IRTrees for a six-point rating scale. The trees differ in the way they model the middle categories (i.e.,  $Y = 3$  and  $Y = 4$ ). In the first schema (Figure 1c), they are represented independently from the extremes of the scale, as for the two-stage IRTree (Figure 1a). By contrast, the second schema (Figure 1c) places the middle categories in the same branches of the extremes, as to represent a more graded decision process [82]. There are many possible ways to conceptualize decision processes in terms of IRTrees and the choice of a particular decision schema depends primarily on research-specific hypotheses [81].

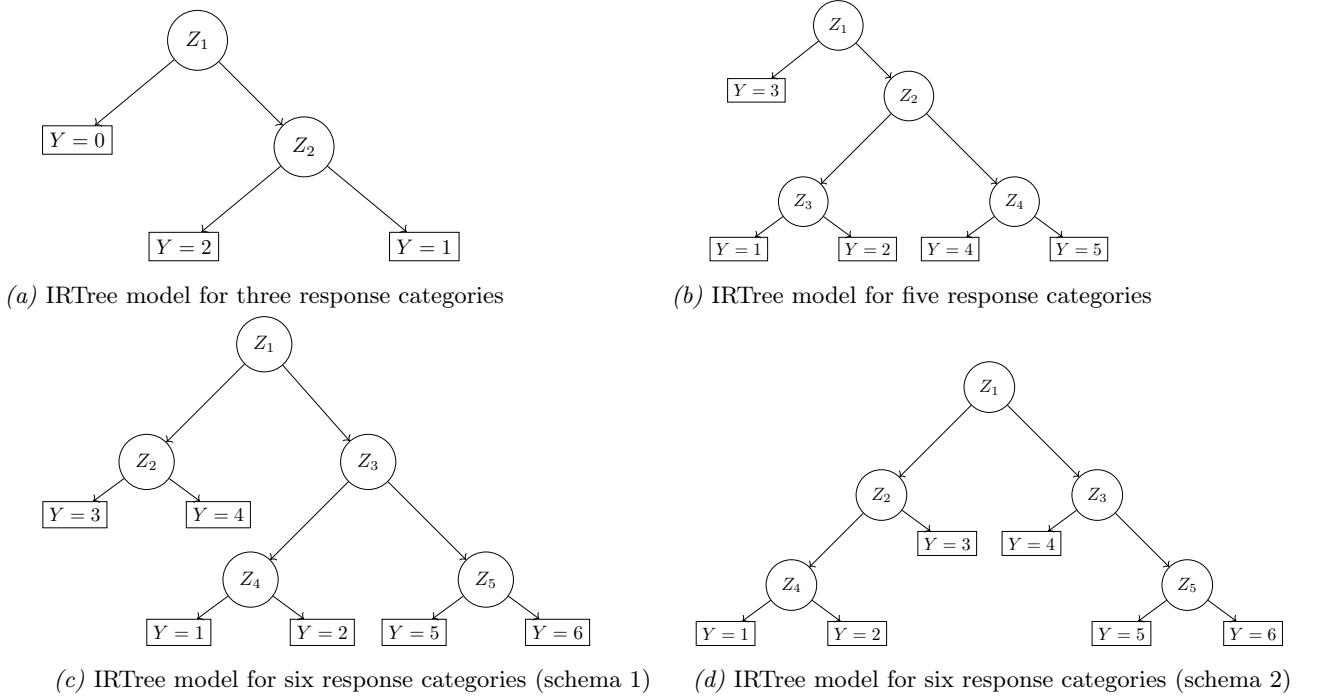


Figure 1. Examples of IRTree models for modeling response processes in rating scales.

By using an IRT parameterization, IRTrees allow for introducing rater-specific and item-specific components for the response process. Hence, the probability to agree or disagree with an item can be represented as a function of a rater's latent trait and the specific content of the item [17]. More formally, let  $i \in \{1, \dots, I\}$  and  $j \in \{1, \dots, J\}$  be the indices for raters and items, respectively. Then, the final response variable  $Y_{ij} \in \{1, \dots, m, \dots, M\} \subset \mathbb{N}$ , with  $M$  being the maximum number



of response categories, can be decomposed in terms of binary responses using  $N$  binary variables  $Z_{ijn} \in \{0, 1\}$ , where  $n \in \{1, \dots, N\}$  denotes the nodes of the tree. For instance, in Figure 1a,  $N = 2$  and the final response  $Y_{ij} = 2$  corresponds to  $Z_{ij2} = 0$ . By following the common Rasch representation [17], for a generic pair  $(i, j)$  the IRTree consists of the following equations:

$$\boldsymbol{\eta}_i \sim \mathcal{N}_N(\mathbf{0}, \boldsymbol{\Sigma}_\eta) \quad (1)$$

$$\pi_{ijn} = \mathbb{P}(Z_{ijn} = 1; \boldsymbol{\theta}_n) = \frac{\exp(\eta_{in} + \alpha_{jn})}{1 + \exp(\eta_{in} + \alpha_{jn})} \quad (2)$$

$$Z_{ijn} \sim \text{Ber}(\pi_{ijn}) \quad (3)$$

where  $\boldsymbol{\theta}_n = \{\boldsymbol{\alpha}_j, \boldsymbol{\beta}_i\}$ , with the arrays  $\boldsymbol{\alpha}_j \in \mathbb{R}^N$  and  $\boldsymbol{\eta}_i \in \mathbb{R}^N$  denoting the easiness of the item and the rater's latent trait. As is usual in IRT models, latent traits for each node are modeled using a  $N$ -variate centered Gaussian distribution with covariance matrix  $\boldsymbol{\Sigma}_\eta$ . For instance, for the two-stage decision process in Figure 1a,  $\alpha_{j1}$  indicates the easiness of choosing the right branch of the tree for the item  $j$  whereas  $\alpha_{j2}$  denotes the easiness of providing an affirmative response ( $Y = 1$  : "yes"). Similarly,  $\eta_{i1}$  indicates the rater's attitude to navigate through the right branch of the tree whereas  $\eta_{i2}$  denotes the rater's attitude to provide an affirmative response. Thus, the probabilities to activate a branch of the tree can be computed using Eq. (2) recursively. For instance, in the two-stage example, the probability of an uncertain response is computed as follows:

$$\mathbb{P}(Y_{ij} = 0) = \mathbb{P}(Z_{ij1} = 0; \boldsymbol{\theta}_1) = 1 - \frac{\exp(\eta_{i1} + \alpha_{j1})}{1 + \exp(\eta_{i1} + \alpha_{j1})}$$

To generalize single-branch probability equations, we first define a  $M \times N$  Boolean matrix  $\mathbf{T}$  indicating how each response category (in rows) is associated to each node (in columns) of the tree. As  $t_{mn} \in \{0, 1\}$ ,  $t_{mn} = 1$  indicates that the  $m$ -th category of response involves the node  $n$ ,  $t_{mn} = 0$  indicates that the  $m$ -th category of response does not involve the node  $n$ , whereas  $t_{mn} = \text{NA}$  indicates that the  $m$ -th category of response is not connected to the  $n$ -th node at all. For instance, considering the simplest two-stage example in Figure 1a, the mapping matrix  $\mathbf{T}_{3 \times 2}$  is defined as follows:

$$\mathbf{T} = \begin{bmatrix} 1 & \text{NA} \\ 1 & 0 \\ 1 & 1 \end{bmatrix}$$

Finally, the probability for a generic rating response can be easily computed as:

$$\begin{aligned} \mathbb{P}(Y_{ij} = m) &= \prod_{n=1}^N \mathbb{P}(Z_{ijn} = t_{mn}; \boldsymbol{\theta}_n)^{t_{mn}} \\ &= \prod_{n=1}^N \left( \frac{\exp(\eta_{in} + \alpha_{jn})}{1 + \exp(\eta_{in} + \alpha_{jn})} \right)^{\delta_{mn}} \end{aligned} \quad (4)$$

where  $\delta_{mn} = 0$  if  $t_{mn} = \text{NA}$  and  $\delta_{mn} = 1$  otherwise.

IRTree models can be estimated either by means of standard methods used for generalized linear mixed models, such as restricted or marginal maximum likelihood [83, 17], or using procedures for multidimensional item response theory models, such as expectation-maximization algorithms [84]. In general, these models are flexible enough to model simple situations like those requiring unidimensional latent variables (a single  $\eta$  for each node of the tree) or common item effects (a single  $\alpha$  for each node of the tree) as well as more complex scenario involving multidimensional high-order latent variables. For further details and implementations, we refer the reader to [17, 84].

### 3.3 Fuzzy numbers

A fuzzy set  $\tilde{A}$  of a universal set  $\mathcal{A}$  is defined by means of its characteristic function  $\xi_{\tilde{A}} : \mathcal{A} \rightarrow [0, 1]$ . It can be easily described as a collection of crisp subsets called  $\alpha$ -sets, i.e.  $\tilde{A}_\alpha = \{y \in \mathcal{A} : \xi_{\tilde{A}}(y) > \alpha\}$  with  $\alpha \in (0, 1]$ . If the  $\alpha$ -sets of  $\tilde{A}$  are all convex sets then  $\tilde{A}$  is a convex fuzzy set. The support of  $\tilde{A}$  is  $A_0 = \{y \in \mathcal{A} : \xi_{\tilde{A}}(y) > 0\}$  and the core is the set of all its maximal points  $A_c = \{y \in \mathcal{A} : \xi_{\tilde{A}}(y) = \max_{y \in \mathcal{A}} \xi_{\tilde{A}}(y)\}$ . In the case  $\max_{y \in \mathcal{A}} \xi_{\tilde{A}}(y) = 1$  then  $\tilde{A}$  is a normal fuzzy set. If  $\tilde{A}$  is a normal and convex subset of  $\mathbb{R}$  then  $\tilde{A}$  is a fuzzy number. The quantity  $l(\tilde{A}) = \max A_0 - \min A_0$  is the length of the support of the fuzzy set  $\tilde{A}$ . The class of all normal fuzzy numbers is denoted by  $\mathcal{F}(\mathbb{R})$ . Fuzzy numbers can conveniently be represented using parametric models that are indexed by some scalars, such as  $c$  (mode) and  $s$  (spread or precision). These include a number of shapes like triangular, trapezoidal, gaussian, and exponential fuzzy sets [85]. A relevant class of parametric fuzzy numbers are the so-called LR-fuzzy numbers [86] and their generalizations like non-convex fuzzy numbers [87], flexible fuzzy numbers [62], and beta fuzzy numbers [88, 89, 90]. The latter represent a special class of fuzzy sets that are defined by generalizing triangular fuzzy sets. In particular, let:

$$\xi_{\tilde{A}}(y) = \left( \frac{y - y_l}{c - y_l} \right) \cdot \mathbb{1}_{(y_l, c)}(y) + \left( \frac{y_u - y}{y_u - c} \right) \cdot \mathbb{1}_{(c, y_u)}(y) \quad (5)$$

be a triangular fuzzy set with  $y_l, y_u, c \in \mathbb{R}$  being lower, upper bounds, and mode parameters, respectively. Then, a beta fuzzy set is of the form:

$$\begin{aligned} \xi_{\tilde{A}}(y) &= \left( \frac{y - y_l}{c - y_l} \right)^a \left( \frac{y_u - y}{y_u - c} \right)^b \cdot \mathbb{1}_{(y_l, y_u)}(y) \\ c &= \frac{ay_u + by_l}{a + b} \end{aligned} \quad (6)$$

where  $y_l, y_u, a, b \in \mathbb{R}$ , with  $y_l$  and  $y_u$  being the lower and upper bounds of the set, and  $c$  the mode of the fuzzy set. Beta fuzzy numbers can be expressed in terms of mode  $c \in \mathbb{R}$  and precision  $s \in \mathbb{R}^+$

parameters, as follows ( $y_l = 0$  and  $y_u = 1$  without loss of generality):

$$\xi_{\tilde{A}}(y) = \frac{1}{C} y^{a-1} (1-y)^{b-1} \quad (7)$$

$$a = 1 + cs$$

$$b = 1 + s(1 - c)$$

$$C = \left( \frac{a-1}{a+b-2} \right)^{a-1} \cdot \left( 1 - \frac{a-1}{a+b-2} \right)^{b-1} \quad (8)$$

with  $C$  being a constant ensuring  $\xi_{\tilde{A}}$  is still a normal fuzzy set. Figure 2 shows some examples of beta fuzzy sets (dashed black curves). Because of their shape, beta-based fuzzy sets can be of particular utility in modeling bounded rating data (e.g., see [91]).

### 3.4 An IRT-map between fuzzy numbers and rating responses

Consider the case where a respondent  $i$  is faced with a  $M$ -choice item  $j$ . In the first stage of the response process, the item content first triggers memories and emotions of past personal experiences. Then, these activate the opinion formation stage, where a coherent opinion representation is formed along with a finite set of potential responses  $\mathcal{U}_{ij}$ . Lastly, the final response  $y_{ij}$  is chosen by trimming the set of possible responses (selection stage). Decision uncertainty emerges as a result of the conflicting demands of the opinion formation stage and it can be quantified by analysing some characteristics of  $\mathcal{U}_{ij}$ . Our approach resorts to using the latter as a source for mapping fuzzy numbers to the latent rater's response process underlying  $y_{ij}$ . To this end, IRT-trees are adopted to estimate a probabilistic model for  $\mathcal{U}_{ij}$  as a function of estimated rater's latent traits  $\hat{\eta}_i$  and item content  $\hat{\alpha}_j$ . In particular, for a given pair  $(i, j)$  the following procedure is used to obtain fuzzy rating data:

1. Define and fit an IRT-tree model to a sample of  $I \times J$  responses  $\mathbf{Y}$  and get the estimates  $\hat{\eta}_{N \times 1}$  and  $\hat{\alpha}_{N \times 1}$ .
2. Plug-in  $\hat{\eta}_{N \times 1}$  and  $\hat{\alpha}_{N \times 1}$  into Eq. (4) to get the estimated probability value  $\hat{\mathbb{P}}(Y = m)$  for each  $m \in \{1, \dots, M\}$ . This is the probabilistic model for  $\mathcal{U}_{ij}$ .
3. Compute the mode of the fuzzy beta number  $\tilde{y}_{ij}$  via the equality:

$$c_{ij} = \sum_{y \in \{1, \dots, M\}} y \cdot \hat{\mathbb{P}}(Y = y) \quad (9)$$

4. Compute the precision of the fuzzy beta number  $\tilde{y}_{ij}$  via the equality:

$$s_{ij} = \frac{1}{v_{ij}} \quad \text{with: } v_{ij} = \sum_{y \in \{1, \dots, M\}} (y - c_{ij})^2 \cdot \hat{\mathbb{P}}(Y = y) \quad (10)$$

In this context,  $\xi_{\tilde{y}_{ij}} : \Omega(y) \rightarrow (0, 1)$ , with  $\Omega(y) = (1, M)$  being the space of the means of  $Y_{ij}$  for each response value. Thus, likewise for latent responses in psychometric models, fuzzy rating data are continuous and bounded instead of being discrete. Note that the above procedure is quite general and can be extended to the more general case of LR-type fuzzy numbers, such as triangular and trapezoidal, by means of any probability-possibility transformations [86] or other general transformations preserving the original information content [92]. For instance, the easiest way to obtain triangular fuzzy numbers from  $\hat{\mathbb{P}}(Y)$  is to compute the core using Eq. (9) whereas lower  $y_{l_{ij}}$  and upper  $y_{u_{ij}}$  bounds can instead be computed using quantiles, such as  $y_{l_{ij}} = \min(\{y \in \{1, \dots, M\} : \hat{\mathbb{P}}(y) \geq 0\})$  and  $y_{u_{ij}} = \max(\{y \in \{1, \dots, M\} : \hat{\mathbb{P}}(y) \geq 0\})$ . Another solution would be to transform fuzzy beta numbers using a kind of moments matching method [93] via the following link equations:

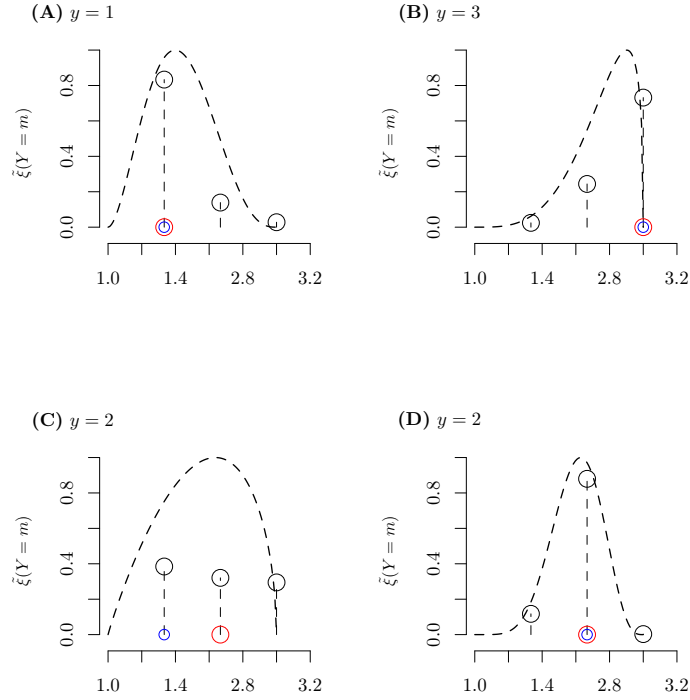
$$\begin{aligned} y_{l_{ij}} &= c_{ij} - h_2, & y_{r_{ij}} &= c_{ij} - h_2 + h_1 \\ h_1 &= \sqrt{3.5v_{ij} - 3(c_{ij} - \mu_{ij})^2} \\ h_2 &= \frac{1}{2}(h_1 + 3c_{ij} - 3\mu_{ij}) \\ \mu_{ij} &= (1 + c_{ij}s_{ij})/(2 + s_{ij}) \end{aligned} \tag{11}$$

The procedure yields regular triangular fuzzy sets defined in terms of lower bound  $y_l$ , mode  $c$ , and upper bound  $y_u$ .

Figure 2 shows some hypothetical examples of fuzzy beta numbers for a two-stage IRTree with  $M = 3$  and  $N = 2$ . As a direct consequence of our modeling approach - which is based upon the use of heterogeneity in rater's pattern of responses - the final response  $y_{ij}$  may not reflect the mode of the fuzzy response  $c_{ij}$  (or similarly other measures like the centroid). This is particularly true for high uncertainty scenarios where two or more responses compete with each other (see Figure 2-c). Thus, decision uncertainty does not necessarily coincide with the choice of the middle or “don't know” response category of the rating scale. Rather, it arises as a result of the transitions probabilities estimated by the IRTree (the easier the transition is, the more certain the response is). This is the case, for instance, shown in Figure 4-d where the middle response category is chosen with little uncertainty.

## 4 Simulation study

The aim of this simulation study is to provide an external validity check on the results provided by the fuzzy IRT-map to recover decision uncertainty from rating tasks. In particular, our model was contrasted against another IRT model for rating data that uses response times (RTs) as a source for modeling decision uncertainty [14, 94]. It is well established that RTs can be used for measuring several cognitive facets such as item/question difficulties and participants' performance on rating and choice tasks [95]. Overall, the findings from the psychometric literature suggest that respondents



*Figure 2.* Examples of hypothetical probability distributions (black dashed vertical lines) and associated fuzzy numbers (black dashed curves) for a two-stage IRTree ( $M = 3$  and  $N = 2$ ). Note that probability masses and fuzzy membership functions are overlapped over the same domain  $\Omega(y)$ , red and blue circles represent observed ( $y$ ) as opposed to most probable responses, respectively.

who are very hesitant and uncertain about their final answers take a relatively long time to make their final choice on a rating scale [94]. Conversely, respondents who are quite sure of their responses are generally fast in providing their final choices. As such, RTs can be considered valuable indirect measures of decision uncertainty in rating tasks [96]. In this study we assessed whether the fuzzy IRT-map can retrieve decision uncertainty from rating data as accurate as response times. To this end, first we will generate rating data and response times according to a dedicated IRT-RTs model, and then we will apply the fuzzy IRT-map on the rating data by evaluating to what extent fuzzy numbers computed via the fuzzy IRT-map will predict response times that were generated using the IRT-RTs model. The whole simulation study has been performed on a remote HPC machine based on 16 Intel Xeon CPU E5-2630Lv3 1.80Ghz, 16x4 Gb Ram whereas computations and analyses have been performed in the R framework for statistical analyses.

*Data generation model.* Discrete rating data  $Y_{ij} \in \{1, \dots, M\}$  and response times  $R_{ij} \in (0, \infty)$  for respondent  $i \in \{1, \dots, I\}$  and item  $j \in \{1, \dots, J\}$  were generated according to the following IRT-RTs model [94]:

$$\eta_i \sim \mathcal{N}(0, \sigma_\eta), \quad \omega_i \sim \mathcal{N}(0, \sigma_\omega), \quad \epsilon_i \sim \mathcal{N}(0, \sigma_\epsilon)$$

$$\mathbb{P}(Y_{ij} = m; \boldsymbol{\theta}) = \frac{\exp(\sum_{k=1}^m (\eta_i - \alpha_j))}{\sum_{h=1}^M \exp(\sum_{k=1}^h (\eta_i - \alpha_j))} \quad (12)$$

$$\ln r_{ij} = \gamma_j + \omega_i + \left( \sum_{m=1}^M \mathbb{P}(Y_{ij} = m; \boldsymbol{\theta})^2 \right) \beta_j + \epsilon_{ij} \quad (13)$$

where  $\boldsymbol{\eta}_{I \times 1}$  and  $\boldsymbol{\alpha}_{J \times 1}$  are respondents' latent traits and item parameters,  $\boldsymbol{\omega}_{I \times 1}$  and  $\boldsymbol{\gamma}_{J \times 1}$  are respondents' speeds and item times,  $\boldsymbol{\beta}_{J \times 1}$  are the time intensity parameters which relate the response data submodel in Eq. (12) to the response time submodel in Eq. (13). The term  $\sum_{m=1}^M \mathbb{P}(Y_{ij} = m; \boldsymbol{\theta})^2$  in the response time submodel can be interpreted as the difficulty for the respondent  $i$  to respond to the item  $j$  (DIFF) and is closely related to the so-called Probability-Difficulty (PD) hypothesis in the IRT literature [14]. The DIFF-based model for RTs states that longer response times occur when DIFF is lower, which is the case where all the  $M$  alternatives for the item  $j$  are equally probable. By contrast, shorter response times are expected when DIFF is higher, which is the opposite case where there is single a response category with probability equals to one [94].

*Design.* The design of the study involved four factors: (i)  $I \in \{50, 100, 150\}$ , (ii)  $J \in \{5, 20\}$ , (iii)  $M \in \{3, 5\}$ , (iv)  $\beta^0 \in \{-10.5, -20.5\}$ . They were varied in a complete factorial design with a total of  $3 \times 2 \times 2 \times 2 = 24$  scenarios. For each combination,  $B = 1000$  samples were generated which yielded to  $1000 \times 24 = 24000$  new data as well as an equivalent number of parameters.

*Procedure.* Let  $i_h, j_t, m_p, \beta_q^0$  be distinct levels of factors  $I, J, M, \beta^0$ . Then, rating data and response times were generated according to the following procedure:

- (a) Respondents' latent traits and speeds were drawn independently as  $\boldsymbol{\eta}_{i_h \times 1} \sim \mathcal{N}(\mathbf{0}_{i_h}, \mathbf{I}_{i_h \times i_h})$  and  $\boldsymbol{\omega}_{i_h \times 1} \sim \mathcal{N}(\mathbf{0}_{i_h}, \mathbf{I}_{i_h \times i_h})$ .
- (b) Item parameters and average response times were generated independently as  $\boldsymbol{\alpha}_{j_t \times 1} \sim \mathcal{N}(\mathbf{0}_{j_t}, \mathbf{I}_{j_t \times j_t})$  and  $\boldsymbol{\gamma}_{j_t \times 1} \sim \mathcal{N}(9\mathbf{1}_{j_t}, \mathbf{I}_{j_t \times j_t})$ .
- (c) For  $i = 1, \dots, i_h$  and  $j = 1, \dots, j_t$ , probabilities for each of the  $m_p$  response categories were computed using the IRT component of the IRT-RTs model:

$$\mathbb{P}(Y_{ij} = m) = \exp\left(\sum_{k=1}^{m_p} (\eta_i - \alpha_j)\right) / \sum_{u=1}^{m_p} \exp\left(\sum_{k=1}^{m_p} (\eta_i - \alpha_j)\right)$$

and response data  $y_{ij}$  were drawn from a Multinomial distribution with probability equals to  $\mathbb{P}(Y_{ij})$ .

- (d) Time intensity parameters were generated as  $\boldsymbol{\beta}_{j_t \times 1} \sim \mathcal{N}(\beta_q^0 \mathbf{1}_{j_t}, \mathbf{I}_{j_t \times j_t})$ .
- (e) Response times were computed using the second component of the IRT-RTs model, which equals to the DIFF-based linear model:

$$\ln r_{ij} = \gamma_j + \omega_i + \text{DIFF}_{ij} \beta_j + \epsilon_{ij}$$

where  $\text{DIFF}_{ij} = \sum_{m=1}^{m_p} \mathbb{P}(Y_{ij} = m)^2$  and  $\epsilon_{ij} \sim \mathcal{N}(0, 0.25)$ , for all  $i = 1, \dots, i_h$  and  $j = 1, \dots, j_t$ .

- (f) The generated matrices of response data  $\mathbf{Y}_{i_h \times j_t}$  and times  $\mathbf{R}_{i_h \times j_t}$  were analysed using the fuzzy IRT-map. For both  $M = 3$  and  $M = 5$  cases, the sequential decision tree (see Figure 1a) was adopted. Since  $\boldsymbol{\alpha}$  and  $\boldsymbol{\eta}$  were simulated using the simplest model where latent traits and item parameters are invariant across nodes (e.g., see [83]), an IRTree with a common latent trait and common parameters was defined using the `IRTrees` R library [83]. The `glmmTMB` R package [97] was used to estimate the model parameters. Once estimates were obtained, fuzzy beta numbers were computed using the procedure described in Section 3.4, which yielded to two new matrices for the modes  $\mathbf{C}_{i_j \times j_t}$  and precisions  $\mathbf{S}_{i_j \times j_t}$  of the fuzzy numbers.

*Measures.* For each condition of the study, we assessed whether the rating uncertainty, as recovered by the precision of the fuzzy set, predicted the response times. Thus, response times were dichotomized into fast responses ( $r_{ij} = 1$ ) and slow ( $r_{ij} = 0$ ) responses by an item median split [98]. Then, for each of the  $j_t$  item, a Binomial linear model with logit link was used to predict the Boolean vector  $\mathbf{r}_{i_h \times 1}^*$  as a function of the precision values  $\mathbf{s}_{i_h \times 1}$ . Finally, predictions of the generalized linear model  $\hat{\mathbf{r}}_{i_h \times 1}^*$  were compared against the observations  $\mathbf{r}_{i_h \times 1}^*$  and the average Area Under the Curve (AUC) index was computed as follows:

$$\text{AUC}_{\text{avg}} = \frac{1}{j_t} \sum_{j=1}^{j_t} \left( \frac{1}{B} \sum_{b=1}^B \text{AUC}(\mathbf{r}_b^*, \hat{\mathbf{r}}_b^*)_j \right)$$

It is expected that the closer the  $AUC_{\text{avg}}$  to one, the more accurate the precisions of the fuzzy numbers will resemble the response times.

*Results.* Table 1 shows the average AUC index as a function of the simulation condition. Overall,  $AUC_{\text{avg}}$  was greater than the threshold for a random classification ( $AUC_{\text{avg}} = 0.5$ ), which indicated that precisions of the fuzzy numbers predicted response times better than random chance. Predictions were more accurate for the cases with  $M = 5$  response categories and a larger number of items ( $J = 15$ ). The number of sample units did not affect the accuracy of prediction. As expected, the greater accuracy was obtained for the cases with stronger time intensity parameters  $\beta^0 = -20.5$ , a condition that occurs if the variation of response times is mainly due to the task (as measured by the DIFF term). By and large, these findings suggest that, if compared to a RTs-based model for decision uncertainty, fuzzy numbers appropriately encode the uncertainty component associated to the choice of the final response in rating scales.

		$\beta = -10.5$		$\beta = -20.5$	
		$J = 5$	$J = 15$	$J = 5$	$J = 15$
$M = 3$	$I = 50$	0.666 (0.059)	0.77 (0.052)	0.694 (0.062)	0.812 (0.054)
	$I = 150$	0.649 (0.036)	0.744 (0.031)	0.683 (0.036)	0.803 (0.032)
	$I = 500$	0.662 (0.019)	0.755 (0.017)	0.697 (0.02)	0.813 (0.017)
$M = 5$	$I = 50$	0.754 (0.059)	0.776 (0.059)	0.772 (0.061)	0.812 (0.061)
	$I = 150$	0.736 (0.035)	0.746 (0.034)	0.765 (0.035)	0.792 (0.035)
	$I = 500$	0.736 (0.022)	0.767 (0.018)	0.766 (0.022)	0.808 (0.019)

Table 1: Simulation study: Average AUC index and its standard deviation (in parenthesis) over the  $B = 1000$  samples as a function of the simulation conditions.

## 5 Applications

In this section we illustrate the features of the proposed approach using four applications to real data. In particular, the first two are based on a controlled scenario in which varying levels of decision uncertainty were experimentally controlled. These two studies offer a way to assess the empirical effectiveness of the fuzzy IRT-map in retrieving decision uncertainty from standard rating data. Instead, the last two studies explore the differences between the proposed fuzzy-IRTree approach and two alternative methods for fuzzy ratings, namely the computerized Fuzzy Rating Scale (FRS) [30] and the Dynamic Fuzzy Rating Scale (DYFRAT) [21].



## 5.1 Case study 1: Rating data under experimental faking condition

The effects of faking behaviors on rating data have been widely studied in the area of psychometrics (e.g., see [4, 5, 99, 100]). Faking is defined as a deliberate behavior through which respondents distort their responses towards ones they consider more favorable in order to give overly positive self-descriptions, to dissimulate vocational interests, to simulate physical or psychological symptoms as a way to obtain rewards, or to have access to advantageous work positions [101]. In all these cases, faking acts as a kind of systematic error which alters the unfolding mechanism of the response process. For instance, in the case of faking-good or faking-bad response styles (i.e., the tendency to use higher or lower response categories in rating procedures, respectively), this results in reducing the overall response variability and increasing the number of stereotype answers. Because of its characteristics, faking can serve as a good candidate for studying uncertainty in rating process. In this application, we resorted to use rating data which were collected under honest and instructed faking-good measurement conditions. The aim is assessing to what extent our approach is sensitive enough to detect variations in decision uncertainty as arise from honest as opposed to faking response patterns. In particular, we expect to observe decreasing levels of decision uncertainty as responses patterns varies from honest to faking-good condition.

*Data and measures.* Data were originally collected and analysed by [5, 102] and refer to a sample of  $n = 484$  undergraduate students (79% females, ages ranged from 18 to 48, with mean age of 20.61 and standard deviation of 2.69) at the University of Padua (Italy). They were administered a personality questionnaire, the Perceived Empathic Self-Efficacy Scale (AEP/A) [103], with items scored on a 5-point scale where 1 denotes that she/he “Cannot do at all” and 5 denotes that she/he “Certain can do” the behavior described by the item. The questionnaire was administered using a paper and pencil format. Participants were randomly assigned to two groups, one ( $n = 237$ ) receiving the instruction to answer the questionnaire items as honest as possible (no faking condition), and the other ( $n = 247$ ) receiving the instruction to answer using a faking good response style. Faking-good was induced by letting participants know that a recruitment company was interested in hiring candidates for a very appealing job position and the questionnaire would have been used as a first method of selection. Following the rational described in [5], for the current analyses we retained a subset of four items only, which guarantee representativeness of the complete item pool, a good factorial structure, and a clear difference between the two groups in response frequencies.<sup>1</sup>

*Data analyses and results.* Table 2 shows the observed frequencies for the four items in the honest (H) and faking (F) conditions as well as the mean response value computed over the five categories. As

---

<sup>1</sup>The items were as follows: Q1. *When you meet new friends, find out quickly the things they like and those they do not like?* Q2. *Recognize if a person is seriously annoyed with you?* Q3. *Understand the state of mind of others when you are very involved in a discussion?* Q4. *Understand when a friend needs your help, even if he/she doesn't overtly ask for it?*

expected, items in the faking condition showed increased frequencies of response categories associated to positive responses (i.e.,  $Y \in \{4, 5\}$ ) as compared to items in the honest condition. A typical IRTree model for 5-point rating scales was defined and adapted to both groups (see Figure 1c). In this case, the decision structure was defined using three nodes, which represent the rating situation where answer using extreme points of the scale ( $Y \in \{1, 2, 4, 5\}$ ) is contrasted to the uncertain response category ( $Y = 3$ ) [17]. Thus, the IRTree model implied four item parameters and three latent traits, with the last trait being the same for lower and higher extreme responses. The model structure was defined using the `IRTrees` R library whereas item and person parameters were estimated via marginal maximum likelihood as implemented in the `glmmTMB` R package [97]. Overall, model fits showed good accuracy in terms of observed as opposed to predicted missclassification error ( $AUC_H = 0.75$ ,  $AUC_F = 0.79$ ). Tables 3-4 show the estimated model parameters for both honest and faking conditions. As expected, the probability to activate the right-branch of the nodes increased in the faking condition, especially for nodes 1 and 2. Similarly, latent traits were more strongly correlated in the faking condition as opposed to the honest condition. Once model’s parameters have been estimated, fuzzy beta numbers for both honest and faking groups were computed using the procedure given in Section 3.4. Thus, for each of the four items, we obtained  $n = 484$  fuzzy numbers expressed in terms of mode ( $m$ ) and precision ( $s$ ). Figure 3 shows an exemplary set of reconstructed fuzzy numbers. In order to compare honest and faking conditions with regards to the decision uncertainty as recovered by fuzzy beta numbers, in addition to mode ( $m$ ) and precision ( $s$ ) we computed fuzzy cardinality  $|\tilde{A}| = \int_{A_0} \xi_{\tilde{A}}(y) dy$  and fuzzy centroid  $\bar{A} = \frac{1+sm}{2+s}$  as well. Figure 4 shows the distribution of these measures for both the experimental conditions. As expected, fuzzy numbers in the faking condition showed higher precision and smaller cardinality as compared to the honest case. Similarly, modes and centroids increased in the faking condition which is in agreement with the previous results on faking experiments [5, 102]. Overall, the reconstructed fuzzy numbers behave according to the faking-good manipulation, which implied a reduction of the rating uncertainty and the choice of high rating scores. This was reflected by a highly increase in precision ( $s$ ) as well as a decrease in the size of fuzzy sets (fuzzy cardinality).

## 5.2 Case study 2: Rating data in moral dilemma scenarios

Moral dilemmas are emotionally salient scenarios in which an agent ought to adopt one of two mutually exclusive alternatives that differ in terms of violation of essential moral principles. Typical moral dilemmas include, for instance, the choice between letting one person die when that is necessary to saving five others (*footbridge*), the choice of smothering the supposedly incurable patient with a pillow in order to get the patient’s life insurance (*smother for dollars*), the choice of handing over one of two children to a doctor for painful experiments (*Sophie’s choice*), the choice of killing an healthy man to transplant his organs and saving five other patients (*transplant*) [104]. In all these cases, the choice between the lesser of two evils involves a tangled web of cognitive and emotional reactions that result in high levels of decision uncertainty. Because of these characteristics,

	$Y = 1$	$Y = 2$	$Y = 3$	$Y = 4$	$Y = 5$	mean response
item1 (H)	0.00	0.08	0.58	0.32	0.02	3.27
item1 (F)	0.00	0.03	0.48	0.44	0.04	3.50
item2 (H)	0.00	0.05	0.30	0.51	0.14	3.75
item2 (F)	0.00	0.04	0.22	0.54	0.20	3.90
item3 (H)	0.02	0.20	0.39	0.33	0.06	3.22
item3 (F)	0.01	0.13	0.36	0.38	0.11	3.45
item4 (H)	0.00	0.04	0.22	0.60	0.14	3.84
item4 (F)	0.00	0.01	0.20	0.52	0.27	4.04

Table 2: Case study 1: Observed frequency tables as a function of item number and type of group (H: honest group; F: faking group).

		node 1		node 2		node 3	
		$\hat{\theta}$	$\sigma_{\hat{\theta}}$	$\hat{\theta}$	$\sigma_{\hat{\theta}}$	$\hat{\theta}$	$\sigma_{\hat{\theta}}$
H	$\alpha_1$	-0.32	0.14	1.99	0.35	-1.60	0.29
	$\alpha_2$	0.85	0.15	3.39	0.42	-1.25	0.22
	$\alpha_3$	0.47	0.14	0.87	0.24	-0.43	0.21
	$\alpha_4$	1.28	0.16	3.69	0.44	-1.52	0.22
F	$\alpha_1$	0.08	0.13	9.19	1.82	-2.12	0.30
	$\alpha_2$	1.30	0.16	10.06	1.80	-1.02	0.19
	$\alpha_3$	0.58	0.14	5.31	1.13	-0.68	0.20
	$\alpha_4$	1.46	0.17	12.55	2.27	-0.77	0.18

Table 3: Case study 1: Estimates ( $\hat{\theta}$ ) and standard errors ( $\sigma_{\hat{\theta}}$ ) for item parameters in the honest (H) and faking (F) conditions.

		$\eta_1$	$\eta_2$	$\eta_3$	$\hat{\sigma}_{\eta}$
H	$\eta_1$	1.00			0.35
	$\eta_2$	-0.27	1.00		1.37
	$\eta_3$	0.38	-0.99	1.00	1.06
F	$\eta_1$	1.00			0.44
	$\eta_2$	0.58	1.00		7.59
	$\eta_3$	0.56	-0.35	1.00	0.90

Table 4: Case study 1: Estimated correlation matrix and standard deviations ( $\hat{\sigma}_{\eta}$ ) for latent traits in the honest (H) and faking (F) conditions.

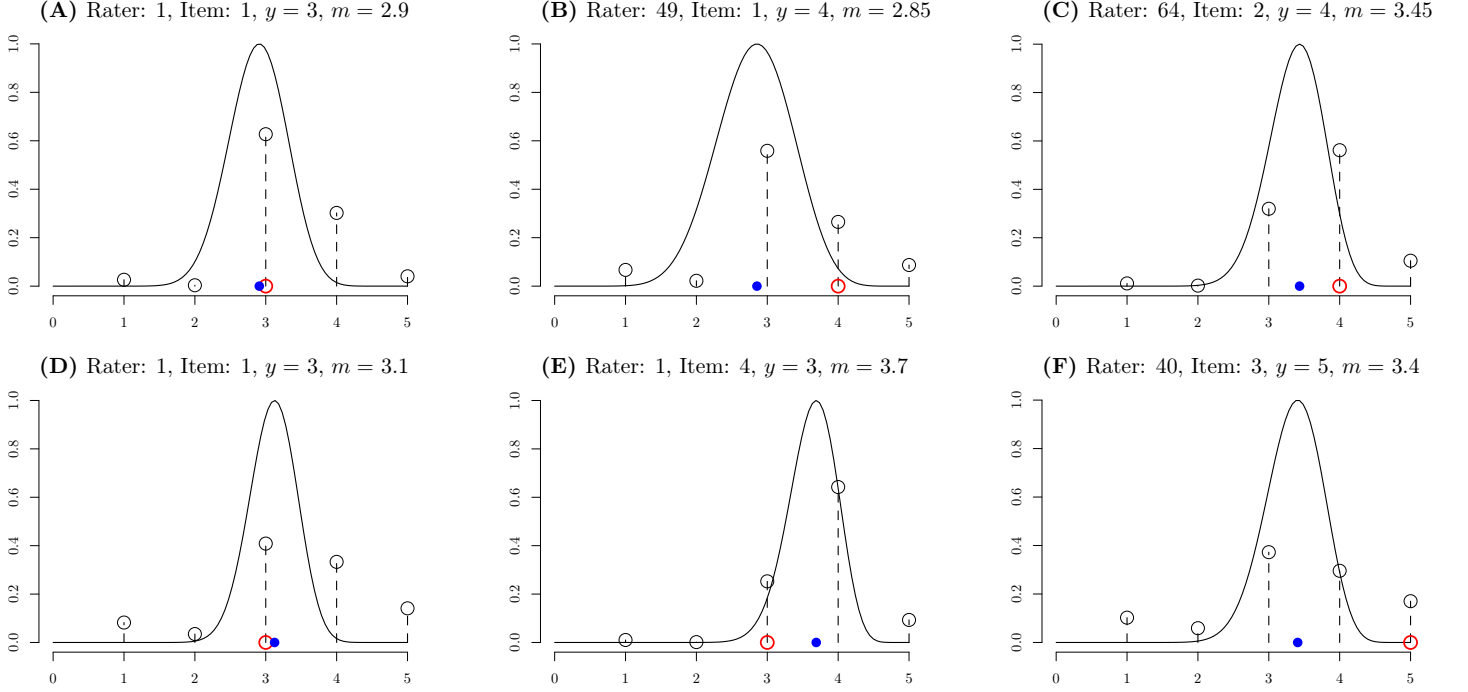


Figure 3. Case study 1: Fuzzy beta numbers (black curves) and estimated probabilities of response categories (black dashed vertical lines) for some raters of the honest (A-C panels) and faking groups (D-F panels). Note that probability masses and fuzzy sets are overlapped over the same domain  $\Omega(y)$ , red and blue circles represent observed response ( $y$ ) and fuzzy mode ( $m$ ), respectively.

moral dilemmas can serve as a framework for studying how ratings behave as a function of decision uncertainty. In this application, we used two moral dilemmas - i.e. *footbridge* and *transplant* - and assessed how they impacted the intensity of raters' negative emotions towards the scenario's protagonist. In both dilemmas, the protagonist must choose between the sacrifice of one person (a stranger in the *footbridge* case, a victim's physician in *transplant*) in order to save a larger group. However, these scenarios differ because of an additional role conflict that results from the different method of killing [105]: while in *footbridge* the perpetrator is an anonymous pedestrian with no relationship to the victim, in *transplant* the perpetrator is a doctor with moral duties. As such, we expect a higher degree of uncertainty in assessing negative emotions for the *transplant* case as opposed to *footbridge*.

*Data and measures.* Data were originally collected by [105] in a large project assessing many aspects of moral decision making, including several cognitive scales and personality surveys. For the purposes of this study, we selected a subset of the entire dataset. The final sample consisted of  $n = 500$  participants (54% females, ages ranged from 18 to 58, with mean age of 25.06 and standard deviation

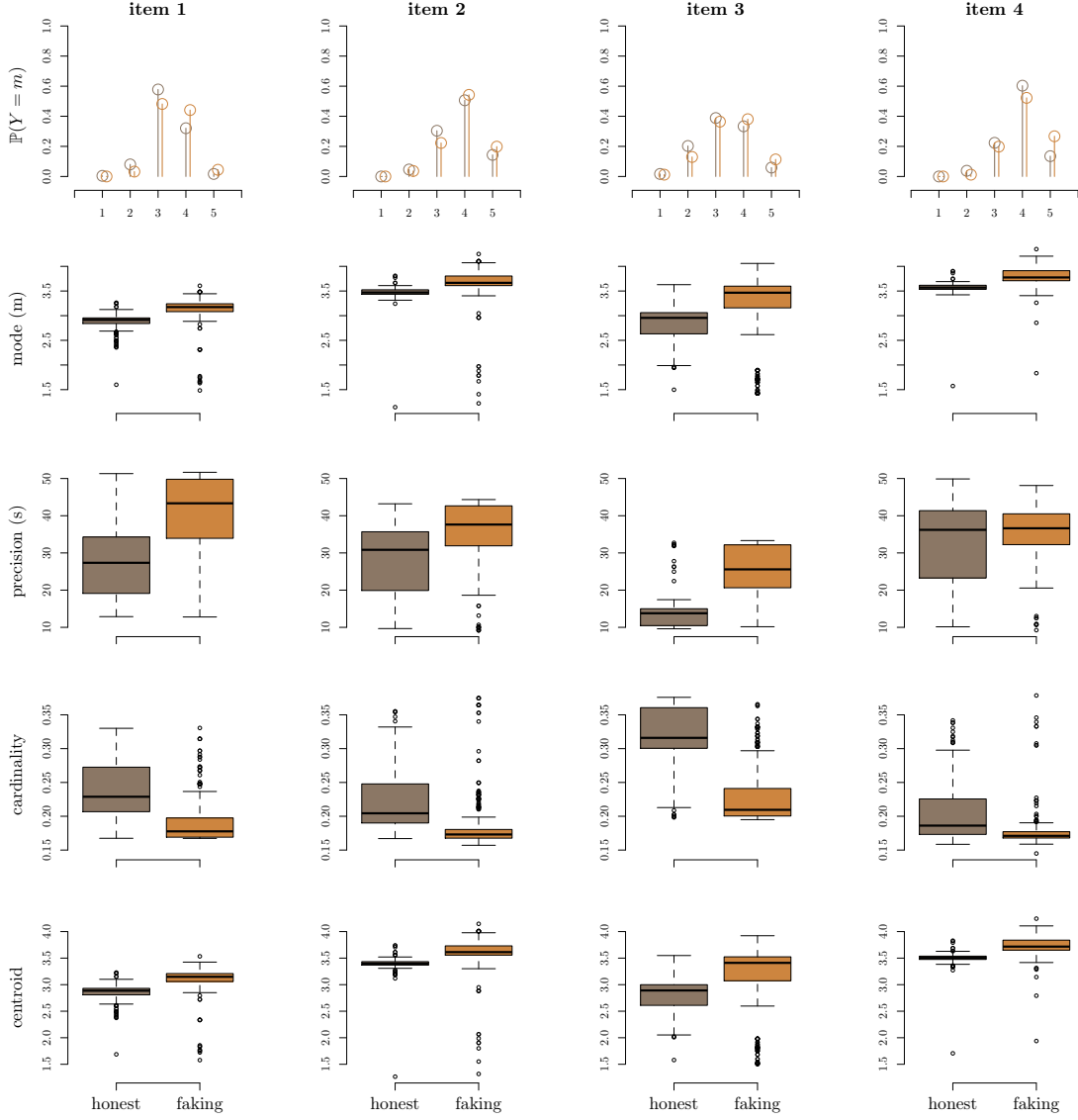


Figure 4. Case study 1: Distribution of summary statistics (mode, precision, cardinality, centroid) for fuzzy numbers computed for each participant in honest (in light brown color) and faking (in orange color) conditions. Note that plots in the first row show the observed frequencies as a function of the experimental conditions.

of 3.96), mainly composed of German speakers. They read both dilemma scenarios and rated the intensity of their negative emotions toward the scenario’s protagonist using a 5-point scale. A total of four emotional items was presented along with the question “When I think of the protagonist and his/her decision, I feel [disappointment, disgust, contempt, anger]”. The texts used for the moral

dilemmas were as follows:

*Footbridge.* A runaway trolley with malfunctioning breaks is heading down the tracks towards a group of five workmen. A pedestrian observes this from a footbridge. If nothing is done the trolley will overrun and kill the five workmen. The only way for the pedestrian to avoid the deaths of the five workmen is to push a large stranger who is standing next to him off the bridge onto the tracks below where his large body will stop the trolley but which will also kill the stranger. Outcome: The pedestrian decided to push the stranger off the bridge. Due to this decision the trolley was stopped and the five workmen were saved; but the stranger was killed.

*Transplant.* Five patients are treated in a hospital. Each of whom is in critical condition due to organ failing. A healthy man consults the head physician for routine checkup. If nothing is done the five patients will die due to a shortage of available transplants. The only way for the head physician to save the lives of the first five patients is to kill the healthy man (against his will) and to transplant his organs into the bodies of the other five patients. Outcome: The head physician decided to kill the healthy man and to transplant the organs. Due to this decision five patients were saved; but the healthy man was killed.

*Data analysis and results.* Two IRTree models with sequential structure (see Figure 1a) were separately defined and adapted to *footbridge* (F) and *transplant* (T) data. The IRT models required  $J = 4$  item parameters and  $N = 4$  number of nodes and latent traits. The model structure was defined using the `IRTrees` R library whereas model parameters were estimated via marginal maximum likelihood as implemented in the `glmmTMB` R package [97]. Tables 5-6 show the estimated model parameters for both footbridge and transplant scenarios. Once model's parameters have been estimated, fuzzy beta numbers were computed using the procedure given in Section 3.4. The final models showed a satisfactory fit ( $AUC_F = 0.89$ ,  $AUC_T = 0.86$ ). Finally, for each of the four items, we obtained  $n = 500$  fuzzy numbers expressed in terms of mode ( $m$ ) and precision ( $s$ ). Likewise for the first case study, also in this context *footbridge* and *transplant* were compared in terms of modes, precisions, fuzzy cardinalities, and fuzzy centroids. Figure 5 shows the distribution of these measures for both the moral scenarios. As expected, unlike for the footbridge scenario, ratings in *transplant* were characterized by higher levels of decision uncertainty. Overall, fuzzy numbers showed larger modes and centroids, precisions of the fuzzy sets were higher in median and more variable, and fuzzy cardinalities were smaller in median. Finally, Figure 6 shows a subset of estimated fuzzy beta numbers for both dilemma scenarios. We can observe how fuzzy sets for *transplant* showed larger support than fuzzy sets associated to *footbridge*. Interestingly, because of the different levels of decision uncertainty underlying rating responses, the estimated modes often differ from the observed final responses.

### 5.3 Case study 3: fuzzy-IRTree and Fuzzy Rating Scale (FRS) in modeling rating data

The computerized Fuzzy Rating Scale (FRS) is a direct rating method which allows raters to express their responses by using fuzzy sets (for further details, see Sect. ??). The main difference

		node 1		node 2		node 3		node 4	
		$\hat{\theta}$	$\sigma_{\hat{\theta}}$	$\hat{\theta}$	$\sigma_{\hat{\theta}}$	$\hat{\theta}$	$\sigma_{\hat{\theta}}$	$\hat{\theta}$	$\sigma_{\hat{\theta}}$
F	$\alpha_1$	8.75	0.59	0.46	0.17	-1.04	0.26	-9.65	1.68
	$\alpha_2$	7.81	0.54	0.16	0.18	-0.16	0.27	-9.49	1.69
	$\alpha_3$	7.71	0.53	0.08	0.18	-0.86	0.28	-9.34	1.68
	$\alpha_4$	9.32	0.61	0.76	0.18	0.08	0.26	-8.89	1.58
T	$\alpha_1$	6.45	0.61	3.87	0.50	3.38	0.41	-1.97	0.31
	$\alpha_2$	6.41	0.61	4.16	0.51	3.81	0.46	-1.39	0.30
	$\alpha_3$	8.42	0.74	4.59	0.55	3.76	0.44	-0.99	0.28
	$\alpha_4$	9.25	0.79	5.38	0.61	4.64	0.51	-0.54	0.27

Table 5: Case study 2: Estimates ( $\hat{\theta}$ ) and standard errors ( $\sigma_{\hat{\theta}}$ ) of item parameters in the footbridge (F) and transplant (T) scenarios.

		$\eta_1$	$\eta_2$	$\eta_3$	$\eta_4$	$\hat{\sigma}_{\eta}$
F	$\eta_1$	1.00				9.49
	$\eta_2$	0.06	1.00			2.23
	$\eta_3$	0.06	0.81	1.00		2.65
	$\eta_4$	0.32	0.39	0.58	1.00	6.14
T	$\eta_1$	1.00				6.26
	$\eta_2$	0.14	1.00			3.09
	$\eta_3$	0.05	0.62	1.00		2.72
	$\eta_4$	0.21	0.60	0.66	1.00	4.13

Table 6: Case study 2: Estimated correlation matrix and standard deviations ( $\hat{\sigma}_{\eta}$ ) for latent traits in the footbridge (F) and transplant (T) scenarios.

between FRS and the proposed fuzzy-IRTree method is that FRS is based on a direct elicitation of the respondent’s fuzziness for each item being rated. By contrast, fuzzy-IRTree is an indirect rating method and computes the fuzziness of a rating response using a psychometric model (IRTree), which in turn formalizes the response process underlying the observed rating response. Thus, in the first case, the fuzziness of a rating response reflects to what extent the rater is uncertain about his/her response to a specific question, or rather the degree of confidence he/she has in the final response [30]. Instead, in the second case, the fuzziness of a rating response reflects the rater’s decision uncertainty as resulting from the conflicting demands of the opinion formation stage, which comes before the responding stage. As a result, the fuzzy-IRTree method computes the fuzziness of a specific item  $j$  as a function of the entire rater’s response pattern  $\mathbf{y}_{i,J \times 1}$  (by means of the estimated IRTree parameters), instead of being computed on the  $j$ -th item only. Hence, with regards to the fuzzy

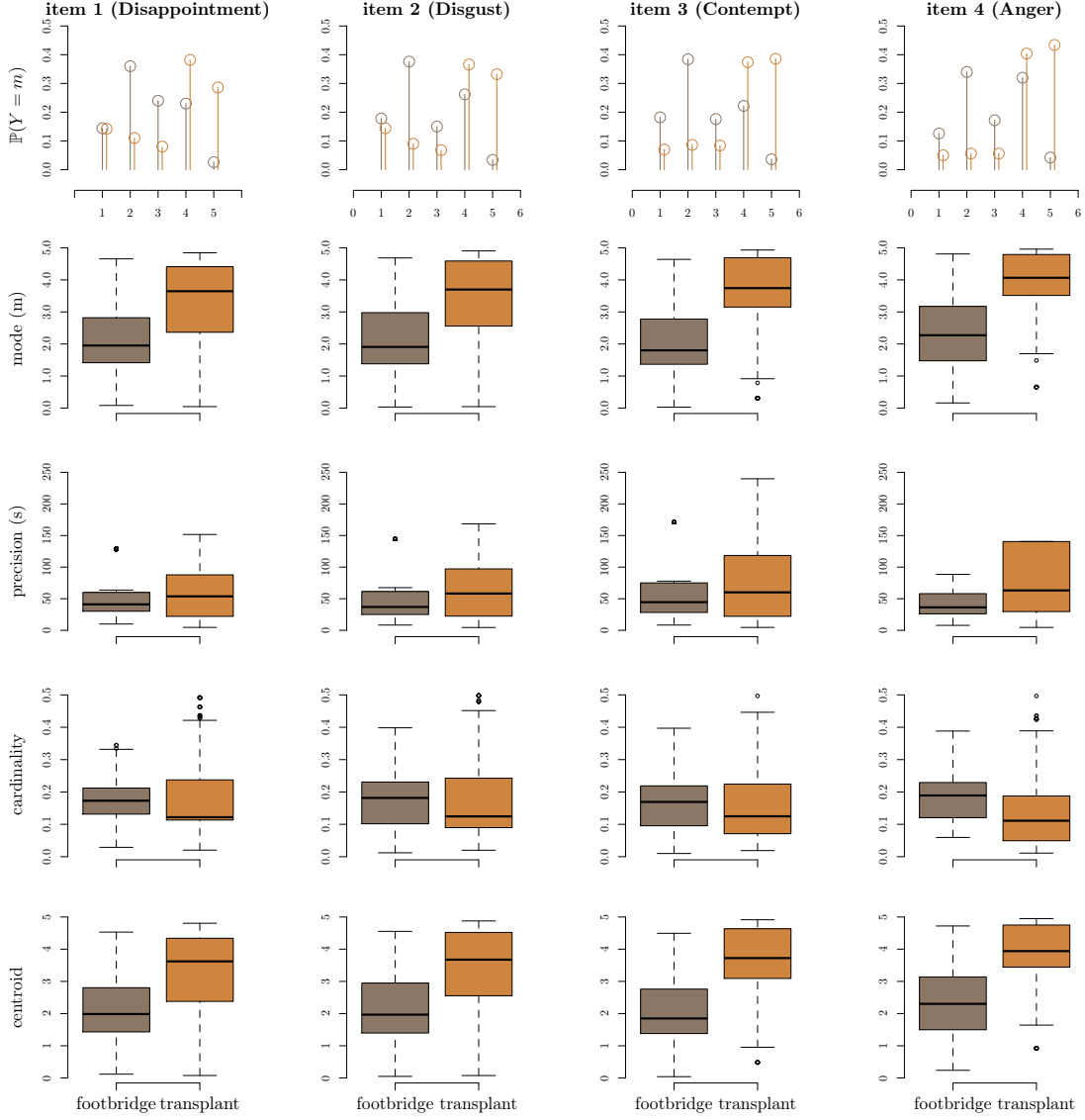


Figure 5. Case study 2: Distribution of summary statistics (mode, precision, cardinality, centroid) for fuzzy numbers computed for each participant in footbridge (in light brown color) and transplant (in orange color) scenarios. Note that plots in the first row show the observed frequencies as a function of the dilemma scenarios.

sets produced by the two methods, we expect no differences in terms of modes (as they reflect the final rating responses) and substantial differences in terms of fuzzy cardinalities (as they reflect the fuzziness of the final rating responses). In particular, we expect that fuzzy-IRTree produces larger fuzziness as it models the entire response pattern  $\mathbf{y}_{i_{J \times 1}}$ .



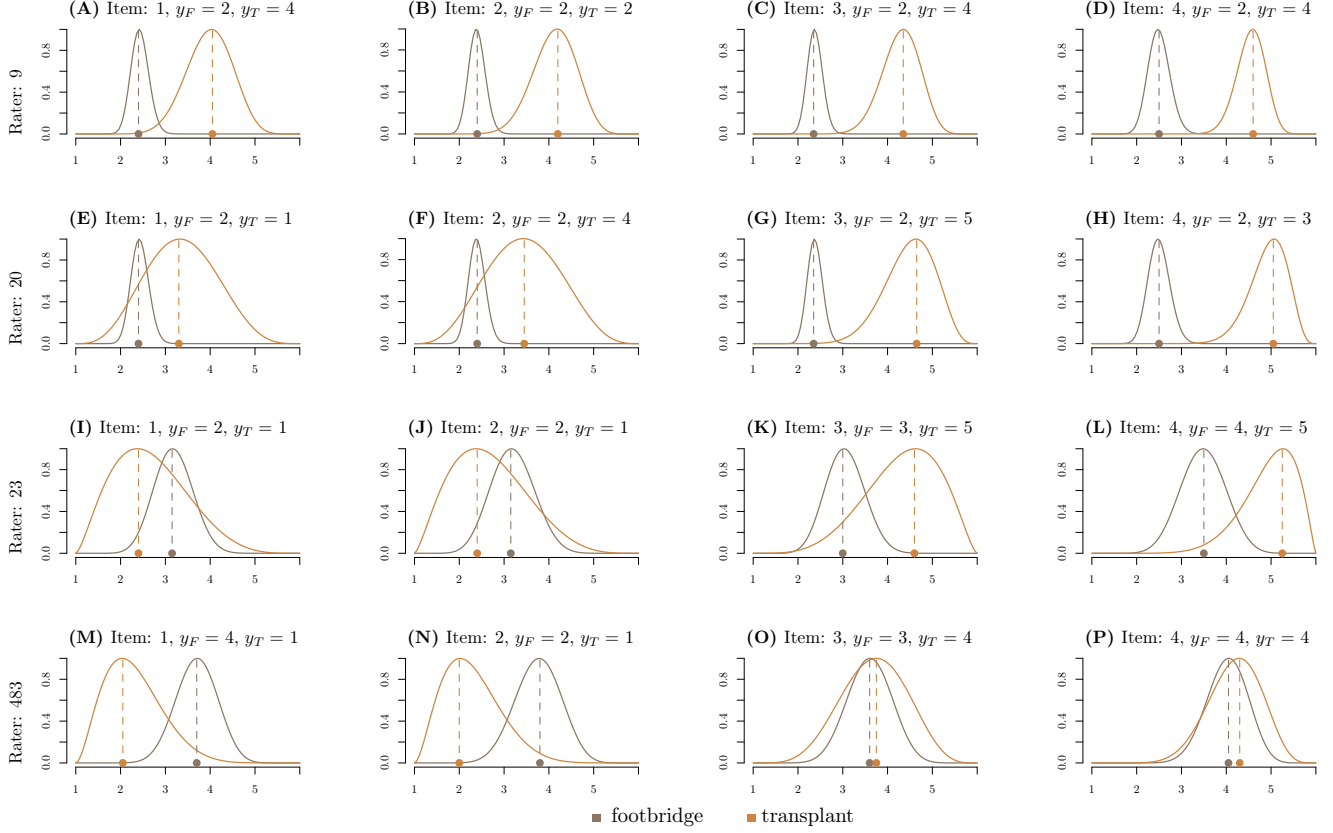


Figure 6. Case study 2: Fuzzy beta numbers for some raters of the footbridge (light brown curves) and transplant (orange curves) scenarios along with their estimated modes (filled circles). Note that  $y_T$  and  $y_F$  indicate the observed crisp responses for *footbridge* and *transplant*, respectively.

*Data and measures.* Data were originally collected by [22] and refer to a survey of  $J = 13$  items administered to a sample of  $n = 70$  raters about restaurant and service quality. The respondents provided their responses using two different rating scales, namely a crisp Likert-type scale with  $M = 5$  levels (from 1: “Strongly Disagree” to 5: “Strongly Agree”) and the computerized Fuzzy Rating Scale (FRS). Thus, the final dataset consisted of crisp Likert-type responses as well as trapezoidal fuzzy responses.

*Data analysis and results.* The fuzzy-IRTtree method was applied on the  $n \times J$  dataset of Likert-type responses. To this end, the simplest IRTtree model for 5-point scales was defined and adapted to the observed data (see Figure 1b). The model implied  $J = 13$  item parameters and  $N = 4$  nodes with a single latent trait. The model structure was defined using the `IRTrees` R library whereas

model parameters were estimated via marginal maximum likelihood as implemented in the `glmmTMB` R package [97]. The estimated model showed a satisfactory fit ( $\text{AUC} = 0.71$ ). Finally, for each of the thirteen items,  $n = 70$  beta fuzzy numbers were obtained. To adequately compare fuzzy-IRTree and FRS, the fuzzy sets produced by the two methods were linearly rescaled in  $[0, 1]$ . Next, they were summarized in terms of modes (i.e.,  $m$  for beta fuzzy sets,  $(m_1 + m_2)/2$  for trapezoidal fuzzy sets), support lengths (i.e.,  $\max(A_0) - \min(A_0)$  for beta fuzzy numbers,  $ub - lb$  for trapezoidal fuzzy sets), and fuzzy cardinalities (i.e.,  $|\tilde{A}| = \int_{A_0} \xi_{\tilde{A}}(y) dy$ ). Finally, since for each rater  $J \times 2 = 26$  fuzzy sets were available (i.e.,  $J$  items for fuzzy-IRTree and  $J$  items for FRS), summary measures were averaged across items and comparisons were made over the  $n$  independent raters. Figure 7 shows the distributions of these measures for both methods. As expected, fuzzy sets computed through fuzzy-IRTree showed larger cardinalities and wider supports as opposed to fuzzy sets computed via FRS whereas no differences can be seen for modes. Moreover, the distribution of these measures were less variable for fuzzy-IRTree. This is potentially due to the fact that beta fuzzy sets were computed as a function of the estimated parameters of the IRTree statistical model (i.e., they are computed using denoised observed data). To evaluate whether the sample differences were statistically significant, three Beta linear models were run with the type fuzzy rating method being used as categorical predictor (note that Beta linear models were chosen because of the distribution characteristics of the involved outcome variables) [106]. In particular, with regards to the modes of fuzzy sets there was no statistically significant difference between the two methods ( $\hat{\beta}_{\text{type:FRS}} = 0.085$ ,  $\hat{\sigma}_{\hat{\beta}_{\text{type:FRS}}} = 0.086$ ,  $z_{\hat{\beta}_{\text{type:FRS}}} = 0.988$ ,  $\alpha = 0.05$ ). On the contrary, cardinalities of fuzzy sets were statistically different for both methods ( $\hat{\beta}_{\text{type:FRS}} = -0.243$ ,  $\hat{\sigma}_{\hat{\beta}_{\text{type:FRS}}} = 0.064$ ,  $z_{\hat{\beta}_{\text{type:FRS}}} = -3.787$ ,  $\alpha = 0.05$ ). Similarly, support lengths differed for both methods ( $\hat{\beta}_{\text{type:FRS}} = -0.515$ ,  $\hat{\sigma}_{\hat{\beta}_{\text{type:FRS}}} = 0.075$ ,  $z_{\hat{\beta}_{\text{type:FRS}}} = -6.790$ ,  $\alpha = 0.05$ ). Overall, the results suggest that the assessment of restaurant and service quality involved a higher level of fuzziness, in particular that referring to the decision uncertainty component which has been quantified through the fIRTree method. This would indicate the presence of a particular pattern of uncertainty in selecting the final response across all the  $J = 8$  items being considered.

#### 5.4 Case study 4: fuzzy-IRTree and Dynamic Fuzzy Rating Scale (DYFRAT) in modeling rating data

The Dynamic Fuzzy Rating Scale (DYFRAT) is an indirect method which computes fuzziness using implicit biometric measures such as hand movements and response times (for further details, see Sect. ??). Although both fuzzy-IRTree and DYFRAT are indirect fuzzy rating methods, DYFRAT is more similar to FRS in the way it computes respondent's fuzziness: it does not use a statistical model to represent the respondent's rating process and rater's fuzziness is based on an item-level analysis. Hence, likewise for the previous case study, we expect to observe no differences in terms of modes of the final fuzzy sets and larger fuzziness for those fuzzy sets produced by fuzzy-IRTree.

*Data and measures.* Data refer to a survey of  $J = 8$  items which were administered to a sample

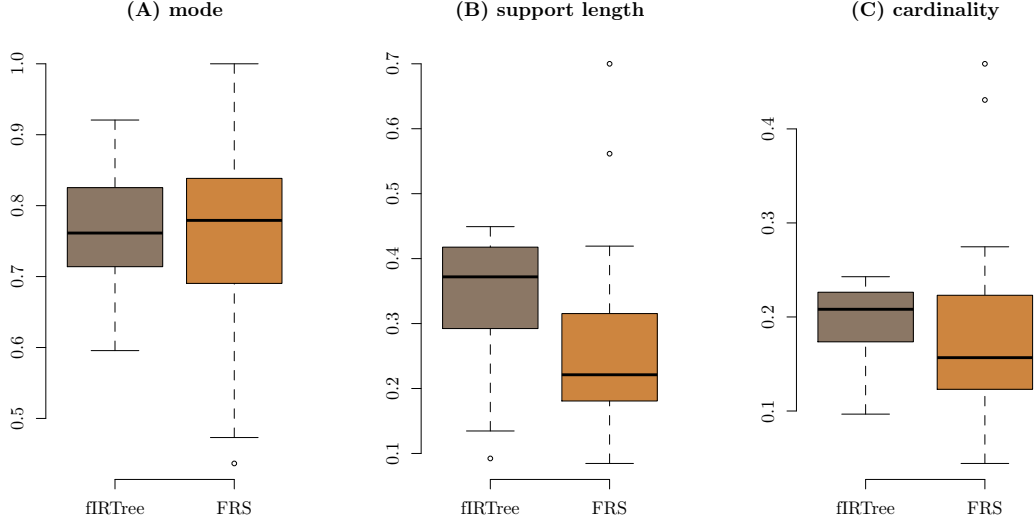


Figure 7. Case study 3: Distribution of summary statistics (mode, support length, cardinality) for beta fuzzy numbers (fIRTree) and trapezoidal fuzzy numbers (FRS) computed across items for fuzzy-IRTree (in light brown color) and Fuzzy Rating Scale (in orange color) methods.

of  $n = 72$  young drivers in Trentino region (north-est of Italy). The items were part of the short version of the Driving Anger Scale (DAS) [107], used to assess driving anger provoked by someone else's behaviors like slow driving and discourtesy. The items were administered using DYFRAT [21] on a pseudo-circular rating scale with  $M = 5$  levels. In this case, fuzzy responses were represented using beta fuzzy numbers. For the sake of comparison, crisp Likert-type responses were collected as for the FRS method. Thus, the final dataset consisted of crisp Likert-type responses as well as beta fuzzy responses.

*Data analysis and results.* The fuzzy-IRTree method was applied on the Likert-type dataset. As for the previous case, the simplest IRTree model for 5-point scales was defined and adapted to the observed data (see Figure 1b). The model implied  $J = 8$  item parameters and  $N = 4$  nodes with a single latent trait. The `IRTrees` and `glmmTMB` R libraries were used for model definition and parameters estimation [97]. The estimated model showed a satisfactory fit ( $AUC = 0.72$ ). Next, for each of the eight items,  $n = 72$  beta fuzzy numbers were obtained. Finally, to adequately compare fuzzy-IRTree and DYFRAT, fuzzy sets produced by the two methods were linearly rescaled in  $[0, 1]$ . Three measures were used in order to summarize fuzzy sets: modes (i.e.,  $m$ ), support lengths (i.e.,  $\max(A_0) - \min(A_0)$ ), and fuzzy cardinalities (i.e.,  $|\tilde{A}| = \int_{A_0} \xi_{\tilde{A}}(y) dy$ ). They were averaged across items so that comparisons were made over  $n$  independent raters. Figure 8 shows

the distributions of these measures for both methods. The results were in the expected directions, namely the two methods showed differences in terms of cardinalities and support lengths, with fuzzy-IRTree based measures being less variable (this is potentially due to the fact the fuzzy-IRTree works on denoised data). In order to evaluate these results statistically, three Beta linear models were run with the type fuzzy rating method being used as categorical predictor [106]. In particular, the modes of fuzzy sets were no statistically significant across methods ( $\hat{\beta}_{\text{type:DYFRAT}} = -0.053$ ,  $\hat{\sigma}_{\hat{\beta}_{\text{type:DYFRAT}}} = 0.077$ ,  $z_{\hat{\beta}_{\text{type:DYFRAT}}} = -0.686$ ,  $\alpha = 0.05$ ). On the contrary, both cardinalities ( $\hat{\beta}_{\text{type:DYFRAT}} = -0.625$ ,  $\hat{\sigma}_{\hat{\beta}_{\text{type:DYFRAT}}} = 0.044$ ,  $z_{\hat{\beta}_{\text{type:DYFRAT}}} = -14.170$ ,  $\alpha = 0.05$ ) and support lengths ( $\hat{\beta}_{\text{type:DYFRAT}} = -1.247$ ,  $\hat{\sigma}_{\hat{\beta}_{\text{type:DYFRAT}}} = 0.056$ ,  $z_{\hat{\beta}_{\text{type:DYFRAT}}} = -22.140$ ,  $\alpha = 0.05$ ) differed across methods. Overall, the results indicate that self-assessing the anger provoked by driving behavior induced a certain level of fuzziness. However, this was not completely represented by the patterns of hesitation to provide the final rating response - i.e., that component of fuzziness quantified by the DYFRAT method. Rather, fuzziness was mainly due to raters' decision uncertainty, namely that component of fuzziness which emerges as a stable response style across all the items being assessed.

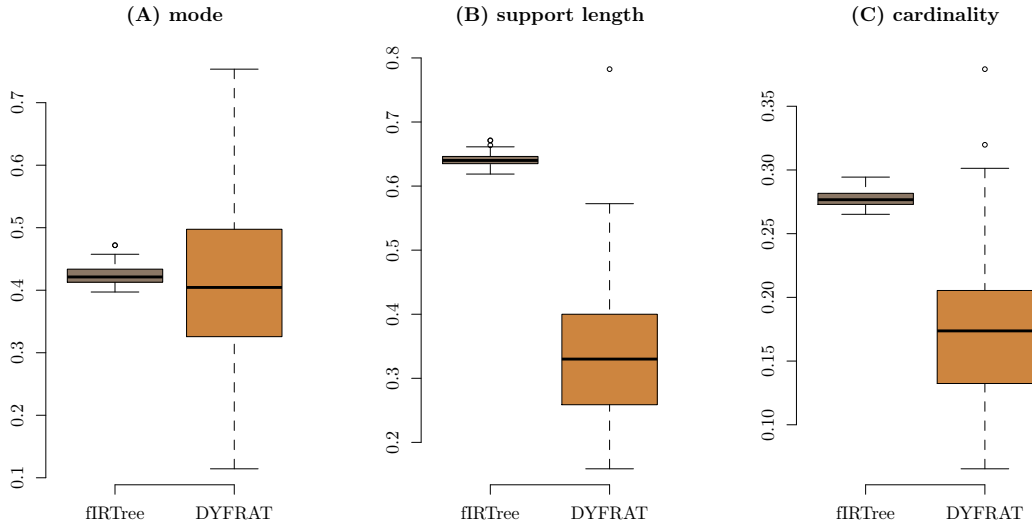


Figure 8. Case study 4: Distribution of summary statistics (mode, support length, cardinality) for beta fuzzy numbers computed across items for fuzzy-IRTree (in light brown color) and Dynamic Fuzzy Rating Scale (in orange color) methods.

## 6 Discussion and conclusions

In this paper we described a novel procedure to represent rating responses in terms of fuzzy numbers. Similarly for other types of fuzzy conversion scales, our approach followed a two-step process by means of which fuzzy numbers are computed based on a previously estimated psychometric model for rating data. To this end, Item Response Theory-based trees (IRTrees) have been used, which provide a formal representation of the stage-wise cognitive process of answering survey questions [80, 108]. Unlike traditional IRT models, IRTrees allows for a flexible modeling of the rating response where item contents relate to latent traits by means of a priori specified response styles, which include decision nodes for the tendency to choose moderate as opposed to extreme response categories as well as for the tendency to agree versus disagree with a given item content. As a consequence, fuzziness of rating responses has been recovered from the characteristics of rater’s response pattern  $\mathbf{y}_i$ , instead of being computed as a byproduct of the item-based direct rating. This offered a coherent meaning system in which fuzzy responses  $\tilde{\mathbf{y}}_i$  can be interpreted in terms of decision uncertainty that characterized the rater’s response process. To this end, although other type of fuzzy sets have been suggested for rating data (e.g., triangular, trapezoidal [30, 74]), we resorted to adopt two-parameter fuzzy beta numbers since beta-like models have been proved to adequately represent the characteristics of asymmetry of bounded rating data [91]. Simulation and real case studies were adopted to evaluate the characteristics and properties of our proposal. In particular, the simulation study was designed in order to provide converging results about the effectiveness of our proposal to recover decision rating uncertainty. To this purpose, a controlled scenario was used and our model was contrasted against a standard IRT-RTs model which uses response times (RTs) to quantify decision uncertainty in rating responses [94]. The results showed the ability of the fuzzy IRT-map to detect decision uncertainty when it is present in rating data. This was also confirmed by the results of the first two case studies, which involved two empirical situations characterized by ratings under uncertainty. Two additional case studies were also used to highlight the differences of the proposed method in relation with two existing methods, namely the computerized Fuzzy Rating Scale (FRS) and the Dynamic Fuzzy Rating Scale (DYFRAT). The results showed that fuzzy-IRTree recovers fuzziness of rating responses differently from standard fuzzy rating methods. In particular, when compared to FRS and DYFRAT, the proposed method produced less variable fuzzy responses, with fuzzy sets having a higher degree of fuzziness. By and large, this difference can be explained in light of three characteristics that make fuzzy-IRTree different from existing fuzzy rating methods: (i) It represents fuzziness in terms of the rater’s decision uncertainty which results from the conflicting demands of the opinion formation stage instead of the conflict provoked by the final response stage; (ii) It is grounded on a model-based approach which uses the IRTree model as a formal representation of the rater’s response process, with the consequence that fuzziness is computed as a function of the entire rater’s response pattern  $\mathbf{y}_i$ ; (iii) It uses a statistical model (i.e., IRTree) which acts by denoising the observed data in advance, with the consequence that fuzzy sets are computed using the estimated IRTree parameters  $\hat{\boldsymbol{\alpha}}$  and  $\hat{\boldsymbol{\eta}}$  instead of being derived from the observed data directly.

Hence, the fuzzy-IRTree method might be of particular utility when data analysts need to quantify the uncertainty that arises through the rater’s decision process on the whole, from the opinion formation stage to the choice of the final response. Instead, other fuzzy rating methods such as FRS and DYFRAT might be used in all those cases whether the interest concerns the quantification of the degrees of hesitation or uncertainty at a single-item level, regardless of the other items or questions. As a result, in the first case fuzziness is calculated by considering the rater’s response style to a set of questions or items whereas in the second case fuzziness is quantified as discrepancy between the final response and the other competing responses for a given item.

Some advantages of the proposed fuzzy IRT-map are as follows. First, since the procedure does not require a dedicated measurement setting, it is applicable over a wide range of survey situations including different rating formats (e.g., Likert-type, forced-choice, funnel response-format) [81, 109]. Second, it avoids using direct rating scales which can often provide distorted responses because of cognitive biases underling numerical and intensity estimation [110]. Third, it uses a flexible psychometric model to represent the cognitive stages of the response process, which can each time be adapted by researchers to model specific rating situations. Moreover, the use of a statistical model as a first processing step allows fuzzy responses to be computed on a kind of denoised data.

However, as for other statistical-based fuzzy quantification procedure, also the proposed fuzzy IRT-map can potentially suffer from some limitations. For instance, as it is based on a psychometric model for rating responses, sample size or the number of items should be large enough to provide reliable results for the estimates  $\hat{\theta} = \{\hat{\alpha}, \hat{\eta}\}$  [111, 112]. In addition, the hypothesized IRTree rating model should also be valid for the sample being analyzed. For instance, in empirical cases for which a rating model cannot be determined in advance, it may be advised to define and test several IRTrees, the best of which can be chosen by means of minimum Akaike Information Criteria (AIC) [17]. Similarly, for studies involving huge samples, MCMC based algorithms should be preferred to estimate IRTrees over standard marginal maximum likelihood-based algorithms [113]. To this end, several methods and implementations are available nowadays (e.g., see [114]).

Our proposal may be extended in several ways. For instance, IRTree models including response times in the computation of raters’ decision uncertainty [14] may also be adopted and generalized fuzzy numbers may be used accordingly [62, 115]. In conclusion, modeling uncertainty in rating data is a crucial task in all those research contexts involving human subjects as source of information such as social surveys, formative and teaching evaluation, decision support systems, quality control, psychological assessment, medical and health decision making, military promotion screening, etc. We believe that our proposal may offer an ecological but reliable procedure to address the problem of measuring subjective evaluations.

## References

- [1] Eunike Wetzel and Samuel Greiff. The world beyond rating scales. *European Journal of Psychological Assessment*, 34(1):1–5, 2018.
- [2] Adrian Furnham. Response bias, social desirability and dissimulation. *Personality and individual differences*, 7(3):385–400, 1986.
- [3] Adam W Meade and S Bartholomew Craig. Identifying careless responses in survey data. *Psychological methods*, 17(3):437, 2012.
- [4] Michael Eid and Michael J Zickar. Detecting response styles and faking in personality and organizational assessments by mixed rasch models. In *Multivariate and mixture distribution Rasch models*, pages 255–270. Springer, 2007.
- [5] Luigi Lombardi, Massimiliano Pastore, Massimo Nucci, and Andrea Bobbio. Sgr modeling of correlational effects in fake good self-report measures. *Methodology and Computing in Applied Probability*, 17(4):1037–1055, 2015.
- [6] Carolyn C Preston and Andrew M Colman. Optimal number of response categories in rating scales: reliability, validity, discriminating power, and respondent preferences. *Acta psychologica*, 104(1):1–15, 2000.
- [7] Jonathan Rabinowitz, Nina R Schooler, Brianne Brown, Mads Dalsgaard, Nina Engelhardt, Gretchen Friedberger, Bruce J Kinon, Daniel Lee, Felice Ockun, Atul Mahableshwarkar, et al. Consistency checks to improve measurement with the montgomery-asberg depression rating scale (madr). *Journal of affective disorders*, 256:143–147, 2019.
- [8] Timothy Johnson, Patrick Kulesa, Young Ik Cho, and Sharon Shavitt. The relation between culture and response styles: Evidence from 19 countries. *Journal of Cross-cultural psychology*, 36(2):264–277, 2005.
- [9] Philip J Rosenbaum and Jaan Valsiner. The un-making of a method: From rating scales to the study of psychological processes. *Theory & Psychology*, 21(1):47–65, 2011.
- [10] Ozlem Ozkok, Michael J. Zyphur, Adam P. Barsky, Max Theilacker, M. Brent Donnellan, and Frederick L. Oswald. Modeling measurement as a sequential process: Autoregressive confirmatory factor analysis (AR-CFA). *Front. Psychol.*, 10, sep 2019.
- [11] Boaz Shulruf, John Hattie, and Robyn Dixon. Factors affecting responses to likert type questionnaires: introduction of the impexp, a new comprehensive model. *Social Psychology of Education*, 11(1):59–78, 2008.
- [12] Roger Tourangeau, Lance J Rips, and Kenneth Rasinski. *The psychology of survey response*. Cambridge University Press, 2000.

- [13] Norbert Schwarz and Daphna Oyserman. Asking questions about behavior: Cognition, communication, and questionnaire construction. *The American Journal of Evaluation*, 22(2):127–160, 2001.
- [14] Pere J Ferrando and Urbano Lorenzo-Seva. A measurement model for likert responses that incorporates response time. *Multivariate Behavioral Research*, 42(4):675–706, 2007.
- [15] Kaiwen Man, Jeffery R. Harring, Yunbo Ouyang, and Sarah L. Thomas. Response time based nonparametric kullback-leibler divergence measure for detecting aberrant test-taking behavior. *International Journal of Testing*, 18(2):155–177, feb 2018.
- [16] John Zaller and Stanley Feldman. A simple theory of the survey response: Answering questions versus revealing preferences. *American journal of political science*, pages 579–616, 1992.
- [17] Paul De Boeck and Ivailo Partchev. IRTrees: Tree-based item response models of the GLMM family. *J. Stat. Soft.*, 48(Code Snippet 1), 2012.
- [18] Pere J Ferrando and Cristina Anguiano-Carrasco. Assessing the impact of faking on binary personality measures: An irt-based multiple-group factor analytic procedure. *Multivariate Behavioral Research*, 44(4):497–524, 2009.
- [19] Cheng-Han Leng, Hung-Yu Huang, and Grace Yao. A social desirability item response theory model: Retrieve–deceive–transfer. *Psychometrika*, 85(1):56–74, nov 2019.
- [20] Michael Schulte-Mecklenbeck, Anton Kühberger, and Joseph G Johnson. *A handbook of process tracing methods for decision research: A critical review and user’s guide*. Psychology Press, 2011.
- [21] Antonio Calcagni and L Lombardi. Dynamic fuzzy rating tracker (dyfrat): a novel methodology for modeling real-time dynamic cognitive processes in rating scales. *Applied soft computing*, 24:948–961, 2014.
- [22] Sara de la Rosa de Saa, María Ángeles Gil, Gil Gonzalez-Rodriguez, María Teresa López, and María Asunción Lubiano. Fuzzy rating scale-based questionnaires and their statistical analysis. *IEEE Transactions on Fuzzy Systems*, 23(1):111–126, 2014.
- [23] Tim Hesketh, Robert Pryor, and Beryl Hesketh. An application of a computerized fuzzy graphic rating scale to the psychological measurement of individual differences. *International Journal of Man-Machine Studies*, 29(1):21–35, 1988.
- [24] Paotjai Vonglao. Application of fuzzy logic to improve the likert scale to measure latent variables. *Kasetsart Journal of Social Sciences*, 38(3):337–344, 2017.
- [25] Renato Coppi, Paolo Giordani, and Pierpaolo D’Urso. Component models for fuzzy data. *Psychometrika*, 71(4):733–761, 2006.



- [26] Heungsun Hwang, Wayne S DeSarbo, and Yoshio Takane. Fuzzy clusterwise generalized structured component analysis. *Psychometrika*, 72(2):181–198, 2007.
- [27] María Ángeles Gil Álvarez, María Asunción Lubiano Gómez, Sara de la Rosa de Sáa, Beatriz Sinova Fernández, et al. Analyzing data from a fuzzy rating scale-based questionnaire: a case study. *Psicothema*, 2015.
- [28] Georg E Matt, Maria R Turingan, Quyen T Dinh, Julie A Felsch, Melbourne F Hovell, and Christine Gehrman. Improving self-reports of drug-use: numeric estimates as fuzzy sets. *Addiction*, 98(9):1239–1247, 2003.
- [29] Isabella Morlini. Fuzzy methods for the analysis of psychometric data: An application for measuring reading disability. *Statistica & Applicazioni*, 16(1), 2018.
- [30] María Asunción Lubiano, Sara de la Rosa de Sáa, Manuel Montenegro, Beatriz Sinova, and María Ángeles Gil. Descriptive analysis of responses to items in questionnaires. why not using a fuzzy rating scale? *Information Sciences*, 360:131–148, 2016.
- [31] María Ángeles Gil and Gil González-Rodríguez. Fuzzy vs. likert scale in statistics. In *Combining experimentation and theory*, pages 407–420. Springer, 2012.
- [32] Concepción San Luis Costas, Pedro Prieto Maranon, and Juan A Hernandez Cabrera. Application of diffuse measurement to the evaluation of psychological structures. *Quality and Quantity*, 28(3):305–313, 1994.
- [33] Itziar García-Honrado, Miquel Ferrer, and Angela Blanco-Fernandez. A tentative fuzzy assessment of the quality of teaching and opportunities to learn mathematics in a classroom discussion. In *2015 Conference of the International Fuzzy Systems Association and the European Society for Fuzzy Logic and Technology (IFSA-EUSFLAT-15)*. Atlantis Press, 2015.
- [34] Ana M Castaño, M Asunción Lubiano, and Antonio L García-Izquierdo. Gendered beliefs in stem undergraduates: A comparative analysis of fuzzy rating versus likert scales. *Sustainability*, 12(15):6227, 2020.
- [35] Inés M Gómez-Chacón. Emotions and heuristics: The state of perplexity in mathematics. *Zdm*, 49(3):323–338, 2017.
- [36] Patricia Conde-Clemente, Jose M Alonso, Éldman O Nunes, Angel Sanchez, and Gracian Trivino. New types of computational perceptions: Linguistic descriptions in deforestation analysis. *Expert Systems with Applications*, 85:46–60, 2017.
- [37] María Asunción Lubiano Gómez, Pilar González Gil, Helena Sánchez Pastor, Carmen Pradas, Henar Arnillas, et al. An incipient fuzzy logic-based analysis of the medical specialty in uence on the perception about mental patients. *The Mathematics of the Uncertain: A Tribute to Pedro Gil*, 2018.

- [38] Adrian Castro-Lopez and Jose M Alonso. Modeling human perceptions in e-commerce applications: A case study on business-to-consumers websites in the textile and fashion sector. In *Applying Fuzzy Logic for the Digital Economy and Society*, pages 115–134. Springer, 2019.
- [39] Ana Belén Ramos-Guajardo, Ángela Blanco-Fernández, and Gil González-Rodríguez. Applying statistical methods with imprecise data to quality control in cheese manufacturing. In *Soft Modeling in Industrial Manufacturing*, pages 127–147. Springer, 2019.
- [40] Qing Li. Indirect membership function assignment based on ordinal regression. *Journal of Applied Statistics*, 43(3):441–460, 2016.
- [41] Yuan Horng Lin and Jeng Ming Yih. Comparisons on reliability of likert scale between crisp and fuzzy data. In *Applied Mechanics and Materials*, volume 635, pages 874–877. Trans Tech Publ, 2014.
- [42] Jyh-Rong Chou. A psychometric user experience model based on fuzzy measure approaches. *Advanced Engineering Informatics*, 38:794–810, 2018.
- [43] Muluken Yeheyis, Bahareh Reza, Kasun Hewage, Janaka Y Ruwanpura, and Rehan Sadiq. Evaluating motivation of construction workers: A comparison of fuzzy rule-based model with the traditional expectancy theory. *Journal of Civil Engineering and Management*, 22(7):862–873, 2016.
- [44] M.A. Lazim and M.T. Abu Osman. Measuring teachers’ beliefs about mathematics: a fuzzy set approach. *International Journal of Social Sciences*, 4(1):39–43, 2009.
- [45] M.A. Lazim, M.T. Abu Osman, and W.A. Wan Salihin. Fuzzy set conjoint model in describing students’ perceptions on computer algebra system learning environment. *International Journal of Computer Science Issues (IJCSI)*, 8(2):92, 2011.
- [46] Konul Memmedova. Quantitative analysis of effect of pilates exercises on psychological variables and academic achievement using fuzzy logic. *Quality & Quantity*, 52(1):195–204, 2018.
- [47] Rahib H Abiyev, Tulen Saner, Serife Eyupoglu, and Gunay Sadikoglu. Measurement of job satisfaction using fuzzy sets. *Procedia Computer Science*, 102:294–301, 2016.
- [48] Pierpaolo D’Urso, Marta Disegna, Riccardo Massari, and Linda Osti. Fuzzy segmentation of postmodern tourists. *Tourism Management*, 55:297–308, 2016.
- [49] Marta Disegna, Pierpaolo D’Urso, and Riccardo Massari. Analysing cluster evolution using repeated cross-sectional ordinal data. *Tourism Management*, 69:524–536, 2018.
- [50] Pierpaolo D’Urso, Marta Disegna, and Riccardo Massari. Satisfaction and tourism expenditure behaviour. *Social Indicators Research*, pages 1–26, 2020.

- [51] Mehmet Ozer Demir, Murat Alper Basaran, and Biagio Simonetti. Determining factors affecting healthcare service satisfaction utilizing fuzzy rule-based systems. *Journal of Applied Statistics*, 43(13):2474–2489, 2016.
- [52] Toni Lupo. A fuzzy servqual based method for reliable measurements of education quality in italian higher education area. *Expert systems with applications*, 40(17):7096–7110, 2013.
- [53] Dian-Fu Chang, An Chen Chiu, and Berlin Wu. Fuzzy correlation among student engagement and interpersonal interactions. *ICIC Express Letters, Part B: Applications*, 9(1):17–22, 2018.
- [54] Shahid Hussain, Prashant K Jamwal, Muhammad T Munir, and Aigerim Zuyeva. A quasi-qualitative analysis of flipped classroom implementation in an engineering course: from theory to practice. *International Journal of Educational Technology in Higher Education*, 17(1):1–19, 2020.
- [55] Hong Tau Lee and Sheu Hua Chen. Using cpk index with fuzzy numbers to evaluate service quality. *International Transactions in Operational Research*, 9(6):719–730, 2002.
- [56] Ming-Tien Tsai, Hsueh-Liang Wu, and Wen-Ko Liang. Fuzzy decision making for market positioning and developing strategy for improving service quality in department stores. *Quality & Quantity*, 42(3):303–319, 2008.
- [57] Hung-Tso Lin. Fuzzy application in service quality analysis: An empirical study. *Expert systems with Applications*, 37(1):517–526, 2010.
- [58] Hsiu-Yuan Hu, Yu-Cheng Lee, and Tieh-Min Yen. Service quality gaps analysis based on fuzzy linguistic servqual with a case study in hospital out-patient services. *The TQM Journal*, 2010.
- [59] Michele Lalla, Gisella Facchinetti, and Giovanni Mastroleo. Ordinal scales and fuzzy set systems to measure agreement: an application to the evaluation of teaching activity. *Quality and Quantity*, 38(5):577–601, 2005.
- [60] Maria Symeonaki and Aggeliki Kazani. Developing a fuzzy likert scale for measuring xenophobia in greece. *ASMDA, Rome*, 2011.
- [61] Zsuzsanna E Tóth, Gábor Árvai, and Rita V Dénes. Are the ‘illnesses’ of traditional likert scales treatable? *Quality Innovation Prosperity*, 24(2):120–136, 2020.
- [62] Zsuzsanna E Tóth, Tamás Jónás, and Rita Veronika Dénes. Applying flexible fuzzy numbers for evaluating service features in healthcare—patients and employees in the focus. *Total Quality Management & Business Excellence*, 30(sup1):S240–S254, 2019.
- [63] Tamás Jónás, Zsuzsanna Eszter Tóth, and Gábor Árvai. Applying a fuzzy questionnaire in a peer review process. *Total Quality Management & Business Excellence*, 29(9-10):1228–1245, 2018.

- [64] Jan Stoklasa, Tomáš Talášek, and Pasi Luukka. Fuzzified likert scales in group multiple-criteria evaluation. In *Soft computing applications for group decision-making and consensus modeling*, pages 165–185. Springer, 2018.
- [65] Elvira Di Nardo and Rosaria Simone. A model-based fuzzy analysis of questionnaires. *Statistical Methods & Applications*, 28(2):187–215, 2019.
- [66] Donata Marasini, Piero Quatto, and Enrico Ripamonti. Evaluating university courses: intuitionistic fuzzy sets with spline functions modelling. *Statistica & Applicazioni*, 15(1), 2017.
- [67] Sen-Chi Yu and Min-Ning Yu. Fuzzy partial credit scaling: A valid approach for scoring the beck depression inventory. *Social Behavior and Personality: an international journal*, 35(9):1163–1172, 2007.
- [68] Sen-Chi Yu and Berlin Wu. Fuzzy item response model: a new approach to generate membership function to score psychological measurement. *Quality and Quantity*, 43(3):381, 2009.
- [69] María Asunción Lubiano, Manuel Montenegro, Beatriz Sinova, Sara de la Rosa de Sáa, and María Ángeles Gil. Hypothesis testing for means in connection with fuzzy rating scale-based data: algorithms and applications. *European Journal of Operational Research*, 251(3):918–929, 2016.
- [70] María Asunción Lubiano, Antonia Salas, Carlos Carleos, Sara de la Rosa de Sáa, and María Ángeles Gil. Hypothesis testing-based comparative analysis between rating scales for intrinsically imprecise data. *International Journal of Approximate Reasoning*, 88:128–147, 2017.
- [71] María Asunción Lubiano, Antonia Salas, Sara de la Rosa de Sáa, Manuel Montenegro, and María Ángeles Gil. An empirical analysis of the coherence between fuzzy rating scale-and likert scale-based responses to questionnaires. In *International Conference on Soft Methods in Probability and Statistics*, pages 329–337. Springer, 2016.
- [72] Irene Arellano, Beatriz Sinova, Sara de la Rosa de Sáa, María Asunción Lubiano, and María Ángeles Gil. Descriptive comparison of the rating scales through different scale estimates: Simulation-based analysis. In *International Conference Series on Soft Methods in Probability and Statistics*, pages 9–16. Springer, 2018.
- [73] Ana Belén Ramos Guajardo, María José González López, and Ignacio González Ruiz. Analysis of the reliability of the fuzzy scale for assessing the students’ learning styles in mathematics. In *2015 Conference of the International Fuzzy Systems Association and the European Society for Fuzzy Logic and Technology (IFSA-EUSFLAT-15)*, pages 727–733. Atlantis Press, 2015.
- [74] María Asunción Lubiano, Antonio L García-Izquierdo, and María Ángeles Gil. Fuzzy rating scales: Does internal consistency of a measurement scale benefit from coping with imprecision and individual differences in psychological rating? *Information Sciences*, 2020.

- [75] Po-Yi Chen and Grace Yao. Measuring quality of life with fuzzy numbers: in the perspectives of reliability, validity, measurement invariance, and feasibility. *Quality of Life Research*, 24(4):781–785, 2015.
- [76] Ernesto Araujo and Susana Abe Miyahira. Unidimensional fuzzy pain intensity scale. In *2009 IEEE International Conference on Fuzzy Systems*, pages 185–190. IEEE, 2009.
- [77] R Darrell Bock and Irini Moustaki. 15 item response theory in a general framework. *Handbook of statistics*, 26:469–513, 2006.
- [78] Wim J Van der Linden. *Handbook of item response theory: Volume 1: Models*. CRC Press, 2016.
- [79] Wim J van der Linden. *Handbook of Item Response Theory: Statistical Tools*. Chapman and Hall/CRC, 2017.
- [80] Ulf Böckenholt. Modeling multiple response processes in judgment and choice. *Decision*, 1(S):83–103, 2013.
- [81] Ulf Böckenholt. Measuring response styles in likert items. *Psychological Methods*, 22(1):69–83, 2017.
- [82] Thorsten Meiser, Hansjörg Plieninger, and Mirka Henninger. IRT tree models with ordinal and multidimensional decision nodes for response styles and trait-based rating responses. *British Journal of Mathematical and Statistical Psychology*, 72(3):501–516, feb 2019.
- [83] Paul De Boeck, Marjan Bakker, Robert Zwitser, Michel Nivard, Abe Hofman, Francis Tuerlinckx, Ivailo Partchev, et al. The estimation of item response models with the lmer function from the lme4 package in r. *Journal of Statistical Software*, 39(12):1–28, 2011.
- [84] Minjeong Jeon and Paul De Boeck. A generalized item response tree model for psychological assessments. *Behavior Research Methods*, 48(3):1070–1085, jul 2015.
- [85] Kwang Hyung Lee. *First course on fuzzy theory and applications*. Springer Science & Business Media, 2004.
- [86] Didier Dubois and Henri Prade. *Fundamentals of fuzzy sets*, volume 7. Springer Science & Business Media, 2012.
- [87] Antonio Calcagni, Luigi Lombardi, and Eduardo Pascali. Non-convex fuzzy data and fuzzy statistics: a first descriptive approach to data analysis. *Soft Computing*, 18(8):1575–1588, 2014.
- [88] Adel M Alimi. Beta neuro-fuzzy systems. *TASK Quarterly Journal, Special Issue on " Neural Networks*, 7(1):23–41, 2003.

- [89] Nesrine Baklouti, Ajith Abraham, and Adel M Alimi. A beta basis function interval type-2 fuzzy neural network for time series applications. *Engineering Applications of Artificial Intelligence*, 71:259–274, 2018.
- [90] William E Stein. Fuzzy probability vectors. *Fuzzy sets and Systems*, 15(3):263–267, 1985.
- [91] Sonia Migliorati, Agnese Maria Di Brisco, Andrea Ongaro, et al. A new regression model for bounded responses. *Bayesian Analysis*, 13(3):845–872, 2018.
- [92] Efendi N Nasibov and Sinem Peker. On the nearest parametric approximation of a fuzzy number. *Fuzzy Sets and Systems*, 159(11):1365–1375, 2008.
- [93] T. M. Williams. Practical use of distributions in network analysis. *Journal of the Operational Research Society*, 43(3):265–270, mar 1992.
- [94] Xiang-Bin Meng, Jian Tao, and Ning-Zhong Shi. An item response model for likert-type data that incorporates response time in personality measurements. *Journal of Statistical Computation and Simulation*, 84(1):1–21, 2014.
- [95] Patrick C Kyllonen and Jiyun Zu. Use of response time for measuring cognitive ability. *Journal of Intelligence*, 4(4):14, 2016.
- [96] Christopher Donkin and Scott D Brown. Response times and decision-making. *Stevens’ Handbook of Experimental Psychology and Cognitive Neuroscience*, 5:1–33, 2018.
- [97] Mollie E. Brooks, Kasper Kristensen, Koen J. van Benthem, Arni Magnusson, Casper W. Berg, Anders Nielsen, Hans J. Skaug, Martin Maechler, and Benjamin M. Bolker. glmmTMB balances speed and flexibility among packages for zero-inflated generalized linear mixed modeling. *The R Journal*, 9(2):378–400, 2017.
- [98] Dylan Molenaar and Paul de Boeck. Response mixture modeling: Accounting for heterogeneity in item characteristics across response times. *psychometrika*, 83(2):279–297, 2018.
- [99] Michael J Zickar. Modeling faking on personality tests. 2000.
- [100] Philseok Lee, Seang-Hwane Joo, and Shea Fyffe. Investigating faking effects on the construct validity through the monte carlo simulation study. *Personality and Individual Differences*, 150:109491, 2019.
- [101] Michael J Zickar, Robert E Gibby, and Chet Robie. Uncovering faking samples in applicant, incumbent, and experimental data sets: An application of mixed-model item response theory. *Organizational Research Methods*, 7(2):168–190, 2004.
- [102] Massimiliano Pastore, Massimo Nucci, Andrea Bobbio, and Luigi Lombardi. Empirical scenarios of fake data analysis: The sample generation by replacement (sgr) approach. *Frontiers in psychology*, 8:482, 2017.

- [103] Gian Vittorio Caprara. *La valutazione dell'autoefficacia. Costrutti e strumenti*. Edizioni Erickson, 2001.
- [104] Joshua D Greene, R Brian Sommerville, Leigh E Nystrom, John M Darley, and Jonathan D Cohen. An fmri investigation of emotional engagement in moral judgment. *Science*, 293(5537):2105–2108, 2001.
- [105] Alexander Behnke, Anja Strobel, and Diana Armbruster. When the killing has been done: Exploring associations of personality with third-party judgment and punishment of homicides in moral dilemma scenarios. *Plos one*, 15(6):e0235253, 2020.
- [106] Achim Zeileis, Francisco Cribari-Neto, Bettina Grün, and I Kos-Midis. Beta regression in r. *Journal of statistical software*, 34(2):1–24, 2010.
- [107] Jerry L Deffenbacher, Eugene R Oetting, and Rebekah S Lynch. Development of a driving anger scale. *Psychological reports*, 74(1):83–91, 1994.
- [108] Ulf Böckenholt. Modeling motivated misreports to sensitive survey questions. *Psychometrika*, 79(3):515–537, 2014.
- [109] Eunike Wetzel, Susanne Frick, and Samuel Greiff. The multidimensional forced-choice format as an alternative for rating scales. *European Journal of Psychological Assessment*, 36(4):511–515, 2020.
- [110] Valerie F. Reyna and Charles J. Brainerd. Numeracy, ratio bias, and denominator neglect in judgments of risk and probability. *Learning and Individual Differences*, 18(1):89–107, jan 2008.
- [111] David Preinerstorfer and Anton K Formann. Parameter recovery and model selection in mixed rasch models. *British Journal of Mathematical and Statistical Psychology*, 65(2):251–262, 2012.
- [112] Thomas R O'Neill, Justin L Gregg, and Michael R Peabody. Effect of sample size on common item equating using the dichotomous rasch model. *Applied Measurement in Education*, 33(1):10–23, 2020.
- [113] Anton A Béguin and Ceec AW Glas. Mcmc estimation and some model-fit analysis of multi-dimensional irt models. *Psychometrika*, 66(4):541–561, 2001.
- [114] Paul-Christian Bürkner. brms: An r package for bayesian multilevel models using stan. *Journal of statistical software*, 80(1):1–28, 2017.
- [115] József Dombi and Tamás Jónás. Approximations to the normal probability distribution function using operators of continuous-valued logic. *Acta Cybernetica*, 23(3):829–852, 2018.