



Inference of preference heterogeneity from choice data

Annie Liang¹

University of Pennsylvania, United States

Received 4 October 2016; final version received 16 August 2018; accepted 23 September 2018
Available online 23 October 2018

Abstract

Suppose an analyst observes inconsistent choices from either a single decision-maker, or a population of agents. Can the analyst determine whether this inconsistency arises from choice error (imperfect maximization of a single preference) or from preference heterogeneity (deliberate maximization of multiple preferences)? I model choice data as generated from imperfect maximization of a small number of preferences. The main results show that (a) simultaneously minimizing the number of inferred preferences and the number of unexplained observations can exactly recover the number of underlying preferences with high probability; (b) simultaneously minimizing the *richness* of the set of preferences and the number of unexplained observations can exactly recover the choice implications of the decision maker's underlying preferences with high probability.

© 2018 Elsevier Inc. All rights reserved.

JEL classification: D01; D11; D80; C52; D90

Keywords: Choice theory; Multiple rationales; Heterogeneity; Revealed preference; Identifiability

1. Introduction

Let X be a finite set of choice alternatives, and consider an analyst who observes choices (by either a single decision-maker, or a population of subjects) from various subsets of X . Empirical

E-mail address: anliang@upenn.edu.

¹ I am especially grateful to Jerry Green. I would also like to thank David Ahn, Emily Berger, Gabriel Carroll, Ian Crawford, Aluma Dembo, Drew Fudenberg, Ben Golub, David Laibson, Jay Lu, Erik Madsen, Eric Maskin, Jose Montiel Olea, Krishna Pendakur, Gleb Romanyuk, Andrei Shleifer, Ran Shorrer, Pedro Brandeo Solti, Ran Spiegler, Tomasz Strzalecki, and Yufei Zhao for useful comments and discussions.

choice data of this nature is often inconsistent, and cannot be explained as perfect maximization of a single preference.

There are two different perspectives for how to interpret such inconsistency. One view is that inconsistency emerges from *preference heterogeneity*. There is abundant evidence that choices depend on details about the choice context—for example, Einav et al. (2012) find that just over 30% of their subject pool makes decisions across six financial domains that can be rationalized using a common risk preference. Additionally, choice data aggregated over a population of decision-makers often exhibits cross-sectional heterogeneity in preferences—for example, Crawford and Pendakur (2012) study household consumption decisions over different kinds of milk, and find that no more than two-thirds of observations in their data set can be rationalized using a single utility function. In both of these cases, choice inconsistencies are understood to reflect intentional maximization which is welfare-relevant.

Another view is that inconsistencies reflect *errors*, e.g. the decision-maker may be inattentive, and the analyst may make mistakes while recording observations. In these cases, inconsistency reflects choices that are not indicative of preference.

Preference heterogeneity and error are distinct sources of inconsistency, with different implications for welfare-assessment and for prediction. To take a stark example, compare two hypothetical choice data sets: one generated by *perfect maximization of two different preferences*, and one generated by *maximization of a single preference with trembling error*. Application of classical approaches such as Houtman and Maks (1985) can fail to distinguish between these data sets, especially if the same fraction of both data sets is rationalizable using a single preference. Nevertheless, the underlying choice mechanics are quite different: the inconsistency represented in the first data set can be expected to be stable across future observations, while the inconsistency represented in the second is idiosyncratic.

A basic question then is *how many* choice domains or subpopulations are present in the data, where a special case of interest is whether there is evidence of multiple preferences or a single preference with error. (The question of “how many preferences” is pursued, for example, in Crawford and Pendakur (2012) in the case of household consumption decisions and Dean and Martin (2010) for individual choices over lotteries.)

At two ends for interpreting the data are the classical approaches of Houtman and Maks (1985) and Kalai et al. (2002), both of which rule out one of the two sources of inconsistency described above. Specifically, we can find a “best-fit” single preference (rationalizing the largest fraction of observations) and interpret the remaining observations as choice errors (Houtman and Maks, 1985), or find the smallest number of preferences that perfectly rationalizes the choice data (Kalai et al., 2002). When preference multiplicity and choice errors are simultaneously present in the data, the Houtman and Maks (1985) solution (weakly) underestimates the number of preferences, while the Kalai et al. (2002) solution (weakly) overestimates. The consequences for welfare evaluation and out-of-sample predictions can be significant.²

The purpose of this paper is to develop a method to determine the “best” intermediate solution. I consider data generated according to a (generalized) random utility model. The decision-maker (DM) chooses from choice set $A \subseteq X$ by sampling a preference according to a distribution μ_A , and maximizing the sampled preference. I suppose that each μ_A is in fact a perturbation of a

² See Appendix A for examples in which use of these approaches to predict choice behaviors leads to suboptimal prediction accuracy.

“sparse” μ_A^* , whose support is a small number of preferences K that are constant across choice sets.³ The goal is to recover from the choice data the underlying number of preferences K .

The proposed approach, presented in Section 4, minimizes a weighted sum of the number of preferences attributed to the decision-maker, and the number of unexplained observations (choices that cannot be rationalized by any of the recovered preferences). The approach thus imposes a cost on each recovered preference, so that a preference is recovered if and only if it explains sufficiently many observations that would otherwise be considered error. The classic Houtman and Maks (1985) and Kalai et al. (2002) solutions are returned for special choices of weights—the former is returned when the cost of preferences relative to unexplained observations is sufficiently high, and the latter is returned when the cost of preferences relative to unexplained observations is sufficiently low.

The main result in Section 5 provides a set of weights (which depend on primitives of the choice model) given which the proposed approach exactly recovers the “true” number of preferences with sufficiently many observations.⁴ Informally, these conditions require that the K preferences are sufficiently differentiated in the sampled data, so that choice inconsistencies that emerge from genuine preference heterogeneity resemble other inconsistencies in the choice data, whereas choice inconsistencies that emerge from error appear idiosyncratic. The Kalai et al. (2002) approach is shown to recover the number of underlying preferences when the probability of choice error is zero; to the best of my knowledge, this is the first statistical justification for the Kalai et al. (2002) approach. The special case of discerning between choice data that is generated by imperfect maximization of a single preference, versus choice data that reflects multiplicity of preference, is considered in Section 5.4.

The set of weights which allow for recovery depends on primitives of the choice model. Next, I explain a way in which we can “test” particular assumptions about these unobservables based on the data. Since the main theorem provides an interval of weights that recover the same solution, we can use the data to determine the actual range of weights over which our inferred solution remains stable. This range can then be used to bound the key primitives (the extent of differentiation of the underlying preferences, and the probability of error), as shown in Corollary 1.

Section 7 revisits an analysis conducted in Crawford and Pendakur (2012), in which the Kalai et al. (2002) approach is used to discover the number of preference types among 500 subjects. Crawford and Pendakur (2012) find that five preferences are needed to perfectly rationalize their data set. I show how the proposed approach can be used to identify some of these preferences as noise.

Section 8 turns to the question of recovering the preferences themselves. Inference of multiple preferences from choice data is an ill-posed problem, and Section 5.1 presents several negative results that help to clarify the reasons for this. In Proposition 1, I show that most sets of orderings are indistinguishable based on their choice implications, so that even in the absence of choice error, most sets of multiple preferences cannot be recovered. This result is very much in the spirit of Ambrus and Rozen (2013), which studies a broad (but different) class of multi-self models and shows that these models have no testable implications without prior restrictions on the number of selves involved in a decision.⁵

³ This paper takes a nonparametric, or “model-free” approach in the spirit of Varian (1982), Famulari (1995), Houtman and Maks (1985), and Kalai et al. (2002)).

⁴ I assume that the DM may be presented with the same choice problem multiple times.

⁵ Ambrus and Rozen (2013) study several choice-set independent aggregation rules over preferences, whereas I consider a specific aggregation rule (in which one preference is assigned “dictator”) that varies across choice problems.

In view of these results, I suggest that a more appropriate object of recovery is the set of *choice implications* of the decision maker’s preferences—that is, the choice observations that are consistent with maximization of one of these preferences. I define equivalence classes for sets of preferences, where two sets belong to the same equivalence class if they have the same choice implications, and ask whether we can recover the equivalence class to which the true set of preferences belongs. Section 8.2 shows that this is indeed possible, but that penalizing the number of inferred preferences is not the appropriate criterion for this goal. This is because penalizing only the number of preferences results in inference of sets of preferences whose choice implications are as diverse as possible. I propose an alternative criterion: minimizing a weighted sum of the number of unexplained observations and the *richness* of the set of preferences, as measured through the number of unique choice implications. Proposition 2 shows that under certain conditions on the choice model described above, this approach will exactly recover the equivalence class of choice implications containing those of the true model.

Finally, Section 8.3 considers a richer kind of data set, which includes auxiliary information on the choice contexts active during different observations. I show that with this additional information, we can (under certain conditions) recover the exact set of preferences.

Taken together, these results suggest that appropriately penalizing the complexity of the inferred choice model—for example, via the number of preferences used or the number of choice implications—can be useful for recovery of stable features of preference from inconsistent choice data.

2. Example

In an adaptation on the Luce and Raiffa dinner (Luce and Raiffa, 1957), suppose that a large number of consumers are observed to choose entrées from different restaurant menus. Each menu includes at least two main entrées from the set $\{x_1, x_2, \dots, x_N\}$, and additionally a special of frog legs (denoted x_{N+1}) is sometimes included.

As in Sen (1993), the presence of frog legs signals a high quality chef and encourages consumers to choose entrées that are harder to prepare; let us order entrées x_1, x_2, \dots, x_N from least to most difficult to prepare. If the menu A includes frog legs, then consumers choose each entrée $x \in A$ with probability

$$c_1(x|A) = \frac{e^{\gamma u_1(x)}}{\sum_{x' \in A} e^{\gamma u_1(x')}} \quad \text{when } x \neq x_{N+1}$$

where $\gamma \in \mathbb{R}_+$ is a logit parameter, and the utility function $u_1(x_k) = k$ assigns a higher payoff to entrées that are more difficult to prepare. Fix $c_1(x|A) = 0$ for every entrée $x \notin A$ and also $c_1(x_{N+1}|A) = 0$, so that frog legs are never themselves chosen.

When frog legs are *not* present on the menu, consumers’ choices follow the logit choice rule

$$c_2(x|A) = \frac{e^{\gamma u_2(x)}}{\sum_{x' \in A} e^{\gamma u_2(x')}} \quad \forall x \in A$$

where the utility function $u_2(x_k) = N - k + 1$ assigns higher payoffs to entrées that are easier to prepare. Fix $c_2(x|A) = 0$ for every $x \notin A$.

Choices from n menus are observed, where menus are sampled uniformly at random (with repetitions permitted) from the set of all menus that contain at least two entrées from $\{x_1, \dots, x_N\}$. I make minimal assumptions about the analyst’s knowledge about the consumers’ choice rule;

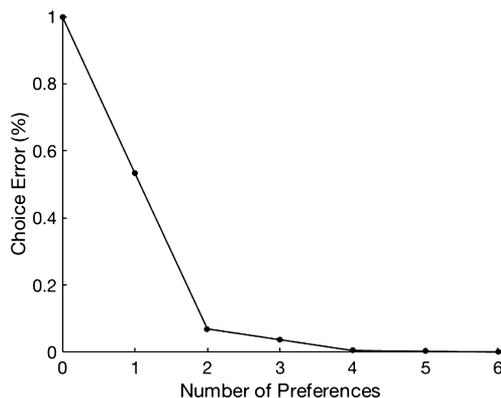


Fig. 1. *Error-preference tradeoff graph*. The (expected) fraction of the choice data that cannot be rationalized using any set of k preferences.

in particular, he does not know c_1 and c_2 , and does not know that there are two choice rules. He simply observes, and seeks to rationalize, the generated choices $\{(x_i, A_i)\}_{i=1}^n$.

The key feature of this example is that there are two distinct reasons why observed choices are unlikely to be consistent with maximization of any single ordering: first, consumers apply different choice rules in different observations; second, choice is stochastic (i.e. maximization is imperfect). At extremes, we can interpret the data in a way that rules out either of these two sources of inconsistency. For example, we can insist on a single preference and find the “best-fit” preference, interpreting the remaining choice observations as error. We can alternatively ascribe to consumers as many preferences as are needed to perfectly rationalize the data. Both of these approaches can lead to misinterpretations of the data, with consequences for welfare evaluation and prediction (discussed further below). The key question for this paper is how to determine instead that there are two primary preferences.

For concreteness, fix the logit parameter to be $\gamma = 2$, so that consumers choose the most preferred alternatives (under respectively u_1 or u_2) with high probability. (For example, given the choice set $\{x_1, x_2, x_3\}$, the most preferred alternative x_3 is chosen 86% of the time.) Consider also $N = 10$, so that there are ten main entrées. Fig. 1 represents the choice data using what I will call an *error-preference tradeoff graph*: For each number of preferences k , it reports the percentage of the choice data that cannot be rationalized using k preferences.⁶

It is possible to rationalize (in expectation) 47% of the choice data using a single preference, while two preferences can rationalize almost all of the data (93%). In contrast, the addition of a third ordering increases the completeness of explanation by only 3%, and the addition of the fourth preference contributes even less. Thus, each preference *up to* the second preference helps to rationalize a large fraction of the data, while each preference *after* the second preference contributes only a marginal improvement in explanation. The substantial drop in marginal explanation after the second ordering suggests the presence of two structural preferences, with additional trembling noise.

This is not the interpretation of either the Houtman and Maks (1985) or the Kalai et al. (2002) approaches. Direct application of Houtman and Maks (1985) produces an (expected) inconsis-

⁶ In practice, we would use the fraction of the actual choice data that is unexplained (interpret these as *choice errors*); for the purpose of this illustration, Fig. 1 reports the expected fractions.

tency measure of 53%, implying a quite irrational DM. Direct application of Kalai et al. (2002) overestimates the number of active preferences to six, and implies a fully rational DM.

We can also see, informally, that prediction of future choices based on these two approaches can be misguided.⁷ The best *single* preference recovered under Houtman and Maks (1985) predicts incorrectly in a majority of new choice problems, while some of the six preferences recovered under the Kalai et al. (2002) approach may perform worse than random guessing. A more rigorous treatment of the topic of out-of-sample prediction accuracy is deferred to Appendix A.

This discussion highlights some cases in which the outputs of the Houtman and Maks (1985) and Kalai et al. (2002) approaches are not well-suited to the analyst's goals. The key question is then how to identify an optimal solution intermediate to the two approaches (in this case, identifying two preferences). The sections below propose a method for doing this.

3. Conceptual framework

3.1. Choice model

Let X be a finite set of alternatives. A *preference* P is a strict linear ordering over X , and the set of all preferences is denoted by \mathcal{P} . A *choice set* is a subset $A \subseteq X$, and 2^X denotes the set of all possible choice sets. I will primarily interpret choices as generated by a single individual, although the framework below applies also when choices are generated by a population of subjects.

The decision-maker (DM) chooses from choice set A by sampling a preference according to a distribution $\mu_A \in \Delta(\mathcal{P})$, and maximizing the sampled preference. This corresponds to a standard generalization of the *random utility model* (RUM), where the decision-maker's distribution over preferences is permitted to vary across choice sets. I refer to $\mu = (\mu_A)_{A \in 2^X}$ as the decision-maker's RUM. Then, the probability that alternative x is selected from choice set A is

$$c(x|A) = \mu_A(\{P : x \text{ is } P\text{-maximal in } A\}). \quad (1)$$

An analyst observes the decision-maker's choices from n choice sets, sampled (with replacement) from a distribution $\pi \in \Delta(2^X)$. A *choice observation* is a pair (x, A) , corresponding to choice of alternative x from set A , and the observed data is a multiset of choice observations $D = \{(x_i, A_i)\}_{i=1}^n$. For simplicity, I will refer to D as simply a *set* of choice observations, although it should be understood that the same observation may appear multiple times. Finally, the ex-ante probability of observing any (x, A) (taking into account both the randomness in which choice sets are presented to the DM, and also the randomness in his choice) is

$$v(x, A) = \pi(A)c(x|A). \quad (2)$$

Notice that by explicitly modeling the sampling of choice sets, I depart from a standard assumption that the analyst knows the stochastic choice rule c , and thus has available to him an "idealized" data set where choices are made infinitely often from each choice set. In this paper, the analyst observes only a finite number of choices. If π assigns positive probability to every

⁷ This discussion is informal, since we have not fixed a model for choosing which preference to use in predicting new choice observations, which is required for applying multiple-preference models to prediction of new choice observations. The additional analysis pursued in Appendix A uses an extension described in Section 8.3 that makes these out-of-sample comparisons possible.

choice set, then the idealized data is returned as a limiting case when we take the number of observations to infinity.

3.2. Separation of preference from error

Consider the general choice framework described in the previous section. When the distributions μ_A are not degenerate, then violations of the Independence of Irrelevant Alternatives axiom are expected.⁸

There are two different perspectives for how to interpret these choice inconsistencies. One view is that the inconsistencies emerge from *preference heterogeneity*. For example, it may be that preference depends on unobserved features of the environment beyond the choice set, so that different preferences are cued in different choice observations. Relatedly, if choices are generated by a population of decision-makers, then inconsistencies may reflect cross-sectional heterogeneity in preferences across the population. According to both of these interpretations, the randomness over outcomes is reflective of intentional maximization. A second view is that these inconsistencies describe *measurement errors* or *choice errors*, which are welfare-reducing and not indicative of genuine preference. The choice model described previously permits both kinds of inconsistency to be simultaneously present.

In general, it will not always be possible to separate preference heterogeneity from error, even conceptually. This paper studies a setting in which there is a *small* set of preferences that are maximized *most* of the time, and suggests that preference heterogeneity and error can be meaningfully distinguished in this case. In particular, we may think of the small stable set of preferences as the “true” preferences—reflecting, for example, different sub-populations or different choice contexts—and the choices that are inconsistent with these preferences as error.

3.3. Underlying “sparse” choice model

Formally, I consider the setting in which the RUM μ is well-approximated by an underlying μ^* , supported on a “sparse” set of preferences.

Specifically, suppose that the DM has a set \mathcal{P} consisting of K preferences. In the absence of choice error, his choice from set A corresponds to maximization of a preference sampled from a distribution μ_A^* , where $\mu_A^*(\mathcal{P}) = 1$. (Interpret non-degenerate distributions μ_A^* as reflecting variation in the activation of preferences. For example, if preferences correspond to different choice contexts, then μ_A^* is the empirical distribution of contexts for the choice problem A .) Thus, when K is small, the RUM $\mu^* = (\mu_A^*)_{A \in 2^X}$ is supported on only a small number of preferences, relative to the complete set of preference orderings \mathcal{P} . In analogy to (1) and (2), define $c^*(x|A) = \mu_A^*({P : x \text{ is } P\text{-maximal in } A})$ for the stochastic choice rule associated with μ^* , and $v^*(x, A) = \pi(A)c^*(x|A)$ for the frequency of observation of (x, A) under RUM μ^* and sampling distribution π .

We do not observe choices generated under RUM μ^* , but rather choices generated under its perturbation μ . The relationship between these RUMs is given as follows: For each choice set A , there is a map $g_A : \mathcal{P} \rightarrow \Delta(\mathcal{P})$ such that

$$\mu_A = \mu_A^* G_A$$

⁸ Here, and throughout the paper, I refer to the classic (deterministic) version of IIA. Naturally, violations of the stochastic version of IIA may also be present.

where μ_A and μ_A^* are $1 \times |\mathcal{P}|$ vectors (choose an arbitrary indexing of preferences) and G_A is the $|\mathcal{P}| \times |\mathcal{P}|$ Markov matrix associated with g_A .⁹ It is important that choice errors do not occur frequently; formally, there is a (uniform) bound p on probability of error such that the diagonal entries in every G_A are at least $1 - p$. Informally, this guarantees that the “right” preference is maximized most of the time. I will refer to p throughout as the *probability of error*.

When K and p are both small, as is the primary case of interest, then most choices under the RUM μ correspond to maximization of a small number of preferences. Section 5 provides conditions on the sampling distribution π , the probability of error p , and the underlying RUM μ^* under which recovery of the number of preferences K is possible.¹⁰

Note the following special cases of the model:

Example 1. If $K = 1$, so that there is only one underlying preference, then each μ_A^* is degenerate on that preference and the observed choice data corresponds to imperfect maximization of a single preference ordering.

Example 2. If $K > 1$ but each μ_A^* is degenerate, then we return the multiple preference model introduced in Kalai et al. (2002), where different preferences are cued in different choice problems.

Throughout, I will take the perspective that μ^* is the DM’s “true” RUM, and K describes the cardinality of the DM’s true set of preferences. An alternative interpretation is that the DM’s true RUM is μ , but we prefer a more parsimonious description—specifically, an approximate representation by an RUM which assigns positive probability only to a small set of preferences. From this perspective, our problem is that of how many preferences are needed to explain most of the choice data generated under μ .

Another interpretation of the proposed framework, building on Rubinstein and Salant (2008) and Bernheim and Rangel (2009), is the following: Let \mathcal{C} be a set of contexts that are relevant to the DM’s preference but unobserved by the analyst.¹¹ Each context is associated with a preference; in a slight abuse of notation, let P_C be the preference associated with the context $C \in \mathcal{C}$. Each choice set A is associated with a distribution μ_A over contexts, and the probability of observing choice of x from set A is given by

$$c(x|A) = \mu_A(\{C : x \text{ is } P_C\text{-maximal in } A\}).$$

The question of interest is whether it is possible to recover the number of unique context-dependent preferences given observation of pairs (x, A) only. (It is also interesting to consider what we can learn when the contexts are observable, and Section 8.3 explores this case further.)

4. Analysis of the choice data

Fix a data set $D = \{(x_i, A_i)\}_{i=1}^n$. A *multiple preference rationalization* of this data is any set of preference orderings $\mathcal{P} \subseteq \mathcal{P}$. The number of *implied choice errors* when using \mathcal{P} to rationalize D is:

$$\varepsilon(D, \mathcal{P}) := \#\{(x, A) \in D : x \text{ is not } P\text{-maximal in } A \text{ for any } P \in \mathcal{P}\}.$$

⁹ Index the preferences $P_1, \dots, P_{N!}$. The i -th row of G_A is $g_A(P_i)$.

¹⁰ See Section 5.4 for application of the approach in an example in which the probability of p is large.

¹¹ These are called *frames* in Rubinstein and Salant (2008) and *ancillary conditions* in Bernheim and Rangel (2009).

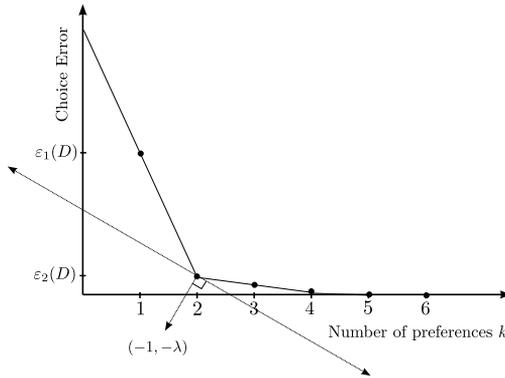


Fig. 2. Define E to be the set of points lying above the linear interpolation of $\{(k, \varepsilon_k(D))\}_{k \in \mathbb{Z}_+}$. Then, the problem in (3) returns a solution with k orderings if and only if the line with normal vector $(-1, -\lambda)$ properly supports E at $(k, \varepsilon_k(D))$.

This counts the number of observed choices that are inconsistent with maximization of any preference in \mathcal{P} .¹² Say that \mathcal{P} constitutes a *perfect rationalization* of D if $\varepsilon(D, \mathcal{P}) = 0$. When restricting to sets of k preferences, the minimal number of implied choice errors is

$$\varepsilon_k(D) := \min_{|\mathcal{P}|=k} \varepsilon(D, \mathcal{P}).$$

Say that the data set D is *k-rationalizable* (Kalai et al., 2002) if there is some set of k preferences that perfectly rationalizes D , so that $\varepsilon_k(D) = 0$.

It is useful to represent D as the linear interpolation of points in $\{(k, \varepsilon_k(D))\}_{k \in \mathbb{Z}_+}$, henceforth its *Error-Preference Tradeoff Graph*. Each point on the convex hull of this graph represents a particular weighted minimization of the number of preferences k (ascribed to the DM), and the number of implied choice errors $\varepsilon_k(D)$.¹³ Formally, every choice of tradeoff λ between these objectives determines a corresponding solution:

Definition 4.1. For every $\lambda \in \mathbb{R}_+$ and data set D , define

$$K_\lambda^*(D) = \operatorname{argmin}_{k \in \mathbb{Z}_+} [k + \lambda \varepsilon_k(D)]. \tag{3}$$

This solution is depicted in Fig. 2. When there are multiple solutions to the problem above, take $K_\lambda^*(D)$ to mean the smallest value of k in the minimizing set.

Intuitively, $1/\lambda$ is the “cost” of each ordering, so that an ordering is attributed to the DM if and only if it explains at least $1/\lambda$ observations that would otherwise be interpreted as choice error. As $\lambda \rightarrow 0$, the cost of errors becomes increasingly small relative to the cost of orderings, so the analyst prefers to attribute a single ordering to the DM and interpret the unexplained observations as choice errors. As $\lambda \rightarrow \infty$, the cost of choice errors becomes increasingly large

¹² Naturally, this is only one of many possible definitions for choice error. In particular, other notions of error may be preferred for different models of preference aggregation (see e.g. the multiple-ordering models of Rubinstein and Salant, 2006; Fudenberg and Levine, 2006; Green and Hojman, 2007; Manzini and Mariotti, 2007, 2009).

¹³ In Section 2 we plotted the fraction of choice errors against the number of preferences, but from here on we will use the number of choice errors. This choice is unimportant, and the resulting analysis can be conducted under either convention.

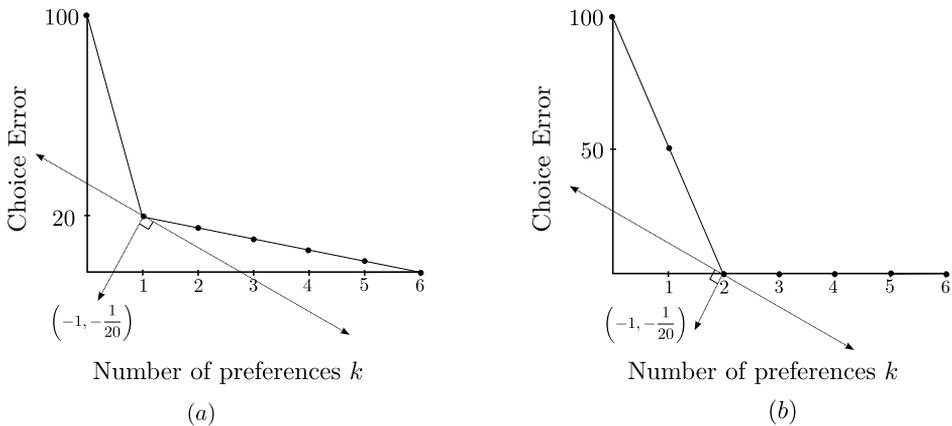
relative to the cost of orderings, so the analyst prefers to use as many orderings as necessary to perfectly rationalize the data.

In particular, if $\lambda < \frac{1}{\varepsilon_1(D)}$ (recall that $\varepsilon_1(D)$ is the necessary number of unexplained observations if the DM is attributed a single preference), then the approach returns the Houtman and Maks (1985) solution, and if $\lambda > 1$, then the approach returns the Kalai et al. (2002) solution.

Observation 1. For every data set D :

- (a) $K_\lambda^*(D) = 1$ for every $\lambda < \frac{1}{\varepsilon_1(D)}$, and
- (b) $K_\lambda^*(D) = L$ for every $\lambda > 1$, where L is the smallest integer such that D is L -rationalizable.

The intervals provided in Observation 1 are sufficient but not necessary for recovery of the Houtman and Maks (1985) and Kalai et al. (2002) solutions. In particular, the same choice of λ can correspond to either solution, depending on the choice data. For example, choice of $\lambda = 1/20$ selects $K_\lambda^*(D) = 1$ (the Houtman and Maks, 1985 solution) in panel (a) of the figure below, and selects $K_\lambda^*(D) = 2$ (the Kalai et al., 2002 solution) in panel (b).



These selections agree with the intuition that a data set represented by (a) resembles imperfect maximization of a single ordering, while a data set represented by (b) resembles perfect maximization of two orderings. The main results below formalize these intuitions, relating the “optimal” choice(s) of λ to the primitives of the choice model described in Section 3.

5. Recovering the number of preferences

5.1. No error baseline: $\mu = \mu^*$

Consider first the recovery problem for an idealized baseline in which the decision-maker’s RUM is exactly μ^* , so that there are no choice errors. Even in this setting, recovery of the number of preferences K is not guaranteed to be feasible. The key condition needed for recovery of K is that preferences are sufficiently differentiated in the data. This differentiation depends *jointly* on the sampling procedure π and also on the underlying RUM μ^* .

Basic challenges to recovery are illustrated in Examples 3–5 below. In Example 3, recovery is not feasible because the DM’s preferences are not sufficiently differentiated by their choice implications; in Example 4, recovery is not feasible because a preference is insufficiently sampled; and in Example 5, recovery is not feasible because preferences agree on the sampled choice sets. Recovering the number of preferences is not obviously meaningful in these cases, and such settings will be subsequently ruled out.

Example 3. The set of choice alternatives is $X = \{x_1, x_2, x_3\}$ and the DM’s preferences are $\mathcal{P} = \{P_1, P_2, P_3\}$, where

$$\begin{aligned} &x_1 P_1 x_2 P_1 x_3 \\ &x_1 P_2 x_3 P_2 x_2 \\ &x_3 P_3 x_2 P_3 x_1 \end{aligned}$$

Fix any RUM μ^* supported on \mathcal{P} , and any sampling distribution π over choice sets. Since every choice observation consistent with maximization of some ordering in \mathcal{P} is also consistent with maximization of an ordering in $\mathcal{P}' = \{P_1, P_3\}$, there is no value of λ given which $K_\lambda^*(D) = 3$ for any data set D generated under this model.

Example 4. The DM’s preferences are $\mathcal{P} = \{P_1, P_2\}$, but only the first preference is sampled; that is, $\mu_A^*(P_1) = 1$ for every choice set A . For all data sets D generated under this model, there is no value of λ given which $K_\lambda^*(D) = 2$.¹⁴

Example 5. The DM’s preferences are $\mathcal{P} = \{P_1, P_2\}$ where

$$\begin{aligned} &x_1 P_1 x_2 P_1 x_3 \\ &x_1 P_2 x_3 P_2 x_2 \end{aligned}$$

The sampling procedure π puts probability 1 on the choice set $\{x_1, x_2, x_3\}$. Since P_1 and P_2 agree on this choice set, for every RUM μ^* (supported on \mathcal{P}), every data set D generated under this model, and every choice of λ , the proposed approach yields $K_\lambda^*(D) = 1$.

These examples highlight that sufficient differentiation of preferences in the data requires first that preferences have different choice implications, and second that there is opportunity for these choice implications to be observed (namely, that the corresponding choice problems and preferences are sampled in the data). The previous examples have the property that the number of underlying preferences cannot be recovered from *any* number of choice observations. Each of these obstacles to recovery may also appear in a more moderate degree: for example, some preference may be rarely sampled, causing its choice implications to appear rarely in the data.

Following, I define a measure for how *differentiated* the K underlying preferences are in the choice data. This measure is defined as a property of the primitives π and μ^* . Sufficient differentiation in preferences serves a dual role: it simultaneously makes recovery of the number of preferences possible, and it also justifies consideration of this number as an object of

¹⁴ I am grateful to an anonymous referee for suggesting this example.

interest. Note that in each of the previous examples, recovery of the number of “true” preferences is an arguably misguided exercise—for example, the set of three preferences \mathcal{P} in Example 3 is a needlessly complex way to rationalize choice data that can also be explained using $\mathcal{P}' \subset \mathcal{P}$.

As a preliminary step, I first define a generalization of IIA¹⁵:

Definition 5.1. For any integer k , say that choice observations $\{(x_i, A_i)\}_{i=1}^k$ are in k -violation of IIA if

- (1) $x_i \neq x_j$ for every $i \neq j$,
- (2) $x_i \in \bigcap_{j=1}^k A_j$ for every $i = 1, \dots, k$.

The first condition requires that every chosen alternative is different, and the second condition requires that each of the chosen alternatives is available in all of the observed choice sets. An immediate implication is that the set of choice observations $\{(x_i, A_i)\}_{i=1}^k$ cannot be rationalized using fewer than k orderings (without introducing choice error). Notice also that every pair of choice observations from $\{(x_i, A_i)\}_{i=1}^k$ constitutes a (standard) violation of IIA. Thus, the set of observations suggests that the decision maker possesses at least k different preferences.

Of special interest are choice observations in K -violation of IIA, where K is the true number of underlying preferences. Below, I define *differentiation* to be the (limiting) fraction of choice observations that can be partitioned into disjoint K -violations of IIA.

Definition 5.2 (Differentiation). For each data set D , define $g(D)$ to be the largest number of disjoint subsets of choice observations in D that are in K -violation of IIA. The *differentiation parameter* for primitives (π, μ^*) is

$$d(\pi, \mu^*) = \liminf_{n \rightarrow \infty} \mathbb{E}_{(v^*)^n} \left[\frac{1}{n} g(D_n) \right] \tag{4}$$

where $(v^*)^n$ is the product measure corresponding to n i.i.d. draws from v , and $D_n \sim (v^*)^n$ is a random data set of size n .

To interpret the differentiation parameter $d(\pi, \mu^*)$, suppose that n choice observations are generated from the choice model described by primitives π and μ^* . Then, in expectation, there is a partitioning of the realized choice data such that at least $n \cdot d(\pi, \mu^*)$ partition elements are in K -violation of IIA. Recalling that choice observations in K -violation of IIA require K preferences for perfect rationalization, large values of $d(\pi, \mu^*)$ imply that use of fewer than K orderings cannot rationalize most of the data, and thus encourages recovery of K preferences.

A basic bound on the size of the differentiation parameter is:

$$0 \leq d(\pi, \mu^*) \leq 1/K \tag{5}$$

for every π and μ^* . The lower bound was attained in Examples 3–5:

¹⁵ Observations (x, A) and (x', B) are in violation of the Independence of Irrelevant Alternatives (IIA) axiom if $x, x' \in B$ and $B \subseteq A$.

Observation 2. Fix any π and μ^* obeying the restrictions described in Example 3, 4, or 5. Then, $d(\pi, \mu^*) = 0$.

Generalizing from Example 3 in particular:

Definition 5.3. Let

$$\mathbb{C}(\mathcal{P}) = \left\{ (x, A) : x \text{ is } P\text{-maximal in } A \text{ for some } P \in \mathcal{P}, A \in 2^X \right\}$$

be the set of unique choice implications of preferences in \mathcal{P} .

Observation 3. If $\mathbb{C}(\mathcal{P})$ does not contain any K -violations of IIA, then $d(\pi, \mu^*) = 0$ for every sampling distribution π and every RUM μ^* supported on \mathcal{P} .

The upper bound $d(\pi, \mu^*) = 1/K$ is attained in the examples below¹⁶:

Example 6. The DM’s preferences are $\mathcal{P} = \{P_1, \dots, P_K\}$. Define x_i^* to be the P_i -maximal element from X , and suppose that all x_i^* are unique. Let every μ_A^* sample uniformly over \mathcal{P} , and assume that π only samples from choice sets containing $\{x_i^*\}_{i=1}^K$. Then, $d(\pi, \mu^*) = 1/K$.

Example 7. Consider $X = \{x_1, \dots, x_N\}$ and define $\mathcal{P} = \{P_1, P_2\}$ where

$$\begin{matrix} x_1 P_1 x_2 P_1 \dots P_1 x_N \\ x_N P_2 x_{N-1} P_2 \dots P_2 x_1 \end{matrix}$$

Let each μ_A^* assign equal probability to both preferences, and let π be an arbitrary distribution over non-singleton choice sets. Then, $d(\pi, \mu^*) = 1/2$.

The claim below says that so long as the differentiation parameter is strictly positive, then we can recover the number of orderings using any $\lambda > 1$ (which implements the Kalai et al., 2002 solution).

Claim 1. Suppose $d(\pi, \mu^*) > 0$ and $\mu = \mu^*$. Then,

$$\Pr(K_\lambda^*(D_n) = K) \rightarrow 1 \quad \text{as } n \rightarrow \infty$$

for every $\lambda > 1$.

The proof is clear and omitted. If the DM imperfectly maximizes, however, then $\lambda > 1$ will not generally be the best choice for recovery of K , and the size of the differentiation parameter $d(\pi, \mu^*)$ will be important. I turn to this case now.

¹⁶ In general, the value of $d(\pi, \mu^*)$ can vary significantly depending on the sampling distribution π . For example, as long as there is a single choice set for which all of the preferences disagree, then exclusive sampling of this choice set will result in $d(\pi, \mu^*) = 1/K$. Similarly, as long as there is a single choice set on which all of the preferences agree, then exclusive sampling of this choice set will result in $d(\pi, \mu^*) = 0$.

5.2. Main case: $\mu \neq \mu^*$

When there are no choice errors, all choice “inconsistencies” are (by assumption) due to preference heterogeneity, so the only obstacle to recovery is representation of all K preferences in the data. This identification problem is also present when there are choice errors. But in addition, the possibility of choice error introduces a second source of inconsistency: Choice errors may artificially inflate the inferred number of preferences (if we mistakenly interpret errors as preference), and they may also artificially reduce the inferred number of preferences (since an imperfectly maximized data set can in some cases be rationalized using fewer orderings than its perfectly maximized counterpart).

The proposed approach separates error from preference by looking for structure in the inconsistencies. A large set of choice inconsistencies that is “internally consistent”—i.e. rationalizable using the same preference—is indicative of preference. Inconsistencies which are “internally *inconsistent*” are taken to represent error. What a “large” set means is governed by choice of the parameter λ .

The main theorem provides values of λ under which the proposed approach recovers the true number of orderings as the number of observations grows large:

Theorem 1. Define

$$\bar{p} = d(\pi, \mu^*)(1 - p)^K. \quad (6)$$

Choose any $\tilde{p} \in (p, \bar{p})$ and set $\lambda = 1/(\tilde{p}n)$. Then,

$$\Pr(K_\lambda^*(D_n) = K) \rightarrow 1 \quad \text{as } n \rightarrow \infty.$$

The condition in (6) defines a value \bar{p} that is increasing in the differentiation parameter $d(\pi, \mu^*)$, decreasing in the probability of error p , and decreasing in the number of preferences K . Theorem 1 says that if each ordering ascribed to the DM is required to uniquely explain at least pn observations, but not more than $\bar{p}n$ observations, then the proposed approach will recover the number of underlying orderings given sufficiently many choice observations. Choosing $\lambda < 1/(\bar{p}n)$ may result in an underestimate of the number of preferences, and choosing $\lambda > 1/(pn)$ may result in an overestimate.

If either the differentiation parameter $d(\pi, \mu^*)$ is too small, or the probability of error p and number of preferences K too large, the condition in (6) will yield $\bar{p} < p$, in which case $K_\lambda^*(D)$ may not recover K for any value of λ . Specifically, it follows from the bounds on the differentiation parameter described in (5) that Condition (6) requires

$$\frac{1}{K} > \frac{p}{(1 - p)^K}.$$

For example, if the DM has more than ten underlying preferences, then the probability of error cannot exceed $p = 0.06$, and if the DM has 5 underlying preferences, then the probability of error cannot exceed $p = 0.11$.¹⁷

The proof of Theorem 1 is deferred to Appendix B.2, but a sketch of the main ideas follows. The key idea is to identify every data set with an undirected hypergraph¹⁸ (henceforth *graph*) in

¹⁷ I have not made efforts to optimize this bound, which can be improved in subsequent work.

¹⁸ A *hypergraph* is a generalization of a graph in which edges may connect more than two vertices.

the following way: every node corresponds to a choice observation, and there is an edge between a set of nodes if and only if the corresponding observations constitute a (minimal) set inconsistent with maximization of any single ordering. The proof notes that a data set is k -rationalizable if and only if the corresponding graph is k -colorable.^{19,20} Thus, the problem in (3) can be re-cast as finding the smallest number of colors k such that a large subset of nodes is k -colorable.

Let us consider a data set generated by repeated sampling from ν^* (which, recall, corresponds to choice without errors) instead of the actual distribution ν . Since by construction, this data set corresponds to *perfect* maximization of K orderings, the corresponding graph must admit a K -coloring. Moreover, since every set of observations in K -violation of PIA creates a complete K -partite subgraph, and at least one such set exists,²¹ the corresponding graph cannot be colored by fewer than K colors. The challenge is to show that even when the graph is perturbed by choice error, with high probability it will remain the case that a large subset of the nodes can be colored by K colors, but no fewer.

To show that K colors are sufficient to color most nodes, I use Hoeffding’s inequality to upper bound the number of imperfectly maximized choice observations by $1/\lambda$ (with high probability) as n gets large. To show that K colors are needed, I use (6) and Hoeffding’s inequality to lower bound the number of disjoint complete K -partite subgraphs by $1/\lambda$ (with high probability) as n gets large. This relies crucially on the differentiation parameter $d(\pi, \mu^*)$ being sufficiently large. Since each complete K -partite subgraph cannot be colored by fewer than K colors, the number of such subgraphs provides an approximate lower bound on the number of nodes that are uniquely colored by each of the first K colors. Thus, each of the first K orderings uniquely explains at least $1/\lambda$ observations, and the marginal $(K + 1)$ -st ordering explains strictly fewer than $1/\lambda$ additional observations, so the proposed approach correctly returns K orderings.

5.3. Evaluating assumptions

Since the number of orderings K , the probability of error p , and the differentiation parameter $d(\pi, \mu^*)$ are neither known nor observable, the expression in (6) cannot be directly determined from the data. Nevertheless, we can infer properties of these unknowns. Theorem 1 provides an *interval* of values of λ that recover the same solution. Thus, if the inferred $K_\lambda^*(D)$ is the “correct” number of underlying preferences, we can use the range of values of λ that induce this solution to bound $d(\pi, \mu^*)$ and p . To ease notation, d is used throughout this section in place of $d(\pi, \mu^*)$.

Formally, for each number of preferences k , define

$$\bar{\lambda}_k(D) = \max \{ \lambda' : K_{\lambda'}^*(D) = k \} \quad \underline{\lambda}_k(D) = \min \{ \lambda' : K_{\lambda'}^*(D) = k \} \tag{7}$$

to be the largest and smallest values of λ' that return the solution k . An implication of Theorem 1 is that eventually each choice of λ' in the interval $[1/(\bar{p}n), 1/(pn)]$ yields the solution K .²² An immediate corollary of Theorem 1 is:

¹⁹ A k -coloring of a graph is a partition of its vertex set V into k color classes such that no edge in E connects two nodes of the same color. A graph is k -colorable if it admits an k -coloring.

²⁰ This equivalence is shown by taking each color class to represent consistency with a distinct ordering.

²¹ This is implied by $d(\pi, \mu^*) > 0$. If $d(\pi, \mu^*) = 0$, then the interval $[1/(\bar{p}n), 1/(pn)]$ is empty, and the theorem holds vacuously.

²² Some care is needed here, since the size of n needed depends on the choice of λ' .

Corollary 1. Define $\underline{\lambda}_K(D_n)$ and $\bar{\lambda}_K(D_n)$ as in (7). Then, for every $\underline{x} > \frac{1}{d(1-p)^{K_n}}$ and $\bar{x} < \frac{1}{pn}$,

$$\Pr([\underline{x}, \bar{x}] \subseteq [\underline{\lambda}_K(D_n), \bar{\lambda}_K(D_n)]) \rightarrow 1 \quad \text{as } n \rightarrow \infty \tag{8}$$

The set of values in (8) is loosely an “inversion” of our main result: Theorem 1 says that if preferences are sufficiently differentiated and error is sufficiently small, then we can recover the number of preferences with an appropriate choice of λ . Corollary 1 asks, if K is the solution, how differentiated could the preferences have been, and how high must the probability of error have been?

Given some conjecture k , we can use (8) to back out implied properties of the (unobservable) primitives $d(\pi, \mu^*)$ and p . Each of $\bar{\lambda}_k := \bar{\lambda}_k(D_n)$, $\underline{\lambda}_k := \underline{\lambda}_k(D_n)$, and n can be computed directly from the data. Thus, if k is indeed the number of underlying preferences, then the interval

$$\left\{ (d', p') : \left[\frac{1}{d'(1-p')^k n}, \frac{1}{np'} \right] \subseteq [\underline{\lambda}_k, \bar{\lambda}_k] \right\} \tag{9}$$

will (eventually) approximately bound the possible values of d and p . Notice also that every d' and p' in the set (9) satisfy

$$\frac{1}{\bar{\lambda}_k(1 - 1/(\bar{\lambda}_k n))^k n} \leq d' \leq \frac{1}{\underline{\lambda}_k(1 - 1/(\underline{\lambda}_k n))^k n} \tag{10}$$

$$\frac{1}{\bar{\lambda}_k n} \leq p' \leq \frac{1}{\underline{\lambda}_k n} \tag{11}$$

Intuitively, when the differentiation parameter d is large, and the probability of error p is small, then we expect the interval $\bar{\lambda}_k - \underline{\lambda}_k$ (over which k is recovered) to be large. If $\bar{\lambda}_k - \underline{\lambda}_k$ is small, then either k is not the underlying number of preferences, or in fact d is small and p is large (in which case recovery may not be meaningful). These bounds are applied in Section 7 to help arbitrate between different solutions.

5.4. Testing rationality ($K = 1$)

The main part of this paper seeks to recover the “true” number of underlying preferences, but an important special case (with long precedence in the literature) considers whether the choice data should be rationalized using a single preference.

Suppose that $K = 1$ in the proposed framework; then, the differentiation parameter $d(\pi, \mu^*)$ is trivially 1, since every choice observation constitutes a 1-violation of IIA. An immediate corollary of Theorem 1 is:

Corollary 2. Set any λ satisfying $0 \leq \lambda \leq 1/(pn)$. Then, $\Pr(K_\lambda^*(D_n) = 1) \rightarrow 1$ as $n \rightarrow \infty$.

It is interesting also to consider choice models that *don't* naturally correspond to maximization of a single preference, and see what the proposed approach recovers in those cases. As one such test, I apply the proposed approach to a class of logit choice rules.

Example 8 (Single preference with logit error). Let $X = \{x_1, x_2, x_3\}$. The DM imperfectly maximizes the preference ordering $x_3 P x_2 P x_1$. Specifically, his probability of choosing alternative x from choice set A is given by

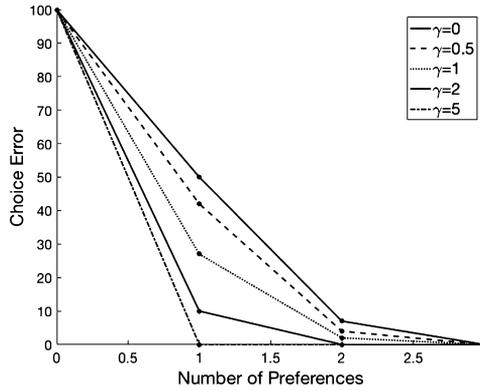


Fig. 3. Error-preference tradeoff graph for various choices of γ in Example 8 (assuming a uniform sampling rule over non-singleton choice sets, and simulating 100 choice observations). The set of choices of λ that recover $K = 1$ is larger for the curves corresponding to higher values of γ .

$$c(x|A) = \frac{e^{\gamma u(x)}}{\sum_{x' \in A} e^{\gamma u(x')}} \quad \forall x \in A,$$

where $\gamma \geq 0$ and $u(x_k) = k$ assigns a higher payoff to alternatives that are higher ranked by P .

Within this class of choice rules, the probability of choosing the most preferred outcome is governed by the logit parameter γ . Lower choices of γ return higher probabilities of error, with the extreme case $\gamma = 0$ corresponding to uniform selection over the available alternatives. Intuitively, this choice rule is better described as (imperfect) maximization of a single preference when γ is large. Claims 2 presents two senses in which the proposed solution aligns with this intuition.

Claim 2. Let D_n^γ be a random data set of n observations when the DM uses the choice rule in Example 8 with logit parameter γ , and choice sets are drawn uniformly at random (excluding singleton choice sets). Then:

- (a) for every n , the expected size of $\bar{\lambda}_1(D_n^\gamma) - \underline{\lambda}_1(D_n^\gamma)$ is increasing in γ .²³
- (b) for each choice of λ and quantity of data n , the expected value of $K_\lambda^*(D_n^\gamma)$ is weakly decreasing in γ .

Part (a) says that the size of the set of values of λ that recover $K = 1$ is monotonically increasing in γ , so that the more concentrated choice behaviors are, the more “slack” there is in the solution. This claim is illustrated in Fig. 3.

Part (b) provides a complementary result: for every fixed value of λ , the recovered number of preferences is (weakly) decreasing with γ . Thus, the expected value of $K_\lambda^*(D_n^\gamma)$ is largest for $\gamma = 0$, again suggesting that $\gamma = 0$ is the choice model in the class above that is most different from maximization of a single ordering.

²³ Note that $\underline{\lambda}_1(D_n) = 0$ for every data set D_n , so a simpler statement of this result says that the expected size of $\bar{\lambda}_1(D_n)$ is increasing in γ .

Another case conceptually distinct from maximization of a single ordering is maximization of multiple preferences. As one such example, we return to the setting from Section 2:

Example 9 (*Two preferences with logit error*). Depending on the choice set, the DM applies either of two logit choice rules, each of which imperfectly maximizes a single preference. For the idealized choice data set described in Section 2, the solution $K_\lambda^*(D_n) = 1$ is recovered for all values of $\lambda \leq 1/(0.46n)$. This interval is in fact smaller than the interval of values of λ that recover $K_\lambda^*(D_n) = 1$ for uniformly random choices ($\gamma = 0$ in Example 8), suggesting that the choice data described in Section 2 is not indicative of a single underlying preference.

6. Extensions: continuous utility

So far, we have considered a decision maker whose preferences are orderings over a discrete set X . I now show that the main results extend to the case in which (X, τ) is a topological space.

Formally, suppose that choice sets are compact subsets $A \subseteq X$, repeatedly sampled according to a distribution π .²⁴ The set $\mathcal{U} = \{u_\theta\}_{\theta \in \Theta}$ is a parametric family of continuous utility functions $u_\theta : X \rightarrow \mathbb{R}$. Conditional on observation of choice set A , the DM maximizes a utility function u_θ , where θ is sampled from a (Borel-measurable) distribution $\mu_A \in \Delta(\Theta)$. Write D for a typical outcome of the choice data.

As before, I assume that the DM possesses an underlying “sparse” set of utility functions. Formally, for each choice set A , there is a map $g_A : \Theta \rightarrow \Delta(\Theta)$ such that $g_A\theta \geq 1 - p$ for each A and θ . I assume that each μ_A can be rewritten as

$$\mu_A = \mu_A^* G_A$$

where G_A is the Markov kernel on (Θ, \mathcal{B}) associated with g_A , and μ_A^* is supported on a finite set of utility functions $\mathcal{U} \subset \mathcal{U}$. The goal is to determine the number of utility functions $K := |\mathcal{U}|$.²⁵

The proposed approach minimizes a weighted average of the number of inferred utility functions and the number of unexplained observations. For every set of utility functions \mathcal{U} , let

$$\varepsilon(D, \mathcal{U}) = \# \left\{ (x, A) \in D : x \neq \max_{x' \in A} u(x') \text{ for any } u \in \mathcal{U} \right\}$$

be the number of choice observations in D that are not consistent with maximization of any utility function in \mathcal{U} . Then,

$$\varepsilon_k(D) = \min_{|\mathcal{U}|=k} \varepsilon(D, \mathcal{U})$$

is the minimal number of observations in D that are unexplained if we rationalize the DM’s choices using k utility functions.

The solution below simultaneously minimizes the number of utility functions k and the implied choice error $\varepsilon_k(D)$:

²⁴ Some care is required in specifying the correct σ -algebra over choice sets; for example, one can take the Borel σ -algebra associated with the product topology of τ .

²⁵ Observe that as before, no parametric assumptions are made regarding the distribution of error; future work may include such assumptions to strengthen the recovery results.

Definition 6.1. For every $\lambda \in \mathbb{R}_+$, define

$$\underline{K}_\lambda^*(\underline{D}) = \operatorname{argmin}_{k \in \mathbb{Z}_+} [k + \lambda \underline{\varepsilon}_k(\underline{D})]. \tag{12}$$

As before, when there are multiple solutions, take $\underline{K}_\lambda^*(\underline{D})$ to mean the minimal value.

Corollary 3 below shows that this solution recovers the “correct” number of utility functions K under conditions that directly parallel the previous section. The statement below follows as a corollary to Theorem 1 (where the differentiation parameter $d(\pi, \mu^*)$ is defined as in Section 3.2):

Corollary 3. Define

$$\bar{p} = d(\pi, \mu^*)(1 - p)^K. \tag{13}$$

Choose any $\tilde{p} \in [p, \bar{p}]$ and set $\lambda = 1/(\tilde{p}n)$. Then,

$$\Pr(K_\lambda^*(\underline{D}_n) = K) \rightarrow 1 \quad \text{as } n \rightarrow \infty.$$

Thus, the proposed approach recovers the number of underlying utility functions as the number of observed choices gets large.

Why do the conditions of Theorem 1 extend to this more general setting? The key observation is that choice data generated in this way can be mapped into discrete choice data, where we reduce \underline{X} to the finite set

$$X = \{x \in \underline{X} : (x, \underline{A}) \in \underline{D} \text{ for some } \underline{A} \subseteq \underline{X}\}.$$

This set consists of all choice alternatives that are observed to be chosen. For example, take $\underline{X} = \mathbb{R}$, and suppose we observe

$$\underline{D} := \{(3, [0, 4]), (2, [1, 4]), (8, [0, 10])\}.$$

Then, labelling ‘3’ as x_1 , ‘2’ as x_2 , and ‘8’ as x_3 , we can redefine the set of choice alternatives as $X = \{x_1, x_2, x_3\}$, and the choice data as

$$D := \{(x_1, \{x_1, x_2\}), (x_2, \{x_1, x_2\}), (x_3, \{x_1, x_2, x_3\})\}.$$

This is a standard mapping in the literature, and yields a data set of the form introduced in Section 3.1.

A lemma in Appendix B.3 shows that the new problem posed in Definition 6.1 is equivalent to the original problem posed in Definition 4.1, in the sense that the solution to

$$\operatorname{argmin}_{k \in \mathbb{Z}_+} [k + \lambda \underline{\varepsilon}_k(\underline{D})],$$

is the same as the solution to

$$\operatorname{argmin}_{k \in \mathbb{Z}_+} [k + \lambda \varepsilon_k(D)].$$

It immediately follows that the conditions for recovery stated in the previous section are also the conditions needed in the present setting.

7. Application

This section describes an example application of the proposed approach, which builds on an analysis from Crawford and Pendakur (2012) (henceforth CP). CP study the consumption decisions of Danish households over six different kinds of milk, aggregated over a month. The relevant choice information is the quantity of each kind of milk purchased during this time (written as a quantity vector $\mathbf{q} \in \mathbb{R}^6$), and the price index at which these purchases were made (written as a price vector $\mathbf{p} \in \mathbb{R}^6$). The main sample in CP consists of 500 households, so the choice data is $\{(\mathbf{p}_1, \mathbf{q}_1), \dots, (\mathbf{p}_{500}, \mathbf{q}_{500})\}$.²⁶

We can map these choice observations into the present framework using the relabelling described in Section 6. That is, index the observations by $i = 1, \dots, 500$, and define $\mathbf{x}_i = (\mathbf{p}_i, \mathbf{q}_i)$. Take $X = \{\mathbf{x}_i\}_{i=1}^{500}$ to be the set of choice alternatives. For each observation i , let

$$A_i = \{\mathbf{x}_j : \mathbf{p}_i \cdot \mathbf{q}_i \geq \mathbf{p}_j \cdot \mathbf{q}_j\}$$

consist of every alternative \mathbf{x}_j that is less costly than the selected alternative \mathbf{x}_i . These are the alternatives in X that could have been chosen when the alternative \mathbf{x}_i was chosen.²⁷ The observed data from CP is now rewritten

$$D = \{(\mathbf{x}_1, A_1), \dots, (\mathbf{x}_{500}, A_{500})\}.$$

This data set is equivalent to the original data in the sense described in the previous section.²⁸

There are many ways to rationalize the choice data D . For example, following the proposal of Houtman and Maks (1985), we could find the single ordering that explains the largest fraction of the data, and interpret the remaining observations as choice error. CP find that no single preference explains more than two-thirds of the observations.

Alternatively—and this is the main approach taken in CP—we can seek the minimal set of preferences that explains every observation (thus following the proposal of Kalai et al. (2002)). CP find that no more than five orderings are needed to perfectly rationalize the data.

The present paper interprets the two solutions above as edge cases among a set of rationalizations of the data, each of which entails a different tradeoff between maximization of fit to the data and minimization of the number of preferences used. Fig. 4 provides approximations from CP for the number of unexplained observations $\varepsilon_k(D)$ (for the purpose of illustration of the approach, I will treat these approximations as exact).^{29,30} For example, with a single preference, we must leave 179 (of the 500) observations unexplained; using two preferences, we must leave 79 observations unexplained; and with five preferences, we can explain all of the observations.

²⁶ Their data also includes a household indicator and covariates describing each household, but these are outside of the proposed framework.

²⁷ That is, if $\mathbf{p}_i \cdot \mathbf{q}_i \geq \mathbf{p}_j \cdot \mathbf{q}_j$, then $\mathbf{x}_j \in A_i$.

²⁸ Suppose there exists a set of k utility functions such that m observations in the original data are consistent with maximization of a utility function from this set; that is,

$$u(\mathbf{q}_i) > u(\mathbf{q})$$

for all quantity vectors \mathbf{q} satisfying $\mathbf{p}_i^T \mathbf{q}_i \geq \mathbf{p}_i^T \mathbf{q}$. Then we can find a set of k preference orderings on X such that m observations in the relabelled data are consistent with maximization of a preference from this set, and vice versa.

²⁹ These values correspond to an algorithm that computes an upper bound on the number of types needed to explain a given number of observations, so the true values of $\varepsilon_k(D)$ are weakly smaller than those reported below. Crawford and Pendakur (2012) provide also lower bounds that show that the errors in this approximation are not large.

³⁰ From CP: $\varepsilon_1(D) = 179$, $\varepsilon_2(D) = 79$, $\varepsilon_3(D) = 26$, $\varepsilon_4(D) = 8$, and $\varepsilon_5(D) = 0$.

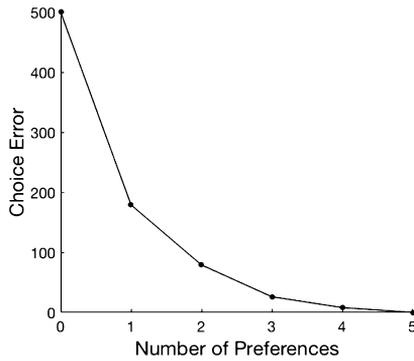


Fig. 4. Error-preference tradeoff graph for the Crawford and Pendakur (2012) data set.

Table 1
Different solutions k (first row) are induced by different intervals of λ .

1	2	3	4	5
$\lambda \in [0, 0.01]$	$[0.01, 0.019]$	$[0.019, 0.056]$	$[0.056, 0.125]$	$[0.125, \infty)$

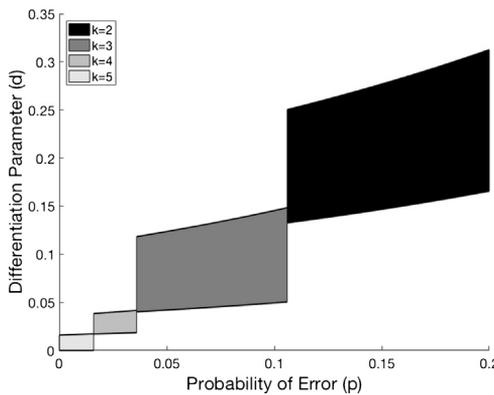


Fig. 5. Different potential solutions $k \in \{2, 3, 4, 5\}$ correspond to different sets of $(p, d(\pi, \mu^*))$ pairs.

Given this data set, there are five possible solutions to (3), each of which holds for a range of choices of λ , as shown in Table 1.

One rule-of-thumb approach is to set $\lambda = 1/(\tilde{p}n)$, where \tilde{p} is a slight overestimate of the probability of error. For example, if the analyst believes that subjects have approximately a 5% probability of error, then the proposed approach uses $\lambda = 1/(0.05 \times 500) = 0.04$ and concludes $K = 3$.

We can also take a more agnostic approach following the discussion in Section 5.3. For each potential solution k , the expression in (9) delivers bounds on the differentiation parameter and probability of error. Fig. 5 plots the sets of $(p, d(\pi, \mu^*))$ values that corresponds to different solutions k .³¹

³¹ The set of values for $k = 1$ can be similarly computed but is not shown in Fig. 5.

Table 2
 Bounds for the differentiation parameter $d(\pi, \mu^*)$ and probability of error p .

k	$d(\pi, \mu^*)$		p	
	Lower bound	Upper bound	Lower bound	Upper bound
1	0.25	1	0.2	1
2	0.13	0.31	0.106	0.2
3	0.04	0.15	0.036	0.106
4	0.017	0.04	0.016	0.036
5	0	0.017	0	0.016

We can additionally use (10) and (11) to bound the probability of error p and the differentiation parameter d , starting from different conjectured solutions k . Table 2 reports these bounds.

Notice that the regions in Fig. 5 corresponding to the solutions $k = 4$ and $k = 5$ are quite small. Specifically, for the candidate solution $k = 5$, the implied bounds are $p \in [0, 0.016]$ and $d(\pi, \mu^*) \in [0, 0.017]$, so that no more than 1.6% of choice observations are in error, and no more than 1.7% of choice observations provide evidence of five orderings. For the candidate solution $k = 4$, the implied bounds are $p \in [0.016, 0.036]$ and $d(\pi, \mu^*) \in [0.017, 0.04]$, so that no more than 3.6% of observations are in error, and no more than 4% of choice observations provide evidence of four orderings. To the extent that an analyst considers these restrictions too severe, this analysis suggests that the solutions $k = 4$ and $k = 5$ should be viewed cautiously. Specifically, the final two Kalai et al. (2002) preferences recovered in Crawford and Pendakur (2012) may be better understood as capturing choice errors.

8. Can we recover more?

Section 5 described conditions under which the problem in (3) recovers the number of underlying preferences with high probability. This section now asks whether it is possible to recover the preferences themselves.

8.1. Non-identifiability of sets of preferences

Say that the set of preferences \mathcal{P} is *identifiable* if there is at least one data set D such that \mathcal{P} is the unique set of $|\mathcal{P}|$ preferences (or fewer) that perfectly explains every observation in D .

Definition 8.1. The set of orderings \mathcal{P} is *identifiable* if there exists some data set D such that $\varepsilon(D, \mathcal{P}) = 0$, and moreover $\varepsilon(D, \mathcal{P}') > 0$ for every $\mathcal{P}' \neq \mathcal{P}$ with $|\mathcal{P}'| \leq |\mathcal{P}|$.

The following proposition says that most sets of orderings are not identifiable.

Proposition 1. No set of orderings \mathcal{P} with $|\mathcal{P}| \geq 3$ is identifiable. Suppose $\mathcal{P} = \{P_1, P_2\}$; then, \mathcal{P} is identifiable if and only if the P_1 -maximal alternative in X is the P_2 -minimal alternative in X , and vice versa. Every singleton set $\mathcal{P} = \{P\}$ is identifiable.

The difficulty of recovery is not specific to the nonparametric nature of the exercise, but to basic issues concerning identifiability for multiple preferences. Consider any set that includes a pair of preferences P_1, P_2 , where some alternative x_1 is ranked first under one ordering (say, P_1) and not ranked last under the other; for example,

$$\begin{matrix} x_1 P_1 x_2 P_1 x_3 \\ x_2 P_2 x_1 P_2 x_3 \end{matrix}$$

We can construct then a new preference P'_2 based on P_2 , with the single difference that x_1 is ranked last:

$$\begin{matrix} x_1 P_1 x_2 P_1 x_3 \\ x_2 P'_2 x_3 P'_2 x_1 \end{matrix}$$

Every choice that can be rationalized using a preference from $\{P_1, P_2\}$ can also be rationalized using a preference from the set $\{P_1, P'_2\}$; this is easily verified for the example above, and the proof of Proposition 1 shows that this follows more generally. Thus, $\{P_1, P'_2\}$ will be a solution whenever $\{P_1, P_2\}$ is.

The obstacle to recovery is that sets of preferences differ in the “richness” of their choice implications. A set of preferences can (strictly) encompass all the choice implications of another set with the same number of preferences. Thus, any approach that penalizes only the *size* of a set of orderings, as both the proposed approach in (3) and the approach suggested in Kalai et al. (2002) do, will be biased towards elicitation of sets with richer choice implications.

I outline below two paths forward. Section 8.2 suggests a modification of the proposed approach which allows for recovery of the *choice implications* of the set of preferences. This is generally understood to be the real content of a set of preferences. Section 8.3 considers a richer kind of data set, which includes auxiliary information on the choice contexts active during different observations. I show that with this additional information, we can (under certain conditions) recover the set of preferences.

8.2. Recovery of choice implications

One approach is to change the “complexity” penalty from *number of preferences* to *number of choice implications*. Formally, define the function $l : \mathcal{P} \rightarrow \mathbb{Z}_+$ such that $l(\mathcal{P})$ counts the number of unique choice observations (x, A) that are consistent with \mathcal{P} ; that is,

$$l(\mathcal{P}) = \#\{(x, A) : x \text{ is } P\text{-maximal in } A \text{ for some } P \in \mathcal{P}\}.$$

The following definition modifies the proposed approach in Definition 4.1 by replacing the metric $|\mathcal{P}|$ with the metric $l(\mathcal{P})$. This solution minimizes a weighted average of the number of choice implications and the number of unexplained choice observations.

Definition 8.2. For every $\lambda \in \mathbb{R}_+$, define

$$\mathcal{P}^*_\lambda(D) = \operatorname{argmin}_{\mathcal{P} \in \mathcal{P}} [l(\mathcal{P}) + \lambda \varepsilon(D, \mathcal{P})]. \tag{14}$$

In this case, $1/\lambda$ can be understood as the cost of each new choice implication, so that a choice implication is attributed to the DM only if it appears at least $1/\lambda$ times in the data. As $\lambda \rightarrow 0$, the cost of errors becomes increasingly small relative to the cost of new choice implications, so that the analyst prefers to attribute to the DM as few choice implications as possible. As $\lambda \rightarrow \infty$, the cost of choice errors becomes increasingly large relative to the cost of new choice implications, so that the analyst prefers to ascribe to the DM as many choice implications as is necessary to perfectly rationalize the data.

The proposition below says that for certain choices of λ , we can recover the choice implications of \mathcal{P} (denoted $\mathbb{C}(\mathcal{P})$ as in Definition 5.3) with probability arbitrarily close to 1 as the quantity of data increases.

Proposition 2. Define α to be the smallest nonzero frequency with which any choice observation occurs under v^* .³² Choose any $\tilde{p} \in (p, \alpha(1 - p))$ and set $\lambda = 1/(\tilde{p}n)$. Then,

$$\Pr(\mathbb{C}(\mathcal{P}_\lambda^*(D_n)) = \mathbb{C}(\mathcal{P})) \rightarrow 1 \quad \text{as } n \rightarrow \infty.$$

Above, the cost $1/\lambda$ of recovering an additional choice observation is chosen to satisfy $1/\lambda = \tilde{p}n < \alpha(1 - p)n$. Thus (with high probability for large data sets) every choice implication $(x, A) \in \mathbb{C}(\mathcal{P})$ is observed sufficiently many times to not be mistaken as error. To see that no choice implication $(x, A) \notin \mathbb{C}(\mathcal{P})$ will be incorrectly recovered, observe that the number of instances of $(x, A) \notin \mathbb{C}(\mathcal{P})$ in the data is upper bounded by the number of choice errors. With high probability, the number of choice errors concentrates below $pn < 1/\lambda$, from which it follows that $\mathbb{C}(\mathcal{P}_\lambda^*)$ will include as few choice implications outside of $\mathbb{C}(\mathcal{P})$ as possible.

8.3. Auxiliary data on contexts

As an alternative approach for recovery of preference, we may turn to richer data sets. Specifically, suppose that there is a set of *observable* choice contexts $\mathcal{C} = \{1, \dots, M\}$. Each context is associated with a preference, and multiple contexts may be associated with the same preference. (For example, suppose that choice data is aggregated over various kinds of financial decisions, and the contexts correspond to different financial domains. Individuals may have the same risk preference over all types of insurance decisions, but a different risk preference over 401(k) savings.)

Formally, there is an unknown map

$$m : \mathcal{C} \rightarrow \mathcal{P},$$

which assigns each choice context to a preference. The primary case of interest is when the image of m is a small set of preferences.

For each choice set A , let the empirical distribution over contexts be given by $\phi_A \in \Delta(\mathcal{C})$. The RUM μ_A^* introduced in Section 3.3 is micro-founded as

$$\mu_A^*(P) = \phi_A(m^{-1}(P)).$$

That is, the probability with which P is sampled (in the absence of choice error) is the probability that a context emerges which cues preference P . As before, there is a sampling distribution $\pi \in \Delta(2^X)$ over choice sets. In the absence of choice error, we can write

$$v^*((x, A), C) = \begin{cases} \pi(A)\phi_A(C) & \text{if } x \text{ is } m(C)\text{-optimal in } A \\ 0 & \text{otherwise} \end{cases}$$

for the probability that choice observation (x, A) and context C are observed. Notice that in this case, the outcome x is deterministic conditional on the choice context C and choice set A .

In the main model, observations are instead sampled from

$$v((x, A), C) = \pi(A)\phi_A(C)q_C(x|A)$$

where each $q_C(\cdot|A) \in \Delta(A)$ is a distribution over choice alternatives, associated with the context C . Assume that each $q_C(\cdot|A)$ assigns probability at least $1 - p$ to the $m(C)$ -optimal alternative in A .

³² That is, let $\alpha = \min_{(x,A) : v^*(x,A) > 0} v^*(x, A)$.

Example 10. There are four kinds of subjects: males over the age of 65, males under the age of 65, females over the age of 65, and females under the age of 65. These categories are indexed (in that order) using $\mathcal{C} = \{1, 2, 3, 4\}$. Only age determines preference, so that $m(1) = m(3) = P$ and $m(2) = m(4) = P'$. The analyst does not know this, however. He observes tuples $((x, A), C)$ indicating that alternative x was chosen from choice set A by a subject from category $C \in \mathcal{C}$. The goal is to back out from the data that there are only two active preferences, and to determine which these are.

We can modify the approach in Section 4 to search for assignments of preferences to contexts. Define

$$\varepsilon(D, \tilde{m}) = \#\{((x, A), C) \in D : x \text{ is not } \tilde{m}(C)\text{-maximal in } A\}$$

to be the number of implied choice errors when \tilde{m} is the mapping from contexts to preferences, and define

$$m_\lambda^*(D) := \operatorname{argmin}_{\tilde{m}: \mathcal{C} \rightarrow \mathcal{P}} [|\tilde{m}(\mathcal{C})| + \lambda \varepsilon(D, \tilde{m})] \tag{15}$$

where $|\tilde{m}(\mathcal{C})|$ is the number of (unique) preferences assigned under \tilde{m} . Thus, $m_\lambda^*(D)$ is the assignment of preferences to contexts that minimizes the number of distinct preferences, and also the associated number of choice errors.

In the proposition below, let α be the smallest nonzero frequency with which any observation $((x, A), C)$ occurs under ν^* .³³

Proposition 3. Choose any $\lambda = 1/(\tilde{p}n)$ where $\tilde{p} \in (p, \alpha(1 - p)/2)$. Then,

$$\Pr(m_\lambda^*(D_n) = m) \rightarrow 1 \quad \text{as } n \rightarrow \infty.$$

Thus, recovery of the exact set of preferences is possible if there is auxiliary information about choice contexts, and sufficient observation of the choice implications for each choice context.

9. Related literature

9.1. Nonparametric preference recovery

This paper builds on a literature regarding nonparametric identification of multiple preferences from choice data. Most directly, it extends Kalai et al. (2002), which defines a set of orderings \mathcal{P} to be a *rationalization by multiple rationales* if for every observation (x, A) , the choice alternative x is P -maximal in A for some ordering $P \in \mathcal{P}$. Kalai et al. (2002) search for the smallest L such that some set \mathcal{P} with $|\mathcal{P}| = L$ is a rationalization by multiple rationales. Using the terminology of this paper, any such set of preferences \mathcal{P} is a perfect rationalization of the data, but it may not correspond to a “best” rationalization of the data as defined in (3). In particular, I suggest that the analyst may prefer an imperfect rationalization of the data using some $K < L$ orderings. The key conceptual difference is that Kalai et al. (2002) is agnostic towards the degree of evidence for orderings, whereas the approach in this paper insists on sufficient evidence for each ordering in order to separate error from preference variation.

³³ That is, $\alpha = \min_{C \in \mathcal{C}, (x, A) \in \mathcal{C}(\{m(C)\})} \nu^*((x, A), C)$.

Ambrus and Rozen (2013) study multiple-preference models in which choice is determined through maximization of a choice-set independent aggregation rule over preferences. They find that without prior restriction on the number of selves involved in a decision, many multiple-preference (“multi-self”) models have no testable implications. Although the class of models considered in their paper is different from the class studied in the present paper,³⁴ their lesson that restricting the number of preferences is necessary for recovery holds here as well (relating especially to the results in Section 8.1), and motivates in part the suggested criterion in (3).

Other nonparametric approaches for preference identification include Houtman and Maks (1985) and Varian (1982). These approaches differ from the present paper, and from Kalai et al. (2002), in finding a *single* best-fit ordering. A separate literature studies related questions under different parametric assumptions—see, for example, the foundational work of Quandt (1956), McFadden and Richter (1970), and Train (1986). Finally, Crawford and Pendakur (2012) and Dean and Martin (2010) apply the approaches of Kalai et al. (2002) and Houtman and Maks (1985) towards recovery of preferences from real choice data.³⁵

9.2. Testing rationality

When choice data is *inconsistent*—meaning that it is incompatible with perfect rationalization by a single ordering—how should we measure the inconsistency of the observed choices? Solutions have been proposed by Afriat (1967), Varian (1982), Echenique et al. (2011), Houtman and Maks (1985), Gross (1989), Famulari (1995), Apesteguia and Ballester (2012), and Dean and Martin (2016) among others. See Apesteguia and Ballester (2012) for a summary and comparison of these approaches.

In view of this literature, one goal of the present paper is to distinguish between choice data that is inconsistent because of choice error, and choice data that is inconsistent because of multiplicity in preference. These two sources for error are confounded in many of the measures above.

The proposed approach offers a new perspective on this question. If the choice data D is generated via approximately perfect rationalization of multiple preferences, then the recovered number of preferences $K_{\lambda}^*(D)$ will exceed 1 for most choices of λ (and the error-preference tradeoff graph will resemble Fig. 1). Otherwise, $K_{\lambda}^*(D) = 1$ for many choices of λ and the error-preference tradeoff graph can be expected to resemble Fig. 3; see Section 5.4 for an extended discussion.

Finally, Halevy et al. (2015) proposes a novel decomposition of error into two sources: *choice inconsistency* and *preference misspecification*. The present paper is related to these ideas at a high level; in particular, the proposed approach also provides a decomposition of inconsistency. The sources of inconsistency considered in the two papers are different, however. For example, Halevy et al. (2015) considers a *single* (continuous) utility function and measures misspecification that arises from parametric restrictions, while the present paper takes a nonparametric approach and considers multiple preferences.

³⁴ In the present paper, the aggregation rule varies across choice problems.

³⁵ See Deb (2009) and Dean and Martin (2010) for computationally efficient approaches for approximating the Kalai et al. (2002) solution.

10. Conclusion

Inconsistencies in choice data may emerge *both* from (unintentional) choice error and also from (intentional) maximization of different preferences. Classic approaches such as Houtman and Maks (1985) and Kalai et al. (2002) focus on either of these sources of error, but—for reasons of welfare evaluation and out-of-sample prediction—we may prefer interpretations of the data that accommodate both.

This paper proposes identification of underlying “structural” preferences that are maximized in the data. The proposed approach looks for the multiple-preference rationalization of the data that simultaneously minimizes the number of preferences and also the number of unexplained observations. Different choices of tradeoffs between these objectives yield different solutions, and the main results relate the optimal choice to primitives of an underlying choice model.

Some of the techniques proposed in this paper may be applied to choice data generated by other choice models as well. For example, consider a DM who has a single partial ordering, and chooses uniformly from those alternatives (potentially multiple) that he most prefers. Such a model is outside the scope of the current paper. Nevertheless, we can differentiate between various kinds of indifference based on the error-preference tradeoff graph proposed in this paper. If the DM is only indifferent between low-ranked alternatives, then (with high probability for large data sets) the error-preference tradeoff graph will take a particular shape—similar to Fig. 3, and different from Fig. 1. In contrast, if the DM is indifferent between highly ranked alternatives, then we may expect the error-preference tradeoff graph to more closely resemble Fig. 1. Thus, the basic analysis of choice data proposed in this paper may additionally be useful towards other goals.

Appendix A. Out-of-sample prediction accuracy

One reason that recovery of the number of underlying structural preferences is important is because it affects prediction of unobserved choice behaviors. Application of Houtman and Maks (1985) towards this goal can underfit the data, missing important structural preferences, and application of Kalai et al. (2002) can overfit the data, interpreting errors as preference. These problems can result in substantial reductions in prediction accuracy, which I demonstrate in the examples below. Since out-of-sample predictions are of interest, these examples follow the extension described in Section 8.3, where choice contexts are observed, and the goal is to correctly assign contexts to preferences.

In both examples, a *training set* (of 20 choice observations) is generated from a fixed choice rule, and different approaches are applied to this data. I then generate a new *test set* (again of 20 choice observations) on which to evaluate the estimated models.³⁶ Prediction accuracy is measured as the fraction of test observations for which the chosen alternative is correctly predicted. This procedure is repeated ten times (with new training and test data), and I report an average of the out-of-sample prediction accuracies.³⁷

³⁶ That is, new observations (x, A) are generated, and the estimated models are applied to predict the chosen alternative x given the choice set A .

³⁷ To keep these examples as simple as possible, in both examples *one of* Houtman and Maks (1985) or Kalai et al. (2002) does achieve the performance of the proposed approach. More complex examples can be constructed in which both perform poorly out-of-sample (see e.g. the example in Section 2).

Example 11. There are four choice alternatives $X = \{x_1, x_2, x_3, x_4\}$. Choice sets $A \subseteq X$ are generated uniformly at random (excluding singleton choice sets). There are two choice contexts, indexed $\mathcal{C} = \{1, 2\}$, both observable.

The DM maximizes a different choice rule in each context. In the first, his probability of choosing alternative x from choice set A is given by

$$c(x|A) = \frac{e^{\gamma u(x)}}{\sum_{x' \in A} e^{\gamma u(x')}}$$

where $\gamma = 10$ and $u(x_k) = k$ (so that higher indexed alternatives are more preferred). In the second, his probability of choosing alternative x from choice set A is given by the same expression but setting $u(x_k) = 5 - k$ (so that higher indexed alternatives are less preferred).

Consider two approaches for making out-of-sample predictions: First, following Houtman and Maks (1985), find the single preference that maximizes the number of rationalized choice observations. Second, following the proposed approach in this paper, identify the solution that solves (15). Out-of-sample performance accuracies are reported and compared in the table below. For robustness, I show the prediction accuracies for a couple of different choices for the tradeoff parameter λ .

	Prediction accuracy
Houtman and Maks (1985)	47%
Proposed approach using $\lambda = 1/(0.1n)$	95%
Proposed approach using $\lambda = 1/(0.2n)$	93%
Proposed approach using $\lambda = 1/(0.3n)$	88%

The proposed approach with choice of $\lambda = 1/(0.1n)$ improves upon Houtman and Maks (1985) by 48% (and the other choices of λ also lead to substantial improvements).

Example 12. Again, there are four choice alternatives $X = \{x_1, x_2, x_3, x_4\}$, and choice sets $A \subseteq X$ are generated uniformly at random (excluding singleton choice sets). There are five choice contexts indexed $\mathcal{C} = \{1, 2, \dots, 5\}$, all observable. In each context, the DM uses the same choice rule in which the probability of choosing alternative x from choice set A is given by

$$c(x|A) = \frac{e^{\gamma u(x)}}{\sum_{x' \in A} e^{\gamma u(x')}}$$

with $\gamma = 1.5$ and $u(x_k) = k$ (so that higher indexed alternatives are more preferred).

I consider two approaches for making out-of-sample predictions. First, following Kalai et al. (2002), I find an assignment of preferences to contexts that minimizes the number of unexplained choice observations.^{38,39} Second, following the proposed approach in this paper, I identify the solution that solves (15). The table below reports and compares out-of-sample performance accuracies for these two approaches. As above, prediction accuracies are shown for a couple of different choices for the tradeoff parameter λ .

³⁸ Note that unlike in the usual setting in which Kalai et al. (2002) is applied, here it may not be possible to achieve a perfect rationalization, because we require preferences to be constant within choice contexts.

³⁹ When there are multiple assignments that achieve the minimal error, I select from these uniformly at random.

	Prediction accuracy
Kalai et al. (2002)	64%
Proposed approach using $\lambda = 1/(0.1n)$	77%
Proposed approach using $\lambda = 1/(0.2n)$	83%
Proposed approach using $\lambda = 1/(0.3n)$	84%

Again, the proposed approach leads to significant improvements in predictive accuracy (up to 20%).

Appendix B. Proofs from main text

B.1. Preliminaries

Following, I collect definitions and results that are used in the proofs of Theorem 1 and Proposition 1.

Let $G : D \mapsto G_D$ be a map that identifies every data set D with a (hyper-) graph $G_D = (V_D, E_D)$, where $V_D = \{1, \dots, n\}$ indexes choice observations, and E_D consists of every set $T \subseteq V_D$ with the property that: (1) the observations $\{(x_i, A_i)\}_{i \in T}$ are not 1-rationalizable, and (2) every proper subset of $\{(x_i, A_i)\}_{i \in T}$ is 1-rationalizable. These concepts are related to our problem as follows.

Observation 4. D is k -rationalizable $\iff G_D$ is k -colorable.

Take each color class to represent consistency with a distinct ordering, and the equivalence follows directly.

This observation further implies that ε_k is the minimum number of vertices that must be removed from G_D in order for the graph to be k -colorable. From here on, I will refer to the vertices of G_D and the observations they represent interchangeably.

B.2. Proof of Theorem 1

First, observe that K is the (unique) solution to the problem in (3) if and only if

$$K + \lambda \varepsilon_K(D) < k + \lambda \varepsilon_k(D) \quad \text{for every integer } k \neq K.$$

We can break this condition up into two parts. For $k > K$, the inequality

$$\varepsilon_K(D) - \varepsilon_k(D) < (k - K)/\lambda$$

requires that preferences beyond the best set of K each rationalize *no more than* $1/\lambda$ choice observations (beyond what would already be rationalized using the best K). For $k < K$, the inequality

$$\varepsilon_k(D) - \varepsilon_K(D) > (K - k)/\lambda$$

requires that each additional preference up to K uniquely rationalizes *at least* $1/\lambda$ choice observations (beyond what would be rationalized using the best k). Lemmas 1 and 2 bound the probability of each of these events.

Lemma 1. *There exists a constant $c_1 > 0$ (uniform across n) such that*

$$v^n \left(\left\{ D_n : \varepsilon_K(D_n) - \varepsilon_k(D_n) < \frac{k - K}{\lambda} \quad \forall k > K \right\} \right) \geq 1 - e^{-c_1 n} \quad \forall n$$

Proof. Suppose some $k > K$ is selected as the solution to (3), so that there are $k - K$ preferences (beyond the best set of K preferences) that together rationalize $(k - K)/\lambda$ additional choice observations. Clearly, a necessary condition for such a k to exist is that the best K preferences leave at least $1/\lambda$ observations unexplained; that is, $\varepsilon_K(D) \geq 1/\lambda$. Equivalently, a *sufficient* condition for

$$\varepsilon_K(D) - \varepsilon_k(D) < (k - K)/\lambda \quad \forall k > K \tag{16}$$

is that the best K preferences leave *fewer* than $1/\lambda$ observations unexplained; that is, $\varepsilon_K(D) < 1/\lambda$. The probability of this event lower bounds the probability of (16).

Recall that \mathcal{P} is the DM’s “true” set of K preferences, and define $T(D)$ to be the number of realized choice observations in the data set D that cannot be rationalized by any preference in \mathcal{P} . Removing these $T(D)$ “bad” observations from D results in a K -rationalizable data set; thus, $\varepsilon_K(D)$ cannot exceed $T(D)$. By assumption, the probability of error in any given observation is no more than p . Thus, the random variable $Y_n \sim \text{Bin}(n, p)$ first-order stochastically dominates $T(D_n)$.

These observations imply that

$$\begin{aligned} v^n(\{D_n : \varepsilon_K(D) - \varepsilon_k(D) < (k - K)/\lambda \quad \forall k > K\}) &\geq v^n(\{D_n : \varepsilon_K(D_n) < 1/\lambda\}) \\ &\geq v^n(\{D_n : T(D_n) < 1/\lambda\}) \\ &\geq \Pr(Y_n < 1/\lambda) \end{aligned}$$

Since $\lambda = 1/(\tilde{p}n)$ is chosen with $\tilde{p} > p$, it follows from Hoeffding’s Inequality that

$$\Pr(Y_n \geq 1/\lambda) = \Pr(Y_n - pn \geq 1/\lambda - pn) \leq \exp\left(-\frac{2(1/\lambda - pn)^2}{n}\right) = e^{-2(\tilde{p}-p)^2 n}$$

Thus

$$v^n \left(\left\{ D_n : \varepsilon_K(D_n) - \varepsilon_k(D_n) < \frac{k - K}{\lambda} \quad \forall k > K \right\} \right) \geq 1 - e^{-c_1 n}$$

with $c_1 := 2(\tilde{p} - p)^2 > 0$, and we are done. \square

Lemma 2. $v^n \left(\left\{ D_n : \varepsilon_k(D_n) - \varepsilon_K(D_n) > \frac{K - k}{\lambda} \quad \forall k < K \right\} \right) \rightarrow 1$ as $n \rightarrow \infty$.

Proof. The basic approach is to study the induced graph G_D , and lower bound the number of disjoint complete K -partite subgraphs.⁴⁰ With some imprecision, a sufficient condition for

$$\varepsilon_k(D_n) - \varepsilon_K(D_n) > \frac{K - k}{\lambda} \quad \forall k < K \tag{17}$$

is that the number of disjoint complete K -partite subgraphs in the induced graph G_D is at least $1/\lambda$. Roughly, the reason this is sufficient is as follows: Consider use of any $k < K$ colors to

⁴⁰ Recall that a *complete K -partite* subgraph is a graph whose nodes can be partitioned into K sets, where no nodes in the same set are connected by an edge.

color the graph G_D . Clearly, only k nodes in each complete K -partite graph can be colored. Moreover, each of the additional $K - k$ colors permits the coloring of at least one more node per complete K -partite subgraph, and hence $(K - k)/\lambda$ more nodes in the graph in total. Thus, going from k colors to K colors reduces the number of uncolored nodes by at least $(K - k)/\lambda$. Using Observation 4, this is equivalent to the statement that going from k preferences to K preferences leaves at least $(K - k)/\lambda$ fewer observations unexplained, and this implies the desired (17).

Instead of directly working with the distribution over graphs induced by ν , it is simpler to first study the related “perfect maximization” graph-generating process, where choice observations are sampled from ν^* . Let $G_{D_n}^* \sim (\nu^*)^n$ denote the typical graph corresponding to n samples from ν^* , and let $T(G_{D_n}^*)$ denote the number of disjoint complete K -partite subgraphs in $G_{D_n}^*$. Since the measure $(\nu^*)^n$ assigns probability 1 to data sets that are K -rationalizable, the graph $G_{D_n}^* \sim (\nu^*)^n$ must be K -colorable. This imposes a particular structure on the relationship between the complete K -partite graphs; in particular, each color up to K colors permits the coloring of at least $T(G_{D_n}^*)$ additional colors.

By definition of the differentiation parameter

$$\Pr(T(G_{D_n}^*) > n \cdot d(\pi, \mu^*)) \rightarrow 1 \quad \text{as } n \rightarrow \infty$$

so for large n the probability that $G_{D_n}^*$ includes at least $nd(\pi, \mu^*)$ disjoint complete K -partite subgraphs is close to 1. Finally, by assumption, the differentiation parameter satisfies $d(\pi, \mu^*) > \frac{1}{\lambda(1-p)^K n}$. So

$$\Pr\left(T(G_{D_n}^*) > \frac{1}{\lambda(1-p)^K}\right) \rightarrow 1 \quad \text{as } n \rightarrow \infty. \tag{18}$$

Thus to summarize, for large n each color up to K colors can (with high probability) color an additional $\frac{1}{\lambda(1-p)^K}$ nodes of $G_{D_n}^*$.

Suppose now that each node in $G_{D_n}^*$ is removed from all of its edges with probability p , and call the resulting graph $\tilde{G}_{D_n}^*$. Because each complete K -partite subgraph of $G_{D_n}^*$ is preserved with independent probability $(1 - p)^K$, the random variable $T(\tilde{G}_{D_n}^*)$ has distribution $\text{Bin}(T(G_{D_n}^*), (1 - p)^K)$, and its expectation is $\mathbb{E}\left(T(\tilde{G}_{D_n}^*)\right) = \mathbb{E}\left(T(G_{D_n}^*)\right) \cdot (1 - p)^K$. Thus, (18) implies

$$\Pr\left(T(\tilde{G}_{D_n}^*) > \frac{1}{\lambda}\right) \rightarrow 1 \quad \text{as } n \rightarrow \infty. \tag{19}$$

Finally, observe that the number of complete K -partite subgraphs in $G_{D_n} \sim (\nu)^n$ first-order stochastically dominates the number of complete K -partite subgraphs in $\tilde{G}_{D_n}^*$. Thus (19) further implies that

$$\Pr\left(T(G_{D_n}) > \frac{1}{\lambda}\right) \rightarrow 1 \quad \text{as } n \rightarrow \infty$$

and we are done. \square

The desired result follows directly from Lemmas 1 and 2.

B.3. Proof of Corollary 3

Fix any data set $\underline{D} = \{(x_i, \underline{A}_i)\}_{i=1}^n$. Let $X = \{x_i\}_{i=1}^n$ consist of all choice alternatives that are observed to be chosen, and define new choice sets $A_i = X \cap \underline{A}_i$. Set $D = \{(x_i, A_i)\}_{i=1}^n$.

Lemma 3. For every $\lambda \in \mathbb{R}_+$,

$$\min_{k \in \mathbb{Z}_+} [k + \lambda \varepsilon_k(D)] = \min_{k \in \mathbb{Z}_+} [k + \lambda \underline{\varepsilon}_k(\underline{D})]$$

Proof. I will show that $\varepsilon(D, \mathcal{P}) = \delta$ for some set of preference orderings $\mathcal{P} = \{P_j\}_{j=1}^k$ if and only if $\underline{\varepsilon}(\underline{D}, \mathcal{U}) = \delta$ for some set of utility functions $\mathcal{U} = \{u_j\}_{j=1}^k$.

Fix any set $\mathcal{P} = \{P_j\}_{j=1}^k$ of k orderings defined on X , and take $\delta := \varepsilon(D, \mathcal{P})$. Every ordering P_j admits representation via a utility function $u_j : X \rightarrow \mathbb{R}$.⁴¹ Moreover, we can extend each u_j to a continuous function \underline{u}_j on \underline{X} such that the maximum of \underline{u}_j on each \underline{A}_i is obtained at $\text{argmax}_{x \in A_i} u_j(x)$. This implies that choice alternative x_j is \underline{u}_j -maximal in \underline{A}_i if and only if it is P_j -maximal in A_i . Now set $\mathcal{U} = \{\underline{u}_j\}$. Then,

$$\underline{\varepsilon}(\underline{D}, \mathcal{U}) = \# \left\{ (x, \underline{A}) \in \underline{D} : x \neq \text{argmax}_{x' \in \underline{A}} \underline{u}_j(x') \text{ for all } j = 1, \dots, k \right\} = \delta.$$

Thus, choice error δ is attainable using a set of k utility functions. It follows that

$$\min_{k \in \mathbb{Z}_+} [k + \lambda \varepsilon_k(D)] \geq \min_{k \in \mathbb{Z}_+} [k + \lambda \underline{\varepsilon}_k(\underline{D})]. \tag{20}$$

In the other direction, fix a set $\mathcal{U} = \{u_j\}_{j=1}^k$ of k continuous functions defined on \underline{X} , and take $\delta = \underline{\varepsilon}(\underline{D}, \mathcal{U})$. For every utility function u_j , let P_j be the ordering on X that satisfies

$$x P_j x' \iff u_j(x) > u_j(x').$$

Then, $x_j = \text{argmax}_{x \in \underline{A}} u_j(x)$ if and only if x_j is P_j -maximal in A . Setting $\mathcal{P} = \{P_j\}_{j=1}^k$, we have that

$$\varepsilon(D, \mathcal{P}) = \#\{(x, A) \in D : x \text{ is not } P_j\text{-maximal in } A \text{ for any } j = 1, \dots, k\} = \delta.$$

Thus, also

$$\min_{k \in \mathbb{Z}_+} [k + \lambda \varepsilon_k(D)] \leq \min_{k \in \mathbb{Z}_+} [k + \lambda \underline{\varepsilon}_k(\underline{D})]$$

as desired. \square

It follows from this lemma that $\underline{K}_\lambda^*(\underline{D}) = K_\lambda^*(D)$ for every choice of λ , so the problem posed in Section 6 can be mapped directly into a corresponding discrete problem of the form described in Section 3. Apply Theorem 1, and the desired result follows.

⁴¹ That is, there exists a utility function u_j such that for every $x, x' \in X$, $x P_j x'$ if and only if $u_j(x) > u_j(x')$.

B.4. Proof of Claim 2

(a) For every number of observations n , the single preference that minimizes the expected value of $\varepsilon_1(D_n)$ is $x_3 P x_2 P x_1$. It can be shown that the pair of preferences that minimizes the expected value of $\varepsilon_2(D_n)$ is $\{P, P'\}$ with $x_3 P x_2 P x_1$ and $x_1 P' x_2 P' x_3$. The associated expected errors are

$$\begin{aligned} \mathbb{E}(\varepsilon_1(D_n)) &= \frac{1}{4} \left(\frac{e^\gamma + e^{2\gamma}}{e^\gamma + e^{2\gamma} + e^{3\gamma}} \right) + \frac{1}{4} \left(\frac{e^\gamma}{e^\gamma + e^{2\gamma}} \right) \\ &\quad + \frac{1}{4} \left(\frac{e^\gamma}{e^\gamma + e^{3\gamma}} \right) + \frac{1}{4} \left(\frac{e^{2\gamma}}{e^{2\gamma} + e^{3\gamma}} \right) \\ \mathbb{E}(\varepsilon_2(D_n)) &= \frac{1}{4} \left(\frac{e^{2\gamma}}{e^\gamma + e^{2\gamma} + e^{3\gamma}} \right) \end{aligned}$$

while $\varepsilon_3(D_n) = 0$ for every data set D_n . From this, we see that the (expected) error-preference tradeoff curve is convex, and thus it is sufficient to show that $\mathbb{E}(\varepsilon_1(D_n)) - \mathbb{E}(\varepsilon_2(D_n))$ is decreasing in γ for all $\gamma > 0$.

Using the above displays,

$$\begin{aligned} \mathbb{E}(\varepsilon_1(D_n)) - \mathbb{E}(\varepsilon_2(D_n)) &= \frac{1}{4} \left(\frac{e^\gamma}{e^\gamma + e^{2\gamma} + e^{3\gamma}} \right) + \frac{1}{4} \left(\frac{e^\gamma}{e^\gamma + e^{2\gamma}} \right) \\ &\quad + \frac{1}{4} \left(\frac{e^\gamma}{e^\gamma + e^{3\gamma}} \right) + \frac{1}{4} \left(\frac{e^{2\gamma}}{e^{2\gamma} + e^{3\gamma}} \right) \end{aligned}$$

Each component of this sum is decreasing in γ for all $\gamma > 0$, so the sum is as well. The desired result follows.

(b) From above, we already have that $\mathbb{E}(\varepsilon_1(D_n)) - \mathbb{E}(\varepsilon_2(D_n))$ is decreasing in γ . It will be sufficient to show additionally that $\mathbb{E}(\varepsilon_2(D_n)) - \mathbb{E}(\varepsilon_3(D_n))$ is decreasing in γ . Using the calculations from part (a), for every data set D_n , we have

$$\mathbb{E}(\varepsilon_2(D_n)) - \mathbb{E}(\varepsilon_3(D_n)) = \frac{1}{4} \left(\frac{e^{2\gamma}}{e^\gamma + e^{2\gamma} + e^{3\gamma}} \right)$$

which is indeed decreasing in γ for all $\gamma > 0$.

B.5. Proof of Proposition 1

First, I will show the following:

Claim 3. *If there exist orderings $P_1, P_2 \in \mathcal{P}$ such that a choice alternative x is ranked first according to P_1 , and not last according to P_2 , then \mathcal{P} is not identifiable.*

Proof. Consider any set of orderings \mathcal{P} that includes the preferences P_1, P_2 , where x_1 is ranked first according to P_1 , x_2 is ranked last according to P_2 , and the alternatives x_1 and x_2 are not the same. Fix any data set D that can be perfectly rationalized using \mathcal{P} (that is, $\varepsilon(D, \mathcal{P}) = 0$). I will show by construction that there exists another set of preferences $\mathcal{P}' \neq \mathcal{P}$ with $|\mathcal{P}'| \leq |\mathcal{P}|$ such that also $\varepsilon(D, \mathcal{P}') = 0$.

Define the ordering P'_2 to agree with P_2 everywhere, except that it ranks x_1 last. Let $\mathcal{P}' = \mathcal{P} - \{P_2\} + \{P'_2\}$, where the operators denote set addition and subtraction. I will now show that $\varepsilon(D, \mathcal{P}') = 0$.

Suppose towards contradiction that there is some choice observation $(x, A) \in D$ where x is not P -maximal in A for any $P \in \mathcal{P}'$. By assumption, there is some ordering $P^* \in \mathcal{P}$ such that x is P^* -maximal in A . If $P^* \neq P_2$, then also $P^* \in \mathcal{P}'$, which yields an immediate contradiction. So it must be that x is P_2 -maximal in A . Now, there are two possibilities: if $x \neq x_1$, then x must also be P'_2 -maximal in A , and we are done. If $x = x_1$, then x is P_1 -maximal in A (by definition of x_1). So we are again done. \square

Every set \mathcal{P} with $|\mathcal{P}| \geq 3$ satisfies the condition in Claim 3, and hence fails to be identifiable.

Now let us consider an arbitrary set $\mathcal{P} = \{P_1, P_2\}$. Index the alternatives such that x_1 is ranked first according to P_1 and x_2 is ranked first according to P_2 . If the condition in Claim 3 is satisfied, it again follows that \mathcal{P} is not identifiable. Suppose otherwise, so that x_1 is ranked last according to P_2 and x_2 is ranked last according to P_1 . Define

$$D = \{(x, A) : x \text{ is } P_1\text{-maximal in } A, A \in 2^X\} \cup \{(x, A) : x \text{ is } P_2\text{-maximal in } A, A \in 2^X\}.$$

Clearly there is no singleton set \mathcal{P}' such that $\varepsilon(D, \mathcal{P}') = 0$. Suppose towards contradiction that there exists some set $\mathcal{P}' = \{P'_1, P'_2\} \neq \mathcal{P}$ such that $\varepsilon(D, \mathcal{P}') = 0$. It must be that x_1 is ranked first according to P'_1 and last according to P'_2 , and that x_2 is ranked first according to P'_2 and last according to P'_1 (otherwise relabel P'_1 and P'_2).

Without loss of generality, suppose that $P'_1 \neq P_1$.⁴² Then there exist alternatives x_i, x_j such that x_i is preferred to x_j under P_1 but not under P'_1 :

$$x_i P_1 x_j \quad \text{and} \quad x_j P'_1 x_i.$$

Take $A := \{x : x_i P_1 x\} \cup \{x_i\}$ to be the set of all alternatives that P_1 ranks weakly lower than x_i . Then $(x_i, A) \in D$. But x_i cannot be P'_1 -maximal in A since $x_j \in A$, and x_1 cannot be P'_2 -maximal in A since $x_2 \in A$. So $\varepsilon(D, \mathcal{P}') > 0$ as desired.

Finally, every singleton set $\mathcal{P} = \{P\}$ is trivially identifiable, taking

$$D = \{(x, A) : x \text{ is } P\text{-maximal in } A, A \in 2^X\}.$$

B.6. Proof of Proposition 2

Consider any choice observation $(x, A) \in \mathbb{C}(\mathcal{P})$. By definition of α , $v^*(x, A) > \alpha$, so also $v(x, A) > \alpha(1 - p)$. Thus, the number of occurrences of (x, A) in the observed data first-order stochastically dominates the random variable $Z_n \sim \text{Bin}(n, \alpha(1 - p))$. Since by assumption, $1/\lambda = \tilde{p}n < \alpha(1 - p)n$, it follows from Hoeffding’s inequality that

$$\begin{aligned} \Pr\left(Z_n < \frac{1}{\lambda}\right) &= \Pr(Z_n - \alpha(1 - p)n < 1/\lambda - \alpha(1 - p)n) \\ &\leq \exp\left(-2(\tilde{p} - \alpha(1 - p))^2 n\right) \end{aligned}$$

⁴² Otherwise, $P'_2 \neq P_2$ and the remainder of the proof is mirrored.

Setting $c_1 = 2(\tilde{p} - \alpha(1 - p))^2 > 0$, we have

$$\Pr\left(Z_n < \frac{1}{\lambda}\right) \leq e^{-c_1 n}.$$

So the probability that (x, A) appears fewer than $1/\lambda$ times in the realized data set is also no more than $e^{-c_1 n}$. Taking a union bound, the probability that any $(x, A) \in \mathbb{C}(\mathcal{P})$ appears fewer than $1/\lambda$ times in the data is no more than $|\mathbb{C}(\mathcal{P})|e^{-c_1 n}$. Thus, the probability that every $(x, A) \in \mathbb{C}(\mathcal{P})$ appears at least $1/\lambda$ times in the realized data set is at least

$$1 - |\mathbb{C}(\mathcal{P})|e^{-c_1 n}$$

which converges to 1 as the quantity of data n increases. This immediately implies that

$$\Pr(\mathbb{C}(\mathcal{P}) \subseteq \mathbb{C}(\mathcal{P}_\lambda^*(D_n))) \rightarrow 1 \text{ as } n \rightarrow \infty. \tag{21}$$

In the other direction, the random variable $Y_n \sim \text{Bin}(n, p)$ first-order stochastically dominates the number of observations of all $(x, A) \notin \mathbb{C}(\mathcal{P})$. (Informally, Y_n is the number of choice observations in error.) Thus, $\Pr(Y_n \leq 1/\lambda)$ is an upper bound on the probability that there are $1/\lambda$ realized observations outside of the set $\mathbb{C}(\mathcal{P})$. Since by assumption $1/\lambda = \tilde{p}n > pn$,

$$\Pr(Y_n \geq 1/\lambda) \rightarrow 0. \tag{22}$$

Combining this with (21), it is optimal to recover preferences whose choice implications are exactly $\mathbb{C}(\mathcal{P})$ whenever feasible. By construction, $\mathbb{C}(\mathcal{P})$ is the set of choice implications corresponding to the set of preferences \mathcal{P} , and \mathcal{P} is a valid solution to the problem in (14). Thus,

$$\Pr(\mathbb{C}(\mathcal{P}_\lambda^*(D_n)) = \mathbb{C}(\mathcal{P})) \rightarrow 1 \text{ as } n \rightarrow \infty$$

as desired.

B.7. Proof of Proposition 3

Below, let

$$\mathcal{D}_g := \{(x, A), C) : C \in \mathcal{C}, (x, A) \in \mathbb{C}(\{m(C)\})\}$$

be the set of all “good” tuples $((x, A), C)$ consistent with (the true) mapping m , and

$$\mathcal{D}_b := \{(x, A), C) : C \in \mathcal{C}, (x, A) \notin \mathbb{C}(\{m(C)\})\}$$

be the set of all “bad” tuples $((x, A), C)$ that are not consistent with m .

The following lemma demonstrates a sufficient condition on λ and the observed data D given which the recovered mapping satisfies $m_\lambda^*(D) = m$.

Lemma 4. *Suppose that*

- (a) every tuple $((x, A), C) \in \mathcal{D}_g$ appears at least $2/\lambda$ times in the data D , and
- (b) there are fewer than $1/\lambda$ (total) instances of tuples from \mathcal{D}_b in the data D .

Then $m_\lambda^*(D) = m$.

Proof. Suppose to the contrary that (a) and (b) are satisfied, but the recovered mapping is some $m' \neq m$. First consider the case

$$|m'(\mathcal{C})| \geq |m(\mathcal{C})| \tag{23}$$

in which there must be some preference P assigned by mapping m but not by mapping m' . Consider any $(x, A) \in \mathbb{C}(\{P\}) \setminus \mathbb{C}(\{P'\})$. By (a), the observation $((x, A), C)$ appears at least $2/\lambda$ times in the choice data. These observations can be rationalized using m but not by m' . Moreover, by (b), the number of observations that can be rationalized using m' but not by m is no more than $1/\lambda$. So $\varepsilon(D, m) < \varepsilon(D, m')$. Combining this with (23), we have that $|m'(\mathcal{C})| + \lambda\varepsilon(D, m') > |m(\mathcal{C})| + \lambda\varepsilon(D, m)$, and hence m' cannot have been the recovered mapping. This yields the desired contradiction.

If instead

$$|m'(\mathcal{C})| < |m(\mathcal{C})|$$

then there must be at least $|m(\mathcal{C})| - |m'(\mathcal{C})|$ contexts where the preference assigned by m is different from the preference assigned by m' . Call the set of such contexts \mathcal{C}^* . For each $C \in \mathcal{C}^*$, there is some (x, A) satisfying $(x, A) \in \mathbb{C}(m(C))$ but $(x, A) \notin \mathbb{C}(m'(C))$. Thus, by (a), there are at least $2/\lambda \cdot (|m(\mathcal{C})| - |m'(\mathcal{C})|)$ observations that can be rationalized under m but not under m' . So $\varepsilon(D, m) - \varepsilon(D, m') < -2/\lambda \cdot |m(\mathcal{C})| - |m'(\mathcal{C})|$. By (b), the number of observations that can be rationalized using m' but not by m is no more than $1/\lambda$. Thus,

$$\begin{aligned} \varepsilon(D, m') - \varepsilon(D, m) &> \frac{2}{\lambda} (|m(\mathcal{C})| - |m'(\mathcal{C})|) - \frac{1}{\lambda} \\ &\geq \frac{1}{\lambda} (|m(\mathcal{C})| - |m'(\mathcal{C})|) \end{aligned}$$

using that $|m(\mathcal{C})| - |m'(\mathcal{C})| \geq 1$. This implies that $|m'(\mathcal{C})| + \lambda\varepsilon(D, m') > |m(\mathcal{C})| + \lambda\varepsilon(D, m)$, so m' is again not the recovered mapping. \square

Now I will show that conditions (a) and (b) are satisfied with probability converging to 1 as the number of observations gets large. By definition of α , every tuple $((x, A), C) \in \mathcal{D}_g$ is observed with probability at least $\alpha(1 - p)$. For each $((x, A), C) \in \mathcal{D}_g$, let $E_{((x,A),C)}$ be the event that $((x, A), C)$ is observed fewer than $2/\lambda = 2\tilde{p}n$ times. By assumption, $2\tilde{p}n < \alpha(1 - p)n$, so application of Hoeffding’s inequality gives

$$\Pr(E_{((x,A),C)}) \leq e^{-2(\alpha(1-p)-2\tilde{p})^2n}.$$

Using a union bound, the probability that any $E_{((x,A),C)}$ (with $((x, A), C) \in \mathcal{D}_g$) occurs is no more than $\kappa \cdot e^{-2(\alpha(1-p)-2\tilde{p})^2n}$, where $\kappa = |\mathcal{C}| \cdot |\mathbb{C}(m(C))|$ is a constant. Or equivalently, the probability that every $((x, A), C) \in \mathcal{D}_g$ is observed at least $2/\lambda$ times is

$$1 - \kappa \cdot e^{-2(\alpha(1-p)-2\tilde{p})^2n}. \tag{24}$$

Thus, the probability that the condition in (a) is satisfied converges to 1 as $n \rightarrow \infty$.

To show that the condition in (b) also converges to 1, let Z_n be the number of tuples $((x, A), C) \in \mathcal{D}_b$ in the observed data D . The random variable Z_n is first-order stochastically dominated by the random variable $Z \sim \text{Bin}(n, p)$. Since $\mathbb{E}(Z) = pn < \tilde{p}n = 1/\lambda$, it follows that

$$\Pr(Z \geq 1/\lambda) \rightarrow 0 \text{ as } n \rightarrow \infty \tag{25}$$

implying the same for Z_n . Combining (24) and (25), the desired result directly follows from Lemma 4.

References

- Afriat, S.N., 1967. The construction of a utility function from expenditure data. *Int. Econ. Rev.* 8, 67–77.
- Ambros, A., Rozen, K., 2013. Rationalizing choice with multi-self models. *Econ. J.*
- Apestequia, J., Ballester, M.A., 2012. A Measure of Rationality and Welfare. Working Paper.
- Bernheim, B.D., Rangel, A., 2009. Beyond revealed preference: choice theoretic foundations for behavioral welfare economics. *Q. J. Econ.*
- Crawford, I., Pendakur, K., 2012. How many types are there? *Econ. J.*
- Dean, M., Martin, D., 2010. How Consistent are your Choice Data? Working Paper.
- Dean, M., Martin, D., 2016. Measuring rationality with the minimum cost of revealed preference violations. *Rev. Econ. Stat.*
- Deb, R., 2009. A testable model of consumption with externalities. *J. Econ. Theory* 144.
- Echenique, F., Lee, S., Shum, M., 2011. The money pump as a measure of revealed preference violations. *J. Polit. Econ.* 119.
- Einav, L., Finkelstein, A., Pascu, I., Cullen, M., 2012. How general are risk preferences? Choices under uncertainty in different domains. *Am. Econ. Rev.*
- Famulari, M., 1995. A household-based, nonparametric test of demand theory. *Rev. Econ. Stat.* 77, 372–383.
- Fudenberg, D., Levine, D., 2006. A dual-self model of impulse control. *Am. Econ. Rev.* 96, 1449–1476.
- Green, J., Hojman, D., 2007. Choice, rationality and welfare measurement.
- Gross, A., 1989. Determining the Number of Violators of the Weak Axiom. Tech. rep. University of Wisconsin–Milwaukee.
- Halevy, Y., Persitz, D., Zrill, L., 2015. Parametric Recoverability of Preferences. Working Paper.
- Houtman, M., Maks, J., 1985. Determining all maximal data subsets consistent with revealed preference. *Kwant. Methoden* 19, 89–104.
- Kalai, G., Rubinstein, A., Spiegler, R., 2002. Rationalizing choice functions by multiple rationales. *Econometrica* 70, 2481–2488.
- Luce, D., Raiffa, H., 1957. *Games and Decisions: Introduction and Critical Survey*. Wiley, New York.
- Manzini, P., Mariotti, M., 2007. Sequentially rationalizable choice. *Am. Econ. Rev.*
- Manzini, P., Mariotti, M., 2009. Categorize then choose: boundedly rational choice and welfare.
- McFadden, D., Richter, M., 1970. Revealed stochastic preference.
- Quandt, R.E., 1956. A probabilistic theory of consumer behavior. *Q. J. Econ.* 70, 507–536.
- Rubinstein, A., Salant, Y., 2006. A model of choice from lists. *Theor. Econ.* 3.
- Rubinstein, A., Salant, Y., 2008. (A, f) : choice with frames. *Rev. Econ. Stud.*
- Sen, A., 1993. Internal consistency of choice. *Econometrica* 61.
- Train, K., 1986. *Qualitative Choice Analysis*. Cambridge University Press.
- Varian, H.R., 1982. The nonparametric approach to demand analysis. *Econometrica* 50, 945–973.