

Tipo de documento: Artículo

Demand Estimation Under Uncertain Consideration Sets

Autoría ditelliana: Vulcano, Gustavo

Otros autores: Jagabathula, Srikanth, Mitrofanov, Dmitry

Fecha de publicación: 2023

Versión final Publicada en: Operations Research (ISSN: 0030-364X)

El presente documento es una versión aceptada del artículo publicado

¿Cómo citar este artículo?

Jagabathula, S., Mitrofanov, D., & Vulcano, G. (2023). Demand Estimation Under Uncertain Consideration Sets. Operations Research. <https://doi.org/10.1287/opre.2022.0006>

El presente documento se encuentra alojado en el Repositorio Digital de la Universidad Torcuato Di Tella bajo una licencia Creative Commons Atribución-No Comercial-Compartir Igual 2.5 Argentina (CC BY-NC-SA 2.5 AR)

Dirección: <https://repositorio.utdt.edu/handle/20.500.13098/12107>

Demand Estimation under Uncertain Consideration Sets

Srikanth Jagabathula

Leonard N. Stern School of Business, New York University, New York, NY 10012, sjagabat@stern.nyu.edu

Dmitry Mitrofanov

Carroll School of Management, Boston College, Chestnut Hill, MA, dmitry.mitrofanov@bc.edu

Gustavo Vulcano

School of Business, Universidad Torcuato Di Tella,
and CONICET, Buenos Aires, Argentina, gvulcano@utdt.edu

Key words

History

Over the last two decades there has been growing interest in the operations management (OM) academic field and in the industry practice to incorporate sophisticated demand models, which provide high-quality inputs for critical tasks such as inventory management, dynamic pricing and assortment planning. Examples of such models include the multinomial logit (MNL), the nested logit, the mixed logit and the latent class MNL, which are traditional in the marketing and economics fields but novel in terms of their applicability in operational contexts. These models have been widely studied resulting in the development of specific estimation (e.g., Newman et al. (2014), Vulcano et al. (2012), Jagabathula et al. (2020)) or assortment optimization algorithms (e.g., Talluri and Van Ryzin (2004), Feldman and Topaloglu (2015), Davis et al. (2014)). More recently, new demand models have been proposed (e.g., the Markov chain model by Blanchet et al. (2016)), and others like the rank list-based model (e.g., Farias et al. (2013)) and the exponential model (e.g., Alptekinoglu and Semple (2016)) have been revisited, jointly with the expansion of the application of choice-based demand models to the operation of online platforms (e.g., Lee and Lee (2012)). Companies in different industry sectors, from airlines to retailers and more recently, online sharing economy platforms, have been testing and incorporating some of these models into their operational capabilities.

The usual data source to calibrate these demand models are records of past transactions, either sales transaction data in the case of retail operations and revenue management, or bookings from past interactions between peers in the case of online platforms. For each transaction, the understanding is that the client (or an agent in a platform) selected one option from a collection of alternatives or *offer set*, which is usually defined as a full category assortment in the retail case, or the full set of available options within an arbitrary radius in spatial choice models (e.g., car sharing). Given the transaction data, most choice models are trained assuming that the chosen option is preferred over all the other products on offer. However, customers may not consider everything on offer. For instance, in retailing, a customer selecting from the coffee category may not evaluate the full assortment and might consider only a subset of products (e.g., decaf) in a choice instance. In online platforms, even if we arbitrarily define the offer set as the cars available within a 0.2-mile radius, an agent may evaluate only compact cars. If we ignore these consideration sets, we make the incorrect inference that the chosen product is preferred over products not even considered, leading to model bias. To deal with this issue, the so-called *consider-then-choose* (CTC) models have been proposed in the literature. These models posit that customers sample a consideration set and then choose the most preferred product from the intersection of the offer set and the consideration set (e.g., Howard and Sheth (1969), Alba and Chattopadhyay (1985), Hauser and Wernerfelt (1990)).

CTC models have been studied quite extensively in the marketing literature (e.g., see [Roberts and Lattin \(1997\)](#)). They also gained recent popularity within the OM literature to make assortment and pricing decisions (e.g., see [Feldman et al. \(2019\)](#); [Aouad et al. \(2020\)](#)).

Despite their richness, CTC models are often difficult to fit in practice because customers' consideration sets are not observed—we know the customer choice and the offer set, but the consideration set itself could be any subset of the full category containing the chosen product. When consideration sets are not observed, CTC models may not be identifiable and it is unclear what, if any, predictive advantage they offer over simply assuming that customers consider everything on offer as classic models do. Existing literature is mostly silent on these issues. It offers broad empirical evidence that customers do indeed form consideration sets and studies specific instances of CTC models. But it mostly takes it as given that firms should fit CTC models over classic choice models when consideration sets are not observed.

In this paper, we systematically study a very general class of CTC models when we only observe customer choices and offer sets. Theoretically, we observe that the general CTC model class is equivalent to the random utility maximization (RUM) class of models (which assumes that customers consider everything on offer), indicating that CTC models span the same modeling scope as classic choice models. We also show that CTC models in general are not identifiable from transaction data alone, but some restricted versions of them indeed are.

To empirically evaluate the CTC models, we develop techniques to estimate them from transaction data. Our numerical analysis on both synthetic and real-world data shows that CTC models outperform classic choice models when offer sets are not observed perfectly (i.e., they are noisy) and the noise is asymmetric between the training and test data. Here, the offer set *noise* associated with a product refers to a product being erroneously recorded as offered or stocked out when in fact it was not. We say noise is *asymmetric* if the degree of the offer set noise associated with a product differs between training and test datasets (e.g., because factors affecting stockouts vary over time). On the other hand, when the noise is symmetric, CTC models have comparable performance to classic models, providing little to no predictive advantage over them.

Noise in offer sets is quite common in practice. In retailing, for instance, offer set descriptions are often unreliable because of potential inventory inaccuracies (e.g., [DeHoratius and Raman \(2008\)](#)). For example, [Kang and Gershwin \(2005\)](#) analyze the accuracy of inventory records of a global retailer and find that only 51% of them match actual inventory on average, with the worst store experiencing more than 67% mismatch. In online sharing platforms, *offer* or *availability* sets are not even clearly defined. For example, in the context of a peer-to-peer car-sharing platform, an offer set can be defined as all the available car listings on the platform or more reasonably, as all

the car listings within a radius of the customer's location. Either definition is arbitrary and prone to errors.

Furthermore, the so-called "test order set" inputs used for predictions tend to be noisier (i.e., more different from the actual ones) than the training order sets. Choice models must be given order sets to make predictions but the firm faces a high degree of uncertainty about future order sets. For example, a retailer is uncertain about next week's order set because existing products may be depleted, stocked-out products may be replenished, or new products might be introduced. Future order sets are even more uncertain for online sharing platforms because availability is determined in real-time not by the platform but by individual providers in the market whose decisions are difficult to predict. These observations make the sequence (i) consider, and (ii) choose, particularly important when modeling customers' choices in practice.

Unlike existing literature which studies specific instances of CTC models, we study the following general class. The population preferences are characterized by a joint distribution over rank lists of the full set of products and over subsets of the product universe (i.e., over the consideration sets). This joint distribution is common across all customers. In each choice instance, a customer samples a consideration set and a preference list, and purchases the most preferred product in the *choice set*, which results from the intersection between the sampled consideration set and the exhibited order set, or does not purchase at all if neither of the considered products is offered.

We make the following remarks about our model. First, we note that we use a product-based as opposed to a feature-based consideration set definition. In a product-based consideration set definition, the model directly specifies the probability of consideration for each product. By contrast, a feature-based definition assumes that customers screen on features, considering only those products whose attributes are within pre-specified acceptable ranges (e.g., see Jagabathula and Rusmevichientong (2017)). Product-based consideration sets are more general and tractable to analyze. They can also readily subsume feature-dependence, as we illustrate in this paper. Second, a key distinguishing aspect of our model is that we allow the distribution over consideration sets to be general, in contrast to the bulk of the existing literature which has generally restricted it to belong to specific classes. Third, while the joint distribution can be general in principle, for the purpose of estimation, we approximate it with an appropriately defined mixture distribution.

The main goal of our proposal is to characterize business environments under which the use of CTC models may provide higher quality estimates to describe the choice behavior of the customer basis, compared to RUM-based estimates. Specifically, we make the following contributions:

- *Statistical properties of CTC models.* Recall that the RUM class is equivalently described by a distribution over product preference lists (Block and Marschak (1960); Strauss (1979); Farias et al. (2013)), so that customers choose the most preferred product among the available ones according to the sampled ranking list.

We first show that the CTC model class is indeed equivalent to the RUM class, in that the set of choice probabilities induced by the CTC model class is the same as those induced by the RUM class (see Proposition 1). Yet, despite the equivalent explanatory power, CTC may be more natural in terms of its practical usage when accommodating consideration set formation becomes relevant (e.g., in a latent way, when customers do not consider all the items on offer; or in an explicit way, when having access to surveys where customers reveal their consideration sets).

Given the equivalency to the RUM class, CTC models are not fully identifiable. But what is somewhat surprising is that the marginal distribution over consideration sets is uniquely identified (see Proposition 3). We also investigate cases when consideration sets are small (see Corollary 1 and Proposition 4), and show that we can compute the marginal distribution over consideration sets efficiently. Furthermore, if we restrict all the customers to choose using the same preference list, then the ranking is also identifiable (see Proposition 5). We also establish conditions to verify whether the observed data are consistent with a specific instance of the CTC model class.

- *Methodology to estimate the parameters of CTC models.* For estimation, we approximate the general CTC model class with a mixture of what we call independent consideration set (ICS) models (see Section 4.1). An ICS model is specified by a single ranking and an independent distribution over consideration sets in which a customer samples a consideration set by including each product independently with a certain probability. We propose an expectation-maximization (EM) algorithm to estimate the mixture. For each mixture component, we propose an outer-approximation algorithm to ensure convergence to the maximum likelihood estimate.
- *Numerical experiments on synthetic data:* Because the consideration set formation is explicitly modeled in consider-then-choose type of frameworks, it is likely that its predictive performance is robust to the noise in the definition of the offer sets in comparison with competitive benchmarks (e.g., MNL and ranking-based models). We verify this conjecture by explicitly adding noise, erroneously including or excluding products from the offer sets. We find that the CTC models outperform classic models when the noise is asymmetric between the training and test offer sets (e.g., the sets of items that are exposed to noise in both training and test sets minimally intersect). Their performance on the other hand is comparable when noise is

symmetric. This result shows that when using choice models in practice, it is not sufficient to choose the model that provides the best predictive performance on the historical data set, but it is important to also understand how the prediction task might differ from the training task in the data generation process (e.g., product availability).

- *Empirical analysis: better demand predictions for the retail industry and online platforms.* To support our findings on synthetic data, we compare choice models under several real-world scenarios in retailing and online platforms where we are likely to face significant noise in the offer set definitions. On real-world grocery transaction data, we find that the relative performance of CTC models over the benchmarks improves as the level of noise in the test offer sets increases, for instance, when making long-term predictions. On a data set obtained from an online car-sharing platform, we show that CTC models outperform classic models because of the significant uncertainty in future availability sets. We also show that CTC models offer the flexibility to use machine learning models, such as decision trees and random forests, in modeling the consideration set distribution. These models are more interpretable and lead to better prediction accuracy than standard choice-based demand models (up to 53.7% improvement in the root mean squared error metric).

We emphasize that CTC models have traditionally been used to account for uncertainty in what the customers actually considered. In practice, decision-makers also face uncertainty about what the customers were actually offered. To the best of our knowledge, we are the first to highlight the uncertainty in offer sets and show that CTC models can successfully deal with both types of uncertainty.

The remainder of this paper is organized as follows. Section 2 positions our work within the existing literature. Section 3 defines our model and some identification results. We describe our data model and estimation algorithm in Section 4. The evaluation of the model starts in Section 5 with synthetic data experiments, followed by experiments on real data in Sections 6 (retailing) and 7 (online car sharing platform). Finally, our concluding remarks are discussed in Section 8.

Consider-then-choose (CTC) models are built upon the key concept of *consideration sets*. This notion has recently gained attention in the OM literature, but is well studied in the marketing and psychology fields, dating back to the papers by Campbell (1969), Howard and Sheth (1969) and Wright and Barbour (1977).

It has long been recognized that consumers usually make choices in a two-stage process (Swait and Ben-Akiva (1987); Lynch et al. (1991); Roberts and Lattin (1997)). First, they identify a

small subset of products for further evaluation, the so-called *consideration set*, and then purchase the most preferred product from this subset. It is hard to observe whether a product that is not purchased has been included or not in a consumer's consideration set, as it might even depend on a number of factors not necessarily related to the consumer's preferences. Nevertheless, there is ample empirical evidence in the literature about the consider-then-choose behavior of customers. In his seminal paper, [Hauser \(1978\)](#) shows that a model based on the consideration set concept accounts for as much as 78% of the explainable uncertainty in purchase transaction data. [Hauser and Wernerfelt \(1990\)](#) empirically observe that customers consider on average only 3 brands of deodorants, 4 brands of shampoos, 4 brands of laundry detergents, and 4 brands of coffee. In a follow-up paper, [Hauser \(2014\)](#) reports that the average size of the consideration set of consumer packaged goods in US is a tenth of the total number of brands in the product category. The aforementioned papers provide a foundation for the belief that the distribution over consideration sets in the CTC model, which encompasses all possible subsets, may be sparse in reality and hence, tractable to estimate.

The notion of consideration sets might arise from the limited information-gathering ability of consumers because they incur a search cost to learn detailed information about the products ([Ratchford \(1982\)](#)). The underlying justification is that consumers keep searching for products until the marginal expected gain from the search is less than the marginal search cost. Another argument to build consideration sets is related to cognitive heuristics, which are popular in the marketing and psychology literature while being of great importance for managerial decisions in advertising, product development, and strategic planning, e.g., conjunctive, disjunctive, compensatory, and elimination by aspects heuristics ([Tversky \(1972\)](#); [Montgomery and Svenson \(1976\)](#); [Hauser \(2014\)](#); [Hogarth and Karelaia \(2005\)](#)). Therefore, the theoretical groundwork from the fields of marketing and psychology strongly suggests that the set of products or services that people consider is highly unlikely to be the same as the set of offered products.

In the OM-related literature, the prevailing assumption aligned with the classical discrete choice literature has been that the consideration set is equivalent to the offer set. It has only been recently that more sophisticated consider-then-choose models of demand have been incorporated. [Aouad et al. \(2020\)](#) study the problem of assortment optimization under several variants of a choice model defined by two elements: a collection of consideration sets and a collection of customer types represented by rank lists. Different constraints in the definition of these collections lead to different special versions of the model (e.g., limiting the number of features that consumers use to filter a subset of alternatives). The authors develop a dynamic programming framework to study the computational aspects of assortment optimization under variants of these consider-then-choose

premises. They show that for many empirically vetted assumptions on how customers consider and choose, their resulting dynamic program is efficient.

Feldman and Topaloglu (2018) consider the assortment optimization problem under the MNL model when consideration sets for different customer types are nested, whereas Feldman et al. (2019) focus on the assortment optimization problem when customers choose in accordance with the rank list model of demand but under small consideration sets. Wang and Sahin (2018) present a consider-then-choose model where the consideration set is formed by balancing the incremental expected utility of a product and the related search cost. The subsequent choice behavior within a consideration set is governed by the MNL model. Given the hardness of the assortment optimization problem, they propose as an approximation a solution that may exclude some high-attractiveness products from the offer set. Jagabathula and Rusmevichientong (2017) propose a model where, first, customers consider the set of products with prices less than a threshold, and then choose the most preferred product from the set considered. They develop a tractable nonparametric expectation maximization (EM) algorithm to fit the model to transaction data and design an efficient algorithm to determine the profit-maximizing combination of offer sets and prices. Jagabathula and Vulcano (2018) propose a framework to estimate individual consumer preferences under some heuristic rules used by consumers to form their consideration sets (e.g., they consider only products under promotion jointly with the ones purchased in the previous store visit). In spite of the recent attention caught by CTC models in the OM field, most of the existing papers mainly focus on various optimization problems under the CTC choice rule, ignoring the question of when this type of model can stand out and outperform classical benchmarks. In contrast, in this paper, we want to better understand the application area of CTC models and develop a methodology to calibrate these models from sales transaction data.

Our proposal here follows a different perspective and builds upon the modeling approach introduced by Manzini and Mariotti (2014). These authors study a choice model where the consideration set formation is stochastic and defined by the realization of the *attention parameter* of every alternative. This attention parameter is equivalent to our *propensity parameter* when the offer set is fully observable. After forming a consideration set, a consumer purchases the product that maximizes a preference relation within considered products. One of their main results states that this random choice rule is the only one for which the impact of removing an alternative on the choice probability of any other alternative, is asymmetric and menu-independent. The potential of the Manzini-Mariotti model in operational contexts was first evaluated by Gallego and Li (2017), who verify in a case study in the airline industry that its ability to fit booking data outperforms both the MNL and mixtures of MNLs in most of the markets evaluated. They also show that the related assortment optimization problem runs in polynomial time even with capacity constraints. In our

paper, we study a more general class of CTC models and focus on understanding their predictive power in different operations contexts.

In this section, we formally introduce *consider-then-choose* (CTC) models, followed by the presentation of related identification conditions.

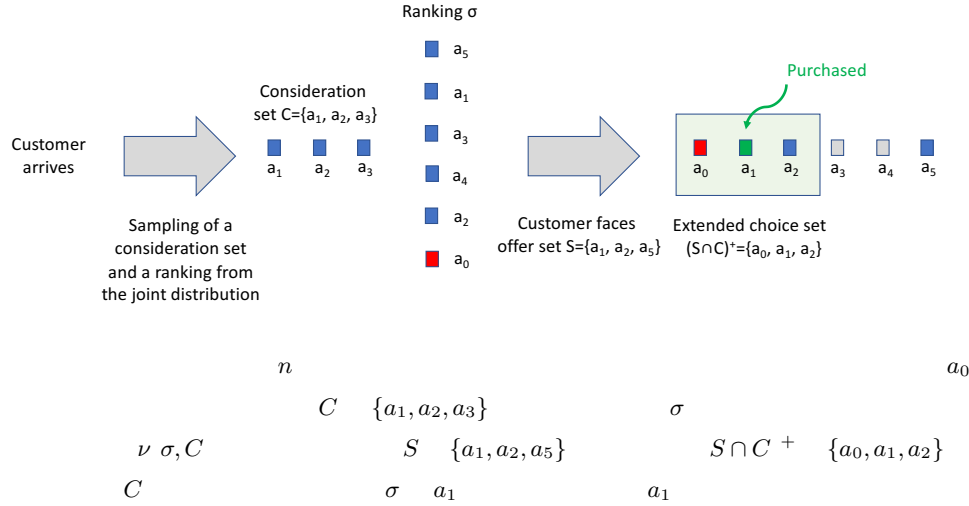
We consider a universe of products $\{ \}$, in addition to the no-purchase or outside option (we use the terms no-purchase option and outside option interchangeably throughout). We use to denote the set $\cup \{ \}$.

Customers arrive at the store sequentially. In each choice instance, a customer is presented with a subset \subseteq of products and chooses either one of the products in or the outside option . We let $\mathbb{P} ()$ denote the probability that a customer chooses product \in and $\mathbb{P} ()$ the probability that the customer chooses the outside option. Our goal is to model this choice process through a stochastic model that specifies all the choice probabilities $\{ \mathbb{P} () : \in \subseteq \}$, where we use to denote the set $\cup \{ \}$. We assume that the choice probabilities satisfy the standard probability laws: $\mathbb{P} () \geq 0$ for all \in and $\sum_{j \in +} \mathbb{P} () = 1$ for all \subseteq .

To explicitly account for the fact that customers may not consider all products on offer before making a choice, we assume that their choice behavior follows a two-stage *consider-then-choose* (CTC) model. In the first stage, the customer forms a consideration set \subseteq , and her own preference is realized (we describe this mechanism below), and in the second stage, selects either a product from the *choice set* \cap or the outside option . In other words, the customer picks one element from the extended choice set $(\cap) = (\cap) \cup \{ \}$.

In this model, for a product to be purchased, it must be both *offered and considered*. The seller restricts customers' choices by deciding the set of products to offer. But the customer further restricts her choices to just the ones in her consideration set . Two usual reasons to justify this approach are: (i) the customer has strong unobserved preferences (which prevent her from ever buying certain products), or (ii) the cognitive overload that prevents her from evaluating all the products on offer before choosing.

To describe customer preferences, let \mathcal{S} denote the set with cardinality ! of all full rankings or permutations of products in in which the outside option is ranked at the bottom. The *preference ordering* or *ranking* of the products in is described by a bijective ranking function $: \rightarrow \{1, \dots, + 1\}$ specifying a preference rank () for each product . The preference ordering induces an antireflexive, antisymmetric, and transitive preference relation \succ , defined



as \succ if and only if $(\cdot) \in \mathcal{S}$. Since we are fixing the position of the outside option to be at the bottom of each ranking, we have that $(\cdot) = \cdot + 1$ for all $\cdot \in \mathcal{S}$. Customer preferences are described by a joint probability distribution $\nu: \mathcal{S} \times 2^{\mathcal{S}} \rightarrow [0, 1]$, where $2^{\mathcal{S}}$ denotes the collection of all subsets of the set \mathcal{S} . In each choice instance, when confronted with an offer set S , a customer samples a ranking σ and a consideration set $C \subseteq \mathcal{S}$ with probability $\nu(\sigma, C)$ and chooses the product $\arg \min \{ (\cdot): \cdot \in (S \cap C)^* \}$.

Figure 1 illustrates the choice process for a particular store visit given a joint distribution function over consideration sets and full rankings from where the customer samples a consideration set $C = \{a_1, a_2, a_3\}$ and a ranking $\sigma = (a_5, a_1, a_3, a_4, a_2, a_0)$, and where we are representing the ranking as a tuple with the products listed in the order of their preference. The seller exhibits set $S = \{a_1, a_2, a_5\}$. In this choice instance, item a_1 is purchased. If the sampled consideration set had been $C = \{a_1, a_2, a_5\}$ instead, then the customer would have gone for the no-purchase option.

Formally, the choice probability $\mathbb{P}(\cdot)$ under this model is given by

$$\mathbb{P}(\cdot) = \sum_{C \in \mathcal{S}_n} \sum_{\sigma \in \mathcal{S}} \nu(\sigma, C) \cdot \mathbf{I}[\cdot \in (S \cap C)^*] \cdot \mathbf{I}[\cdot \succ \forall \cdot' \in (S \cap C)^* \cdot' \neq \cdot] \quad (1)$$

where $\mathbf{I}[\cdot]$ is the standard indicator function taking the value 1 if condition \cdot is satisfied, and the value 0 otherwise. We further assume that the empty condition $\cdot = \emptyset$ is always satisfied. We say that choice data $(\mathbb{P}(\cdot)): \cdot \in \mathcal{S} \rightarrow \mathbb{R}$, for some collection of subsets $\mathcal{S} \subseteq 2^{\mathcal{S}}$, is *consistent* with an underlying CTC model if there exists a distribution $\nu(\cdot, \cdot)$ that satisfies equation (1) for all $\cdot \in \mathcal{S}$ and $C \in \mathcal{S}$.

The sampling of the ranking σ as part of the choice process resembles the general random utility maximization (RUM) framework (Block and Marschak, 1960). The RUM class is the most

studied choice-based demand class in the literature and includes popular models such as the MNL, the nested logit (NL), and the mixture of MNLs (MMNL) models. At the core, it assumes that customers sample utility values for products from some underlying joint distribution and choose the product with the highest utility. A natural question that may arise is how RUM and CTC classes are related. We establish the following result.

PROPOSITION 1. *The collection of choice probabilities $\{\mathbb{P}(\cdot): \mathcal{C} \subseteq \mathcal{I}\}$ is consistent with an underlying RUM model if and only if it is also consistent with an underlying CTC model.*

The proof of the proposition is rather straightforward and is provided in Appendix A2 for the sake of completeness. The above result shows that the CTC model provides an equivalent but alternative parameterization of the RUM class and hence offers the same explanatory power. Yet, as we show below, this parameterization allows for a more compact representation of customers preferences, particularly, when they do indeed form consideration sets. It also allows us to use product and customer-level features to explicitly explain how customers form consideration sets and use data sources other than transaction data, such as survey data, to help infer customers consideration sets. Finally, it also seems particularly appropriate for some specific business contexts, as we will illustrate in subsequent sections.

The CTC model as stated above is not amenable to estimation, and therefore, we parameterize it as a mixture of what we call *independent consideration set* (ICS) models. As shown below, this parameterization is without loss of generality. A single-class ICS model is defined by a *single* preference ordering \succsim and a product form consideration set distribution. As above, the outside option is ranked at the bottom of the preference list; that is, $(\cdot) = \cdot + 1$. Further, each product $i \in \mathcal{I}$ is associated with an inclusion probability (or *propensity*) parameter $\alpha_i \in [0, 1]$, which denotes the probability that a customer includes product i in her consideration set. Customers make product inclusion decisions independently of each other, and therefore, the probability of sampling consideration set $\mathcal{C} \subseteq \mathcal{I}$ is $\mathbb{P}(\mathcal{C}) = \prod_{j \in \mathcal{C}} \alpha_j \prod_{j \in \mathcal{I} \setminus \mathcal{C}} (1 - \alpha_j)$. The ICS model was studied in [Manzini and Mariotti \(2014\)](#), and it can be shown that the probability of choosing product $i \in \mathcal{C}$ from offer set \mathcal{C} is

$$\mathbb{P}(i) = \begin{cases} \prod_{i \in \mathcal{C}} (1 - \alpha_i) & \text{if } i = \cdot \\ \prod_{i \in \mathcal{C}} \alpha_i \prod_{j \succ_{\sigma} i} (1 - \alpha_j) & \text{otherwise} \end{cases}$$

where the first expression corresponds to the probability of the event that none of the offered products are considered (which results in $(\mathcal{C} \cap \mathcal{I}) = \{\cdot\}$, and hence, the selection of the outside option), and the second expression corresponds to the probability of the event that product i is considered but none of the products preferred over i are. Finite mixtures of these models provide more modeling flexibility, and in fact, subsume the entire class of CTC models, as shown in the proposition below.

PROPOSITION 2. *Every CTC model with joint probability distribution (\cdot, \cdot) over consideration sets and full rankings can be represented as a finite mixture of ICS models with preference orders \succ , propensity parameters α , and mixture weights (w_i) for $i \in \{1, \dots, K\}$, and for some value of β .*

The proof of this result is rather intuitive and the argument is as follows. Consider any CTC model with underlying joint distribution (\cdot, \cdot) over consideration sets and preference lists. Let K denote the support size of the distribution (\cdot, \cdot) , so that $(\mathbf{c}, \mathbf{r}) \neq 0$ for some collection of tuples $\{(\mathbf{c}, \mathbf{r}) : \|\mathbf{c}\| = 1, \|\mathbf{r}\| = 1\}$. Then, we can define K independent ICS models with parameters $(\alpha_i, \beta_i, \succ_i)$ and w_i , for every $i \in \{1, \dots, K\}$ such that $w_i = 1$ if item \mathbf{c}_i belongs to the consideration set \mathbf{c}_i , and 0 otherwise. The mixing distribution over the K classes is defined by (w_i) such that $(w_i) = (w_i)$.

Our next goal is to study the problem of identifying customer consideration sets from sales transaction data alone.

In this section, we derive various conditions under which the CTC models are identifiable. To this end, our goal is to investigate if we can uniquely infer the parameters of the underlying CTC model from the collection of choice probabilities $\{\mathbb{P}(\mathbf{c}) : \mathbf{c} \in \mathcal{C} \subseteq \mathcal{U}\}$.

Our first results shows that the marginal distribution over consideration sets, defined as $(\mathbf{c}) = \sum_{\mathbf{c} \in \mathcal{S}_n} (\mathbf{c}, \mathbf{r})$, is uniquely identifiable from the observed choice probabilities alone. Specifically, we have the following result:

PROPOSITION 3. *Suppose that a collection of choice probabilities $\{\mathbb{P}(\mathbf{c}) : \mathbf{c} \in \mathcal{C} \subseteq \mathcal{U}\}$ are consistent with an underlying CTC model (\cdot, \cdot) . Then, the marginal distribution (\cdot) over consideration sets, defined as $(\mathbf{c}) = \sum_{\mathbf{c} \in \mathcal{S}_n} (\mathbf{c}, \mathbf{r})$, is uniquely identified:*

$$(\mathbf{c}) = \sum_{\mathbf{c} \subseteq \mathcal{U}} (-1)^{|\mathcal{U}| - |\mathbf{c}|} \mathbb{P}(\mathbf{c} \setminus \mathcal{U}) \quad \forall \mathbf{c} \subseteq \mathcal{U}$$

This is a rather surprising result. Intuitively, it seems that consideration sets cannot be identified from choice observations alone: if all we know is that a customer chose product \mathbf{c} from a set \mathbf{c} , then the customer may have considered any subset of products containing product \mathbf{c} , all the way from the singleton set $\{\mathbf{c}\}$ to the entire product universe \mathcal{U} ; we would not have any basis to select one subset over the other. While this intuition is correct for all the products in \mathcal{U} , the choice of the outside option reveals more.

To see this, note that in the definition of the CTC model class, we fix the position of the outside option to be at the bottom (i.e., at position $|\mathcal{U}| + 1$) of all the preference lists. This does not restrict

the generality of the CTC model class (cf. Proposition 1), yet provides us with valuable information on the customers' consideration sets. Specifically, because the outside option is the least preferred option, a customer chooses the outside option *only if* the sampled consideration set is disjoint from the offered set. In other words, if S is the offered set, then a customer choosing the outside option must have sampled a consideration set $C \subseteq S^c$. We, therefore, have that $\mathbb{P}(C) = \sum_{C \subseteq S^c} \lambda(C)$. For example, when $S = \emptyset$, it follows that a customer will choose the outside option only if they sample the empty set as the consideration set; therefore, $\mathbb{P}(C) = \lambda(\emptyset)$. Similarly, if $S = S \setminus \{i\}$, then the customer chooses the outside option only if their sampled consideration set is $\{i\}$ or the empty set, which implies that $\mathbb{P}(C) = \lambda(\{i\}) + \lambda(\emptyset)$. From this, we can back out the value of $\lambda(\{i\})$. More generally, we use a particular form of the inclusion-exclusion principle stated in [Graham \(1995\)](#) to back out the consideration set distribution. For any finite set S , if $f: 2^S \rightarrow \mathbb{R}$ and $g: 2^S \rightarrow \mathbb{R}$ are two real-valued set functions defined on the subsets of S such that $f(C) = \sum_{D \subseteq C} g(D)$ for every $C \subseteq S$, then the inclusion-exclusion principle states that $g(C) = \sum_{D \subseteq C} (-1)^{|C|-|D|} f(D)$ for every $C \subseteq S$. Our result then follows from replacing $f(C)$ with $\mathbb{P}(C)$ and defining $g(C) = \lambda(C)$. For completeness, we provide an alternative proof of this result from the first principles in [Appendix A2](#).

Although the marginal distribution over consideration sets is identifiable, the marginal distribution over rankings itself is not. This follows from an existing result in [\(Sher et al., 2011\)](#), which shows that a general distribution over rankings (which is one of the building blocks of our CTC class) is not uniquely identifiable from choice probabilities alone for $n \geq 4$. Therefore, the mixture of ICS models is also non-identifiable given the result stated in [Proposition 2](#).

Empirical evidence in the marketing literature suggests that the size of the consideration sets for most customers in different categories is relatively small, e.g., [Hoyer \(1984\)](#) concludes that the median number of laundry detergents that a consumer considers before making a purchase is one. When the size of consideration sets is bounded above by n , with $n \geq 4$, it follows from [Proposition 3](#) that to recover λ , we need choice probabilities under offer sets of size $n - 1$ or larger.

COROLLARY 1. *Consider a CTC model in which customers sample consideration sets of size at most n for some $1 \leq n \leq S$; that is, $\lambda(C) = 0$ whenever $|C| > n$. Furthermore, suppose that the no-purchase option is the least preferred product in all the preference lists in the support. Then, the distribution λ over consideration sets can be identified using choice probabilities under offer sets of size $n - 1$ or larger, i.e., from the collection $\{\mathbb{P}(C) : |C| \geq n - 1\}$.*

$$\mathbb{P}_0(S) = \sum_{C \subseteq N \setminus S} \lambda(C)$$

When the consideration sets are small, Corollary 1 argues that it is sufficient to collect choice probabilities for the no-purchase alternative from large offer sets. In many applications, however, firms cannot offer very large offer sets to their customers because of space constraints either in a physical store or on the relevant locations (e.g., top slots) of a website. The next proposition shows that when the consideration sets are of size at most k , then the consideration set distribution can be recovered using choice probabilities of offer sets of size at most k .

PROPOSITION 4. *Consider a CTC model in which customers sample consideration sets of size at most k for some $1 \leq k \leq n$. Let $\{\mathbb{P}(C): C \subseteq [n] \mid |C| \leq k\}$ be a collection of choice probabilities that are consistent with such a CTC model. Then, we have*

$$\mathbb{P}(C) = \sum_{\subseteq} \sum_{\supseteq \cup} (-1)^{|C| - |C \cap C'|} \cdot \mathbf{1}[|C \cup C'| \leq k] \cdot \mathbb{P}(C')$$

where $\Delta(C, C')$ denotes the symmetric difference $(C \setminus C') \cup (C' \setminus C)$.

The proof of the proposition is involved. It requires establishing several combinatorial identities. We present it in Appendix A2.

The results in the previous section focus on the recovery of the marginal distribution over consideration sets. To ensure complete identification, we need to restrict our CTC model further. We consider the class of models in which customers are homogeneous in their preference orderings, that is, the model is described by a single preference list, but is heterogeneous in their consideration sets. We call this model the *general consideration set* (GCS) model. Similar to the CTC class, we parameterize the GCS models as a mixture of ICS models sharing the same ranking.

More precisely, we assume that the GCS model is defined by a single preference ordering \succ and a distribution $\mathbb{P}(\cdot)$ over consideration sets. As before, the outside option is ranked at the bottom of the preference list, and a customer samples a consideration set according to \mathbb{P} and then chooses the most preferred product according to \succ from the set $(C \cap [n])$. With this restriction, we show that the choice rule is also identifiable:

PROPOSITION 5. *Suppose that the collection of choice probabilities $\{\mathbb{P}(C): C \subseteq [n] \mid |C| \leq 2\}$ are consistent with an underlying GCS model. Then, for all $1 \leq i \leq n$ and $i \neq j$, we have that if $\mathbb{P}(\{i\}) > \mathbb{P}(\{j\})$, then $\mathbb{P}(C \cup \{i\}) > \mathbb{P}(C \cup \{j\})$.*

The argument is intuitive. Given that there is a single preference list shared by all the consumers, if, for any given C alone, the presence of another i lowers its probability of being chosen, it is because j is preferred over it.

Unsurprisingly, the GCS model is not as rich as the CTC model class. In particular, it is a special case of the RUM choice rule.

PROPOSITION 6. *The GCS choice model is a special case of the RUM choice rule, that is, $\mathcal{C} \subseteq \mathcal{R}$, but $\mathcal{R} \not\subseteq \mathcal{C}$.*

The proof of the proposition is provided in Appendix A2 and uses the result in Proposition 5. It exhibits an example of a choice model that belongs to the RUM class but not to the GCS class.

Going back to Proposition 5, we assumed therein that the collection of observed choice probabilities is consistent with an underlying GCS model. To verify if that is indeed the case, we establish here a set of necessary and sufficient conditions that the observed choice probabilities must satisfy.

PROPOSITION 7. *The collection of choice probabilities $\{\mathbb{P}(S): S \in \mathcal{C}\}$ is consistent with a GCS model with unique parameters θ and consideration set distribution μ such that $\mu(S) > 0$ for all $|S| \leq 3$ and $\mu(S) \geq 0$ whenever $|S| \leq 3$ if and only if it satisfies the following conditions:*

Condition 1. For all offer sets $S \subseteq X$ and $S' \in \mathcal{C}$ such that $S \cap S' \neq \emptyset$: if $\mathbb{P}(S \setminus \{i\}) \neq \mathbb{P}(S')$, then it must hold that $\mathbb{P}(S \setminus \{i\}) = \mathbb{P}(S)$.

Condition 2. For all offer sets $S' \subseteq X$ and $S \in \mathcal{C} \cap S'$ such that $S \cap S' \neq \emptyset$: if $\mathbb{P}(S \setminus \{i\}) > \mathbb{P}(S)$, then it must hold that $\mathbb{P}(S' \setminus \{i\}) > \mathbb{P}(S')$; and if $\mathbb{P}(S \setminus \{i\}) = \mathbb{P}(S)$, then it must hold that $\mathbb{P}(S' \setminus \{i\}) = \mathbb{P}(S')$.

Condition 3. For all offer sets $S \subseteq X$, we have that $\sum_{i \in S} (-1)^{|S|-1} \mathbb{P}(S \setminus \{i\}) \geq 0$ with a strict inequality when $|S| \leq 3$.

Proposition 7 is similar to the set of conditions established in Manzini and Mariotti (2014) (see Theorem 1) for the case when the consideration set distribution μ has the product form due to the independence of the attention (or propensity) parameters. Our result consists of new conditions applied to a general consideration set distribution μ . Condition 1 is similar to the I-Asymmetry assumption in Manzini and Mariotti (2014), which states that either product i increases the sales of product j or vice versa, but not both (note that the increase may either be an increase or decrease). In other words, increase is one-directional and two products cannot increase the sales of each other. Condition 2 states that if product i increases the sales of product j in one offer set, then it must continue to do that in all the offer sets. That is, the direction of increase is consistent across all the offer sets. Condition 3 is a technical restriction to ensure the existence of a valid probability distribution function μ over the consideration sets. The strict inequality in Condition 3 is needed to ensure that the preference list over products in S satisfies the transitivity requirement. The proof of Proposition 7 is presented in Appendix A2. As it can be observed therein, establishing necessity is straightforward, but establishing sufficiency requires significant work.

In this section, we propose techniques to estimate the general CTC model and its restricted versions presented above, from sales transaction data. The building block is an algorithm to fit a single-class, ICS model, which we then extend using the expectation-maximization (EM) framework to fit a finite mixture of them.

Throughout the section, we assume access to sales data consisting of purchase transactions over periods. Every purchasing instance is represented by a tuple (\mathcal{I}_t, p_t) for $t \in \{1, \dots, T\}$, where \mathcal{I}_t denotes the subset of products offered in period t and p_t denotes the product purchased.

We highlight here that the CTC model is amenable to incorporating product features. In Appendix A3.1 we illustrate how to do it using three popular methods in machine learning: logistic, decision tree, and random forest- regressions.

To formulate the likelihood function under this model, we define binary linear ordering variables $\delta_{i,j}$, $\forall i, j \in \mathcal{I}_t, i \neq j$ where $\delta_{i,j} = 1$ if product i is preferred over product j in the preference list \succ (or, equivalently, $\delta_{j,i} = 0$), and $\delta_{i,j} = 0$ otherwise. Note that $\delta_{i,j}$ is an alternative parameterization of the preference order \succ . Then, the log-likelihood function under the single class ICS model can be shown to be

$$\mathcal{L}(\theta) = \sum_t \left[\log p_t + \sum_{\substack{k \in \mathcal{I}_t \\ k \neq p_t}} [\delta_{k,p_t} \log(1 - p_t)] \right]$$

and the maximum likelihood estimation (MLE) problem can be formulated as follows:

$$\begin{aligned} \text{s.t.:} \quad & \delta_{i,j} + \delta_{j,i} = 1 \quad \forall i, j \in \mathcal{I}_t, i \neq j \end{aligned} \quad (2)$$

$$\delta_{i,j} + \delta_{j,k} + \delta_{k,i} \leq 2 \quad \forall i, j, k \in \mathcal{I}_t, i \neq j \neq k \quad (3)$$

$$\delta_{i,j} \in \{0, 1\} \quad \forall i, j \in \mathcal{I}_t, i \neq j \quad (4)$$

$$0 \leq p_t \leq 1 \quad \forall t \quad (5)$$

where constraints (2) and (3) ensure that $\delta_{i,j}$ indeed represents a total order. In particular, the set of constraints (2) ensures that either i is preferred over j or vice versa, and the set of constraints (3) imposes the total ordering among any three products.

4.1.1. Estimation methodology. To be able to solve the above problem, we reformulate it as follows. First, we introduce a new variable $\theta_{i,j}$ defined as $\theta_{i,j} = \delta_{i,j} - \delta_{j,i}$, $\forall i, j \in \mathcal{I}_t, i \neq j$ and rewrite the likelihood function in the following way

$$\mathcal{L}(\theta) = \sum_t \left[\log p_t + \sum_{\substack{k \in \mathcal{I}_t \\ k \neq p_t}} \log(1 - p_t) \right]$$

Note that with this change of variable, the log-likelihood function becomes jointly concave in θ, τ, δ . We can then formulate the MLE problem in terms of the variables (θ, τ, δ) :

$$\begin{aligned} \text{s.t.:} \quad & \theta \tau \delta \leq \mathcal{L}(\theta, \tau, \delta) \quad \forall \theta, \tau, \delta \end{aligned} \quad (6)$$

$$\theta \tau \delta \leq \tau \delta \quad \forall \theta, \tau, \delta \quad (7)$$

$$\theta \tau \delta \geq \theta + \tau - 1 \quad \forall \theta, \tau, \delta \quad (8)$$

$$\theta \tau \delta \geq 0 \quad \forall \theta, \tau, \delta \quad (9)$$

(θ, τ, δ) satisfy (2) – (5)

where linear constraints (6)-(9) ensure that $\theta \tau \delta = \tau \delta$, $\forall \theta, \tau, \delta$, given that θ is a binary variable, and constraints (2) and (3) ensure again a total order on θ, τ, δ . We reformulate the MLE problem to have a linear objective function:

$$\begin{aligned} & \theta \tau \delta \quad (10) \\ \text{subject to } & (\theta, \tau, \delta) \text{ satisfy (2) – (5) and (6) – (9)} \end{aligned}$$

$$\theta \tau \delta \leq \mathcal{L}(\theta, \tau, \delta) \quad (11)$$

Note that if we know the ranking θ, τ, δ , then optimization problem (10) reduces to solving a globally concave maximization problem with a unique, closed form solution given by

$$\theta \tau \delta = \frac{\sum \mathbb{I}[\theta = \tau] + \sum \mathbb{I}[\theta \in \tau \succ \delta]}{\sum \mathbb{I}[\theta = \tau] + \sum \mathbb{I}[\theta \in \tau \succ \delta]} \quad (12)$$

Next, we show how to apply the outer-approximation method of [Duran and Grossmann \(1986\)](#) to solve the optimization problem (10)-(11). The proposed algorithm effectively exploits its structure, where we have linearity of the constraints involving the binary variables θ, τ, δ , and convexity of the non-linear constraint (11) which only depends on continuous variables θ, τ, δ . In order to linearize the optimization problem, we use the outer-approximation of a convex set by the intersection of the collection of its supporting half-spaces. The broad idea of this algorithm is to approximate the convex constraint in the MINLP (i.e., constraint (11)) with a set of linear constraints. As a result, solving the MINLP reduces to solving a sequence of MILPs, where at each iteration we add only one linear constraint to the MILP formulated in the previous iteration. Next, we provide the details of how to apply the outer-approximation algorithm to our MLE problem.

Let $\mathcal{C}[\theta, \tau, \delta]$ denote a constraint which is a linear approximation of the constraint (11) at a point (θ, τ, δ) , i.e.,

$$\mathcal{C}[\theta, \tau, \delta] := \left\{ \theta \tau \delta \leq \mathcal{L}(\theta, \tau, \delta) + \sum \frac{1}{t} \cdot (\theta - \tau) + \sum \sum_{\substack{k \in t \\ / t}} \frac{1}{t-1} \cdot (\tau - \delta) \right\}$$

Then we define the following optimization problem

$$\begin{aligned} & \theta \tau \delta \\ \text{subject to } & (\theta, \tau, \delta) \text{ satisfy (2) – (5) and (6) – (9)} \end{aligned} \tag{13}$$

$$\begin{aligned} & \mathcal{C}[\theta, \tau, \delta] \forall (\theta, \tau, \delta) \in \mathcal{A} \\ & \leq \leq \end{aligned} \tag{14}$$

where we have replaced constraint (11) with the finite collection of linear constraints (14) at points in a set \mathcal{A} that is incrementally built. It follows from the convexity of constraint (11) that every point that satisfies constraint (11) also satisfies the collection of constraints (14) for every finite set \mathcal{A} . We thus obtain an outer approximation. We also add bounds θ and τ to the log-likelihood function, which will be chosen in each iteration to tighten the interval containing the solution. This outer approximation defines the optimization subproblem as an MILP. Because of the potentially many continuous points required for outer-approximation, we solve a sequence of MILPs to build up increasingly tight relaxations of the original MINLP. The algorithm to calibrate the ICS model is provided below.

ICS model calibration algorithm

Input Given sales transaction data, do:

Step 1 Sort products in decreasing number of sales and let $\pi^{(0)}$ (i.e., $\pi^{(0)}$) denote the corresponding ranking. Compute $\mathcal{L}^{(0)}$ using equation (12) given $\pi^{(0)}$.

Step 2 Obtain all possible rankings $\pi^{(1)}, \pi^{(2)}, \dots, \pi^{(m)}$, by swapping positions of any pair or two pairs of items in $\pi^{(0)}$. Compute $\mathcal{L}^{(i)}$ using equation (12) given $\pi^{(i)}$ for all $i \in \{1, 2, \dots, m\}$.

Step 3 Set $\mathcal{L}_{kj}^{(i)} := \mathcal{L}_{kj}^{(i)} / \mathcal{L}_k^{(i)}, \forall k, j$ and for all $i \in \{0, 1, 2, \dots, m\}$. Recall that $\mathcal{L}^{(i)}$ is an alternative parameterization of the ranking $\pi^{(i)}$. Set $\mathcal{L}_L := \min_{0 \leq i \leq m} \mathcal{L}(\pi^{(i)})$ and $\mathcal{L}_U := \max_{0 \leq i \leq m} \mathcal{L}(\pi^{(i)})$. Set $\epsilon := \epsilon$.

Step 4 While $|\mathcal{L}_U - \mathcal{L}_L| > \epsilon$ and running time did not exceed the limit, do:

- Set $\epsilon := \epsilon + 1$.
- Solve optimization problem (13) with set $\mathcal{A} = \{(\pi^{(k)}, \mathcal{L}^{(k)})\}_{k=0}^{i-1}$ and obtain solution $(\theta^{(i)}, \tau^{(i)})$. Note that in each iteration, we only add one constraint to the optimization problem solved in the previous iteration.
- Update $\mathcal{L}^{(i)}$ using equation (12) given $(\theta^{(i)}, \tau^{(i)})$.
- Set $\mathcal{L}_{kj}^{(i)} := \mathcal{L}_{kj}^{(i)} / \mathcal{L}_k^{(i)}, \forall k, j$.
- Set $\mathcal{L}_L := \min \{\mathcal{L}(\pi^{(i)}, \mathcal{L}^{(i)})\}$ and $\mathcal{L}_U = \mathcal{L}(\pi^{(i)}, \mathcal{L}^{(i)})$.

Endwhile

Step 5 Find solution $\pi^* = \pi^{(i^*)}$ and $\mathcal{L}^* = \mathcal{L}^{(i^*)}$ where $i^* := \arg \max_{0 \leq k \leq i} \mathcal{L}(\pi^{(k)}, \mathcal{L}^{(k)})$.

Step 6 Stop.

Overall, the proposed algorithm consists of solving a finite sequence of MILPs. The size of each MILP scales in $\mathcal{O}(n^2)$, quadratically in the number of variables and cubically in the number of constraints. It follows from existing results in Duran and Grossmann (1986) that this algorithm converges to the global optimum in the long run.

Empirically, we analyze the performance of the proposed algorithm to estimate the ICS model on the IRI Academic dataset (to be described later in Section 6). We limited its running time to 3 hours, and the precision was set to $\epsilon = 10^{-6}$. It follows from Figure A4 in Appendix A3 that the optimality gap of the outer approximation algorithm (1) to calibrate the ICS model is 3.3% on average over 20 product categories, determined by the time limit, which indicates that this algorithm provides quite a reasonable performance in our setting. We also implemented a cutting plane method as an alternative algorithm (details also provided in Appendix A3), but the results were of slightly lower quality.

To make estimation tractable, we represent the general distribution $\mathcal{P}(\cdot)$ using a finite mixture of ICS models and then apply the EM algorithm. In a nutshell, this method starts from arbitrary parameter estimates \mathbf{x}^0 . Then, it computes the conditional expected value of the log-likelihood function $E[\log \mathcal{L}(\mathbf{x}) | \mathbf{x}^0]$ (the E, expectation, step). Next, the resulting expected log-likelihood function is maximized to compute new estimates \mathbf{x}^1 (the M, maximization, step), and both steps are repeated within a loop until convergence or a time limit is reached, to get a sequence of estimates $\{\mathbf{x}^k\}_{k=1}^T$.

The log-likelihood function to be maximized can be represented in the following way:

$$\log \mathcal{L}(\mathbf{x}) = \sum_{t=1}^T \log \left(\sum_{j \in \mathcal{I}_t} w_j \prod_{h \in \mathcal{H}_t} (1 - x_{jh})^{I_{jh,t}} \right)$$

where $w_j \geq 0$ is the weight of the class j , s.t. $\sum_{j \in \mathcal{I}_t} w_j = 1$; \mathcal{I}_t denotes the set of offered items at time t ; \mathcal{H}_t denotes the product purchased at time t ; and $I_{jh,t}$ denotes the number of transactions.

We next briefly outline the E-step and M-step of every iteration and how we start the algorithm in the context of our CTC estimation problem. Further specific details are relegated to Appendix A3.5.

Initialization: we randomly allocate sales transaction to one of the K classes which allows us to compute initial mixing point probabilities w_j , and estimate initial parameters: x_{jh} and θ_{jh} for all $j \in \{1, \dots, K\}$. To this end, we use the outer approximation algorithm for the ICS model on each class.

Then, we iterate over the sequence of the following E- and M-steps:

E-step: we compute the probability of every transaction at time t to come from a segment customer, based on the parameter estimates $(\hat{\theta}_t^1, \hat{\theta}_t^2, \dots, \hat{\theta}_t^K)$ of the previous iteration, and the transaction data.

M-step: first, we update the mixing distribution $(\hat{\pi}_t^k)$ of every segment $k \in \{1, 2, \dots, K\}$ based on the E-step membership probabilities, and then optimize the resulting conditional expected value of the log-likelihood function by using the outer approximation algorithm on each segment. This way we obtain $\hat{\theta}_t^k$ and $\hat{\pi}_t^k$ for all $k \in \{1, 2, \dots, K\}$.

Even though it is well acknowledged that the convergence of EM algorithms is not guaranteed a priori, it was verified consistently in all our experiments.

A key observation is that we can calibrate the CTC model by estimating the following parameters at the segment level: (1) membership probability π^k , (2) propensity parameter θ^k , and (3) ranking σ^k . Therefore, the approach is convenient when we have a small to moderate support of customer segments.

In a similar way, we can also calibrate the GCS model with the EM algorithm (see Appendix A3.4 for details). Note that it is very straightforward to build upon Proposition 2 and show that GCS model can also be represented as a mixture of ICS models sharing the same rank list σ .

For both the CTC and the GCS models, in our implementation, we consider a mixture of up to $K = 5$ classes and report out-of-sample results for the number of mixtures that drove to best in-sample results.

In this section, we describe the results of an extensive simulation study, the main purpose of which is to characterize the performance of CTC models relative to the classical RUM models under various noise regimes. We find that CTC models are generally more robust to noise in the order sets and outperform the classical RUM models when noise affects training and test order sets differently.

To streamline the analysis of this simulation study, we assume three ground truth, RUM-based models of demand: the classical MNL, the rank list, and the LC-MNL. To showcase the potential of CTC models, we start here from the restricted, single-class, ICS model. Given the similarity of the insights obtained, here we report results for the MNL ground truth model and defer results based on the rank-based and LC-MNL ground truth models to Appendix A4.1.

In our simulations, customers have perfect information of the order sets and consider all the items on order. Given the order set S , the customer chooses product i with probability $\frac{\theta^i}{1 + \sum_{i \in S} \theta^i}$, where the parameter $\theta^i > 0$ is the weight or the attraction value corresponding to product i , and the 1 stands for the weight of the outside option. The modeler observes customer choices but

does not observe the order sets perfectly. In fact, the assumed order sets could be a superset of the true order set, consistently with the common inventory inaccuracy problems. In the presence of such noise, we compare the predictions of an MNL model against an ICS model, both fitted to the choice observations, to understand the conditions under which one outperforms the other one.

In our setup, the benchmark MNL model does not suffer from model misspecification but does suffer from noise in the order sets. The ICS model, on the other hand, suffers from both model misspecification and noise in the definition of the order sets, but it is designed to handle the latter more effectively.

Our main finding is that the ICS model significantly outperforms the ground-truth MNL model when the noise is *asymmetric* between the training and test data sets. In other words, if product availability is hard to predict (because it might look different from the training data), then models based on the consider-then-choose framework outperform classical RUM models.

We assume that we have $N = 15$ items in the product universe. For each product i , we sample its nominal utility u_i uniformly at random from the interval $[1, 2]$ and set its MNL weight $w_i = \exp(u_i)$. We normalize the MNL weight of the outside option to 1. We parameterize the level of noise in the definition of the order set using two parameters: the *noise exposure* parameter, $\alpha \in [0, 1]$, and the *noise intensity* parameter, $\beta \in [0, 1]$. The noise exposure parameter determines if a product is exposed to noise. Intuitively, the noise exposure parameter is designed to control the degree of asymmetry in the level of noise between the training and testing datasets. The noise intensity parameter specifies the conditional noise level, as described below. Next, for given values of α and β , and realized parameters θ of the MNL model, the data simulation proceeds as follows:

1. We sample 100 order sets, $\{O_j\}_{j=1}^{100}$, uniformly at random, which correspond to the underlying true order sets.
2. For each order set O_j , we generate 10,000 sales transactions according to the MNL model with parameter values θ .
3. We generate a single exposure set E containing products exposed to noise by including each product $i \in I$ with probability α .
4. We generate noisy order set observations by modifying each true order set O_j . To this end, we add extra products from the exposure set E with probability β ; specifically, we obtain the noisy order set \tilde{O}_j by adding each product $i \in I \setminus O_j$ with probability β to the set O_j .

We generated both training and test datasets for 200 different combinations of α and β : $\alpha \in \{0.05, 0.1, \dots, 1\}$ and $\beta \in \{0.1, 0.2, \dots, 1\}$. First, we generated 100 realizations of MNL parameters.

Next, for given values α and β and for each MNL ground truth realization, using the procedure above, we generated both training and test synthetic instances of the same size. Each instance consists of 100 randomly generated order sets, and for each order set, 100,000 transactions, giving a total of 10,000,000 transaction records. All our algorithms were coded in Python (version 2.7.2) using Gurobi (version 7.0) as the optimization engine, and run on a 3.0Ghz processor with 16GB of RAM.

Our simulation setup is designed to capture not only different noise intensities but also different degrees of noise asymmetry between the training and test datasets. The noise exposure parameter captures this asymmetry as follows. Letting \mathcal{T} and \mathcal{S} denote the training and test exposure sets, respectively, the cardinality of the symmetric difference $(\mathcal{T} \setminus \mathcal{S}) \cup (\mathcal{S} \setminus \mathcal{T})$ captures the number of products that are exposed to noise in only one of the datasets. In expectation, this cardinality is equal to $2(1 - \alpha)$ because the probability that a product is exposed to noise in only one of the data sets is equal to $(1 - \alpha) + (1 - \beta) = 2(1 - \alpha)$. Therefore, the degree of asymmetry in the level of noise between the training and testing datasets is highest when $\alpha = 0.5$ and lowest when $\alpha = 0$ or $\alpha = 1$. Furthermore, the order sets are perfectly observed in both the training and test datasets when $\alpha = 0$ or $\alpha = 1$, and the noise is perfectly symmetric when $\alpha = 1$. In practice, noise may be asymmetric when calibrating future demand forecasts because promotion strategies, product availability, or replenishment processes may change over time.

We evaluate our models on two standard metrics, the mean absolute percentage error (MAPE) and the root mean square error (RMSE), defined as follows (in percentage points):

$$\text{MAPE} = \frac{100}{|\mathcal{I}|} \sum_{j \in \mathcal{I}} \frac{|\hat{y}_j - y_j|}{10 + y_j} \quad \text{RMSE} = \frac{100}{\sum_{j \in \mathcal{I}} y_j} \sqrt{\frac{1}{|\mathcal{I}|} \sum_{j \in \mathcal{I}} (y_j - \hat{y}_j)^2} \quad (15)$$

where y_j denotes the observed sales for $j \in \mathcal{I}$ in the test dataset and \hat{y}_j denotes our prediction. We make predictions on noisy test order sets, so $y_j = 10000 \cdot \sum_{o \in \mathcal{O}} p_{o,j}$, where \mathcal{O} is the noisy test order set and $p_{o,j}$ is the probability that product j will be purchased from order set o under the fitted model. Note that we add 10 in the denominator of MAPE to deal with undefined instances, and we divide the RMSE score by the total number of observed sales in the test data set to make it a relative metric so that it is more interpretable.

Intuitively, both scores quantify the power of the model combinations to predict the market shares for each product, with lower scores indicating a better prediction accuracy.

In Figure 2 we present a heatmap of the prediction scores under an MNL model fitted to the training data where each column corresponds to a particular noise intensity δ and each row corresponds to a particular noise exposure ϵ . We focus on the MAPE and RMSE prediction scores in the left and right panels, respectively. Recall that the MNL model is the ground truth model for this simulation study. As expected, the MNL model captures the ground-truth choice probabilities almost exactly when $\delta = 0.05$ and $\epsilon = 0.1$, i.e., there is only a small amount of noise added to the sales transaction data. However, the MNL prediction scores worsen with higher noise intensity for a given level of noise exposure. Interestingly, it can also be seen that MNL prediction scores are not monotonic with respect to the noise exposure level, i.e., the scores first increase and then decrease with the noise exposure level, for a given noise intensity. The fitted MNL model performs the worst at noise exposure level $\epsilon = 0.5$, which as noted above, corresponds to the highest degree of noise asymmetry between the training and test data sets, whereas the performance is relatively good even at high noise intensity levels when the degree of asymmetry is low (δ is close to 0 or 1).

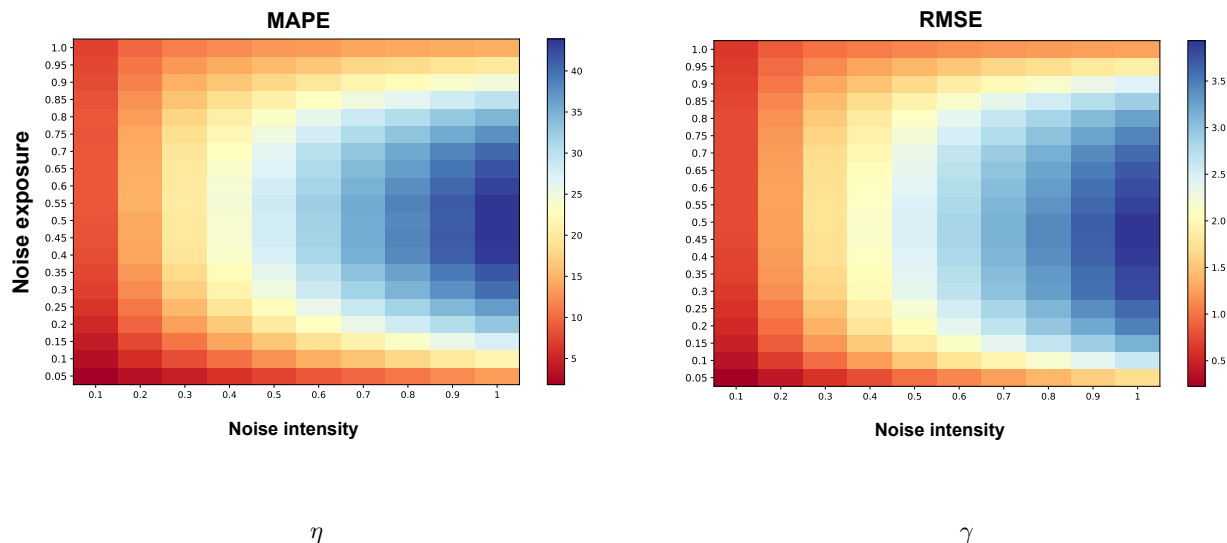
To provide a quantitative understanding of the variation of the model performance with respect to the two noise parameters δ and ϵ , we carry out the following linear regression:

$$Y_{\delta, \epsilon} = \beta_0 + \beta_1 \delta + \beta_2 \delta^2 + \beta_3 \epsilon + \beta_4 \text{Asymm} + \beta_5 \text{Shared} + \epsilon_{\delta, \epsilon} \quad (16)$$

where the index δ stands for the noise intensity $\delta \in \{0.1, 0.2, \dots, 1\}$, and the index ϵ stands for the noise exposure $\epsilon \in \{0.05, 0.1, \dots, 1\}$. The outcome variable $Y_{\delta, \epsilon}$ is the MAPE prediction score of the corresponding cell.

The first two terms in the regression capture dependence on the noise intensity level δ and the last two terms capture dependence on the noise exposure level ϵ . We add a quadratic term in the noise intensity level to capture any potential non-linear response to the noise intensity level. We also separate out the dependence on ϵ into degree of asymmetry and degree of overlap in noise. The covariate Asymm is the probability that an item in the product universe is exposed to noise only in the test dataset or only in the training dataset, and the covariate Shared is the probability that an item in the product universe is exposed to noise both in the test and training datasets. Note that $\text{Asymm} = (\delta - \epsilon) + (1 - \delta) + (1 - \epsilon) = 2(1 - \delta)$ and $\text{Shared} = \delta\epsilon$, where δ is the noise exposure.

The results for the regression (16) are reported in the last column of Table 1. It follows from there that the noise intensity δ deteriorates the predictive performance of the MNL model in a non-linear way, with the coefficient for the linear term being positive and the coefficient for the quadratic term being negative. The variables Asymm (i.e., degree of noise asymmetry) and Shared (i.e., the degree



of the shared noise) are positively correlated with the MAPE score, which also implies that the prediction performance of the MNL model worsens as the number of items in the product universe that are exposed to noise increases. Interestingly, the coefficient of the variable η has more than seven times higher magnitude than the coefficient of the variable γ , even though both variables have a comparable range of possible values (e.g., η varies from 0 to 0.5 and γ varies from 0 to 1), which indicates that the benchmark (i.e., MNL model) struggles the most in making accurate predictions when the impact of the noise is asymmetric between the training and test sales transactions. Note that the independent variables in the regression model (16), included as Model (5) in Table 1, explain most of the variation in the MAPE score under the MNL model, i.e., $R^2 = 0.93$.

These regression results indicate that it is not the noise per se, but the asymmetry in noise that is hurting model performance. The reason is that when noise is symmetric, model estimates are biased but the bias is in the correct direction. For example, if a product is frequently stocked out but the model does not know about it, then the model attributes low sales to a low attraction value as opposed to a stockout, resulting in an MNL weight that is biased downward. If the product continues to be frequently stocked out in the test data set, then the model should continue to forecast low sales, which would happen with a downward-biased MNL weight. On the other hand, if the stockout frequency is asymmetric and, say, reduces in the test data set (perhaps because of better replenishment), then the MNL model forecasts would be incorrectly biased downward.

In Figure 3, we compare the performances of the MNL and the ICS models under different noise regimes. Each cell presents the improvement obtained by the ICS model over the MNL model, computed as the difference in the corresponding prediction scores, so that higher values are better

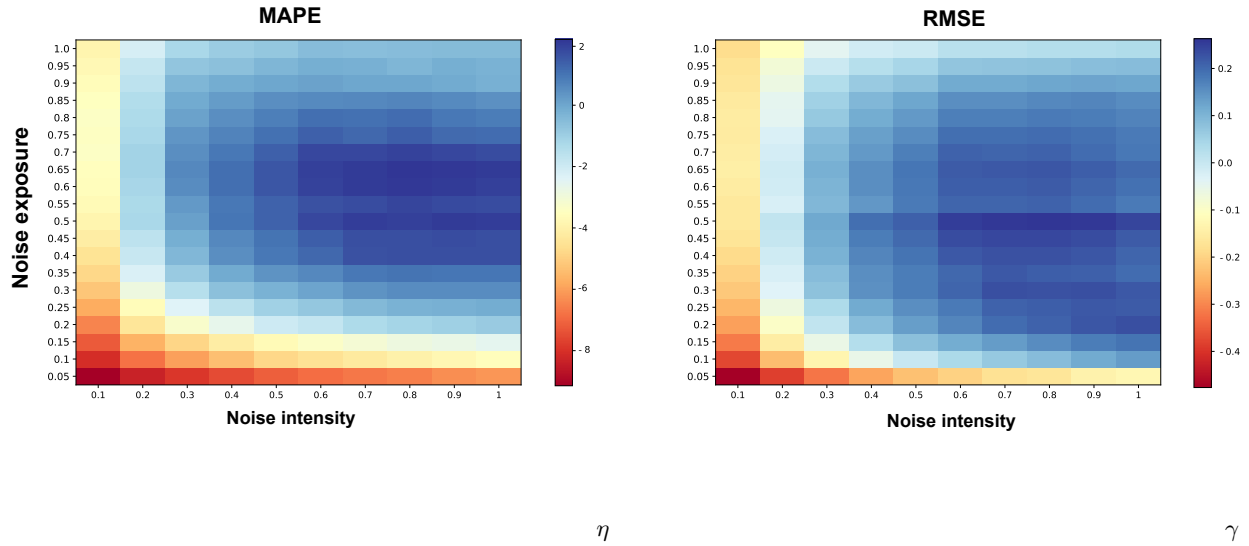
	Model (1) Score	Model (2) Score	Model (3) Score	Model (4) Score	Model (5) Score
2	27.336*** (17.997)	22.929*** (15.755)	34.662*** (8.677)	-2.590 (-1.143)	42.079*** (14.246)
	6.260*** (6.643)	12.467*** (17.021)	9.770*** (6.698)	22.224*** (20.645)	-13.403*** (-5.122)
					39.320*** (28.338)
					5.385*** (8.015)
					-11.694*** (-12.579)
No. Observations:	200	200	200	200	200
R-squared:	0.621	0.556	0.275	0.007	0.929
Adj. R-squared:	0.619	0.554	0.272	0.002	0.928

t
* $p < .$ ** $p < .$ *** $p < .$

for the ICS model. Unsurprisingly, the MNL model is significantly better than the ICS model when the level of noise is small (small values of both α and β) because the ICS model is misspecified. But, at high values of noise intensity α , and values of β close to 0.5, we see that the ICS model outperforms the MNL model. Table 2 presents the results from regressing the improvement scores according to regression equation (16). As above, we find that the benefits of the ICS model are most significant when the noise is asymmetric between the training and test order sets. Although not shown here, the pattern of the heatmap for the ICS model looks similar to that of the MNL model in Figure 2. However, the ICS model is more robust to noise (both the intensity and the degree of asymmetry) and its performance does not deteriorate as much, because of which it outperforms the MNL model in asymmetric and high noise regimes.

We highlight that the results in this section are robust to different asymmetric scenarios of noise as well (see Figure A10 in Appendix A4.1).

In order to shed some light on one potential mechanism that drives the superior performance of ICS over the MNL model, we present a stylized study in Appendix A4.2. The main insight we obtain is that the one-dimensional cannibalization property of the ICS model makes it robust to order set noise. This property states that the presence or absence of lower-ranked products does not affect the demand for higher-ranked products. As a result, order set noise mainly impacts the lower-ranked products, limiting the overall error rates. See Appendix A4.2 for details.



	Model (1) Impr.	Model (2) Impr.	Model (3) Impr.	Model (4) Impr.	Model (5) Impr.
2	4.162*** (7.883)	3.025*** (6.157)	7.311*** (7.112)	2.539*** (4.811)	16.621*** (17.214)
	-2.966*** (-9.055)	-1.842*** (-7.449)	-3.108*** (-8.281)	-1.588*** (-6.332)	-11.327*** (-13.241)
					11.530*** (25.421)
					4.878*** (22.207)
No. Observations:	200	200	200	200	200
R-squared:	0.239	0.161	0.203	0.105	0.874
Adj. R-squared:	0.235	0.156	0.199	0.100	0.871

t
 * $p < .$ ** $p < .$ *** $p < .$

Our analysis so far has been based on synthetic data, establishing that even the restricted ICS model outperforms classic choice models when o er sets are uncertain and the noise is asymmetric between the training and test o er sets. We now present a real-world case study where such conditions are met.

For our study here, we use the household purchase panel and store data from the IRI Academic Dataset (Bronnenberg et al. (2008)). This panel dataset keeps track of the household purchase histories for grocery and drug store chains, collected from the two largest Behavior Scan markets in the US over the years 2001-2011. We also use the GCS model, more general than the previous ICS, as a baseline to conduct our analysis. Because we now have access to panel data, we also implement a variant of the EM algorithm for the estimation of the parameters of the GCS model, as described in Appendix A3.4.2.

The purpose of this empirical study is threefold: (i) provide various real-world scenarios based on the IRI dataset where we are likely to face significant noise in the observed definitions when making the long-term demand predictions, (ii) investigate the prediction performance of choice models under different noise regimes, e.g., quantify the improvement of the GCS model over the latent class MNL (LC-MNL) model on panel data under several real-world scenarios with various noise intensities, and (iii) compare the GCS model studied in this paper with the more restricted ICS model of Manzini and Mariotti (2014).

Our main findings are as follows: (a) the improvements of GCS versus the benchmark LC-MNL are higher for scenarios in which the observed sets have a high level of noise, (b) the predictive performance of the GCS model is robust to the noise level in the observed sets, and (c) the GCS model significantly outperforms the ICS model in prediction accuracy, indicating that the independent consideration set model is indeed restrictive.

In this section, we present all the comparisons with respect to the GCS model (see the GCS model estimation details in Appendix A3.4.2). It is natural to wonder if the more general CTC model (see Section 4.2 above for the CTC model estimation framework), which allows for heterogeneity in customer preference orders, offers additional gains in prediction performance. To this end, we carried out a similar analysis with the CTC model; see Appendix A3.6 for details. From our results, we can not claim dominance of GCS or CTC; that is, the GCS model dominates the CTC model under the first noise generation process whereas the CTC model outperforms the GCS model under the second noise generation process. Based on this finding, we focus our analysis and discussions mainly on the more parsimonious GCS model.

The dataset consists of weekly sales transactions. We analyze a total of 20 categories, presented in Table A1. We focus on sales transaction data from the calendar year 2007. For every store visit, we are given the following information: the Universal Product Code (UPC) and price of the purchased item, a binary indicator if the product is on price or display promotion, the purchased quantity, the customer ID, the store ID, and the week when the purchase was made. Since we

are not given explicit information about the subset of items offered to each individual upon her store visit, we follow existing literature (e.g., Jagabathula and Vulcano (2018)) and construct this subset by aggregating all the transactions made in a particular store within a given category during a particular week. We also aggregate items with the same vendor code (i.e., items that are associated with the same brand name) into a single product due to data sparsity and divide the sales transaction data into two parts: the training set, which consists of the first 26 weeks of sales observations; and the test set, which consists of the last 26 weeks of sales observations. Note that we exploit the panel data structure in the predictive performance analyses below.

We compare the GCS model with the benchmark: the LC-MNL choice model with K latent classes. In this model, each customer belongs to one unobservable class, and customers from class $c \in \{1, 2, \dots, K\}$ make purchases according to the MNL model associated with that class. The model is described by the parameters of the MNL model characterizing each class and by the prior probabilities of customers belonging to each of the classes. Once the model parameters are estimated, we make transaction-level predictions for each customer by averaging the predictions from K single-class models, weighted by the posterior probability of class membership. Similarly to the GCS model, we estimated the model for $K = 1, 2, \dots, 5$, and report the best performance measure from these 5 variants, for each of the performance metrics introduced in Section 5.2. Because we have panel data, we make individual-level predictions and consequently, the number of predicted sales for product i is now calculated as
$$y_i = \sum_{c \in C} \sum_{j \in \mathcal{C}_i} (p_{ij}^c)$$
, where C is the set of all customers, \mathcal{C}_i is the number of transactions of customer c , and p_{ij}^c is the predicted probability that customer c purchases item i from noisy offer set estimate \mathcal{O}_i^c . Note that we do not have access to the no-purchase observations in our dataset and thus our prediction metrics do not include them.

Demand predictions or forecasts using choice models implicitly involve two steps: (a) forecasting the offer set and (b) predicting the demand for each product given the forecasted offer set. Most existing literature on choice models has only focused on the second step, implicitly assuming that the future offer set is accurately specified. In practice, one must also forecast future offer sets and these forecasts often contain errors. In our study, we follow these two steps explicitly to study the impact of forecast errors in the first step on the overall accuracy of predictions.

6.3.1. Different data aggregations towards offer set forecasting. We consider three different prediction tasks that naturally arise in practical retail contexts and which differ in the level of difficulty of forecasting the offer sets.

1. *Short-term forecasts.* Often, store managers want to make short- or immediate-term forecasts, say, for the next week, to help with inventory and promotion planning. For these forecasts, the manager has a reasonably accurate estimate of the product assortment or offer set.
2. *Long-term forecasts.* To be successful in major strategic and investment decisions, a store manager must also make long-term demand forecasts, say, over the next quarter or the next year. For making these forecasts, the retailer often does not have a good sense of how the product assortment evolves over the forecasting horizon because of product replenishment and stockouts, which manifest as errors in offer set forecasts.
3. *Warehouse forecasts.* Another scenario is when the warehouse of the retail chain distributes products to the stores and makes centralized decisions on the inventory level in the warehouse. In this case, the warehouse is likely to make predictions at the centralized level without knowing the up-to-date information on product assortments in every store.

In each of these scenarios, the retailer has access to the same purchase observations to train models on. It is the first step, that of forecasting the offer sets, that these scenarios differ on. As we shall see, the CTC models and the classic choice models differ on their ability to deal with noise in offer set forecasts.

We simulate these scenarios using the purchase observations we have as follows. As mentioned above, we split the purchase transactions into training and test sets. In the training set, the purchased product and the corresponding offer set are known. Because of the way the offer set was inferred, it may contain errors, but we train all the models disregarding any potential errors. Note that this process reflects the standard way in which choice models are trained in the literature and in practice. In addition to a trained choice model, the retailer must build what we call an OSForecaster or an *offer set forecaster*. We abstract away from details of how such a forecaster might be constructed and instead peek in our test dataset to simulate one. Specifically, let \mathcal{O} denote the set of offered products at store s and period t in the test horizon. Further, let \mathcal{T} denote the set of test time periods.

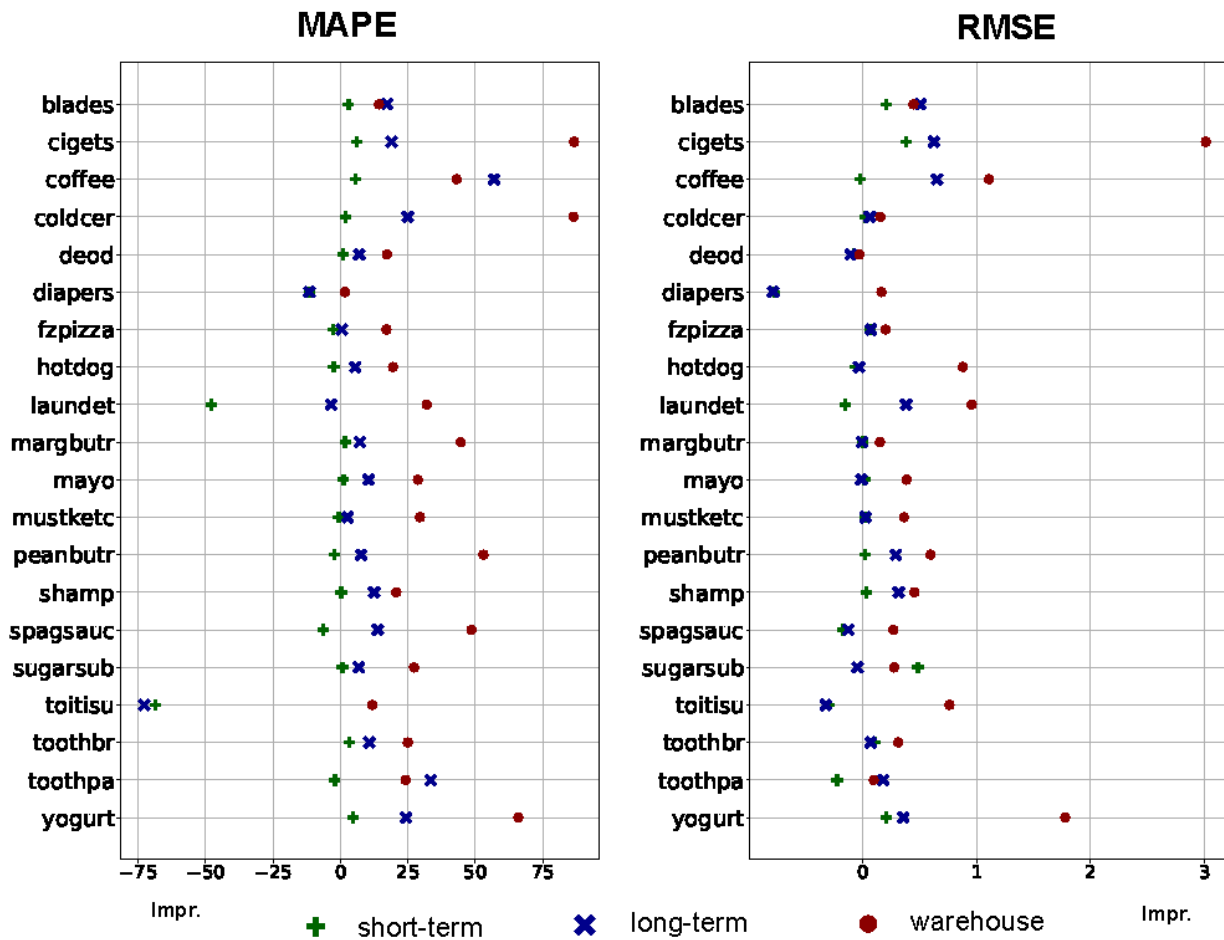
For short-term forecasts, we provide the retailer with the number of purchasing customers at each store s and period t , and the retailer must predict the sales for each of the products. In practice, this will entail making sales predictions for the next week for each of the stores. To make this prediction, the retailer must first forecast an offer set for each store and time period during the test horizon. To reflect the fact that the retailer makes few errors in forecasting the next week's offer set, we assume that the short-term OSForecaster returns the true offer set $\mathcal{O}_{s,t}$ when queried

with a store s and time period t . The predictions of the OSForecaster are more complex for the other two scenarios. In particular, to reflect the difficulty of making such assortment forecasts, we assume that the OSForecaster provides not a point forecast, but a distribution over possible assortments. We describe below how we use the test data to construct these distributions.

For long-term forecasts, we provide the retailer with the total number of purchasing customers at each store over the 26 weeks (≈ 6 months) comprising the entire test horizon, and the retailer must predict the sales for each of the products. To make these predictions, the retailer must first forecast an assortment to use for each store. We assume that when queried with a store s , the OSForecaster returns a uniform distribution over 20 assortments, each of which is constructed as follows: we first construct a random collection \mathcal{A} of assortments by including assortment a for each $a \in \mathcal{T}$ with probability 0.5 and then obtain one assortment by taking the union of all sets in the collection \mathcal{A} . Note that every assortment in the OSForecaster is obtained from sampling a new random collection \mathcal{A} of assortments. This construction is designed to reflect the possibility that the retailer is able to forecast some of the stock out and replenishment events, but not all.

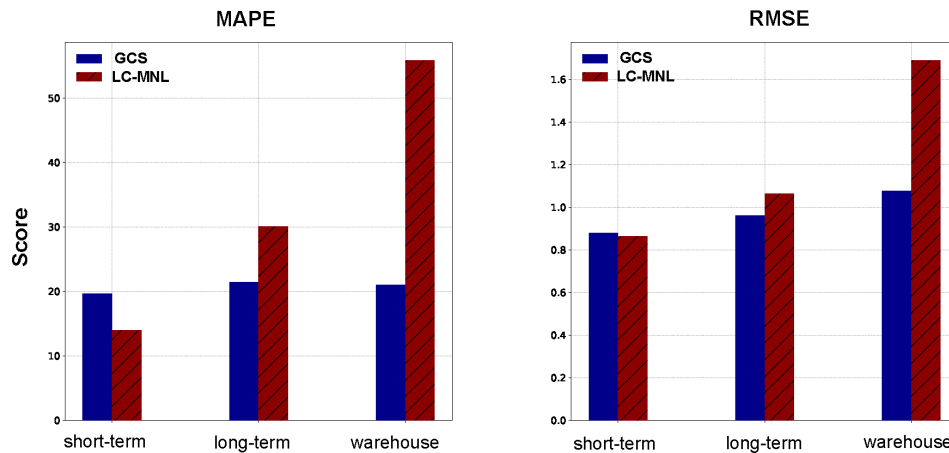
The procedure we use for warehouse forecasts is similar. The retailer must generate sales forecasts for each period t , across all the stores in the corpus, for which the retailer must forecast an assortment for each period t . We assume that when queried with a period t , the OSForecaster returns a uniform distribution over 20 assortments, each of which is constructed as follows: we construct a random collection \mathcal{A} of assortments by including assortment a for each store s with probability 0.5, and then obtain one assortment by taking the union of all sets in the collection \mathcal{A} . This construction is designed to reflect the possibility that the retailer has information on the assortments at some of the stores but not all.

6.3.2. Performance assessment. We evaluate all three scenarios in terms of the accuracy of predicting the total sales of each product across the entire test horizon, as described in Section 5.2. When the OSForecaster outputs a distribution, our predictions will be averaged over all the assortments in the distribution. In Figure 4 we present scatter plots of the improvements of the GCS model versus the LC-MNL model across 20 product categories under the three forecast scenarios discussed above. In the left and right panels we measure the predictive performance of the models under the MAPE and RMSE metrics, respectively, as defined in (15) (see Appendix A4.3 for alternative definitions of these metrics and corresponding results). We observe that GCS outperforms LC-MNL for around half of product categories for short-term forecasts, and for almost all product categories under the second and third forecast types. Note that we have dots located to the right of pluses with crosses being in between, under both MAPE and RMSE scores and across most of the product categories. It reveals that the improvement of GCS over LC-MNL across product



categories increases when we switch from short- to long-term, and from long-term to warehouse forecasts.

Figure 5 exhibits MAPE (left panel) and RMSE (right panel) scores of the GCS and LC-MNL models, averaging across 20 product categories, for the three different scenarios. We observe that the performance of the LC-MNL model deteriorates once we shift from short to long-term, and from long-term to warehouse forecasts. On the other hand, the predictive performance of the GCS model only moderately decreases once we switch to the noisy scenarios, i.e., the performances stays rather flat for all three scenarios. From the panels in Figure 5 we observe that the improvements of GCS over LC-MNL are -5.6% (-0.027%), 5.3% (0.081%), and 33.5% (0.56%) under the first, second, and third scenarios, respectively, based on the MAPE (RMSE) score.



The above observations show that the GCS model does not offer much in terms of performance gains when the test order sets can be forecasted accurately. It is only when one makes significant errors in forecasting the test order sets that we see a deterioration in the performance of the LC-MNL model, while the performance of the GCS model remains robust to noise.

Turning back to Figure 4, we notice that the improvement of GCS over LC-MNL varies across product categories for a given scenario. To better explain this variation, we regress the improvement of GCS over LC-MNL for each category against the noise intensity, which captures how much the forecasted order set differs from the true order set. We only consider the scenarios of long-term and warehouse forecasts for this analysis because the noise intensity for short-term forecasts is zero by definition. We define noise intensity at the transaction level and aggregate the metric across all the transactions. For each transaction, i in the test set, let O_i denote the true order set and \hat{O}_i denote the forecasted order set. For the transaction occurring at store s in period t , the forecasted order set \hat{O}_i takes the value $OSForecaster(\cdot)$ and $OSForecaster(\cdot)$ for long-term and warehouse predictions, respectively. We then use the following natural definition for noise intensity:

$$\text{noise intensity} = \frac{1}{\# \text{ of test transactions}} \sum \mathbb{E} \left[\frac{|O_i \setminus \hat{O}_i|}{|\hat{O}_i|} \right]$$

where the expectation is with respect to the distribution of the order set predictions \hat{O}_i . Intuitively, for long-term forecasts, the noise intensity captures how much the order set at a store varies over time; if there are very few stockouts, then the order set remains stable and the noise intensity is close to zero, but if there are many stockouts then the order set varies a lot and the noise intensity

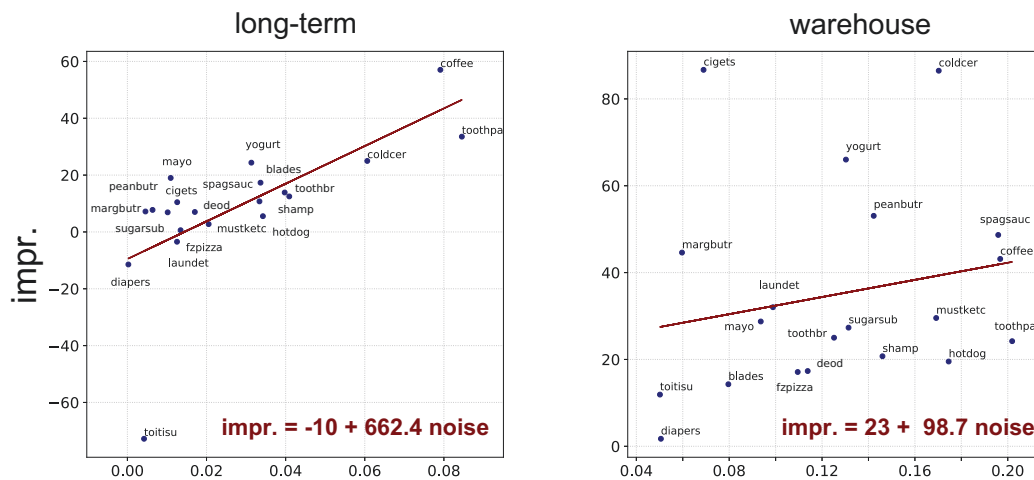


Figure 6 Scatter plots and linear regressions of the MAPE score improvement of GCS vs. LC-MNL over the noise intensity in the definition of offer sets across 20 product categories. Improvements are defined as the difference between two scores. In the left and right panels, we focus on the long-term and warehouse forecast scenarios, respectively.

is large. Similarly, for warehouse forecasts, the noise intensity captures how much the offer sets vary across stores. The left and right panels in Figure 6 illustrate the regression under the long-term and warehouse forecast scenarios, respectively. We see a clear positive correlation between the improvement of GCS over LC-MNL and noise intensity in both panels, suggesting the improvement becomes more significant with higher noise intensity in the product category.

In sum, our results all indicate that the CTC models have an advantage over the classic choice models when we expect the forecasted offer sets to be noisy. We make a few remarks. First, it is natural to wonder if our results are sensitive to the specific OSForecaster model we have used. To alleviate this concern, we repeat the above analysis with a ‘black-box’ noise model for the OSForecaster, where we generate the forecasted offer set \tilde{S}_{rt} for each store r and test period t by randomly adding products to the ‘true’ offer set S_{rt} , as done in the simulation study. The qualitative insights continue to hold; see Appendix A4.4. Our study also highlights the need to invest efforts into forecasting future offer sets more accurately; few, if any, such studies are available in the existing literature, indicating a natural direction for future work.

Second, we highlight the following peculiarity of our study: for long-term and warehouse forecast scenarios, the models are trained on noiseless offer sets but then tested on noisy offer sets. This is clearly in contrast to standard practice where efforts are made to train models on the same setups that they are tested on. Would not it then be better to ignore the fact that we have noiseless offer sets and instead train our models also on ‘forecasted’ offer sets? In Appendix A4.5, we repeat

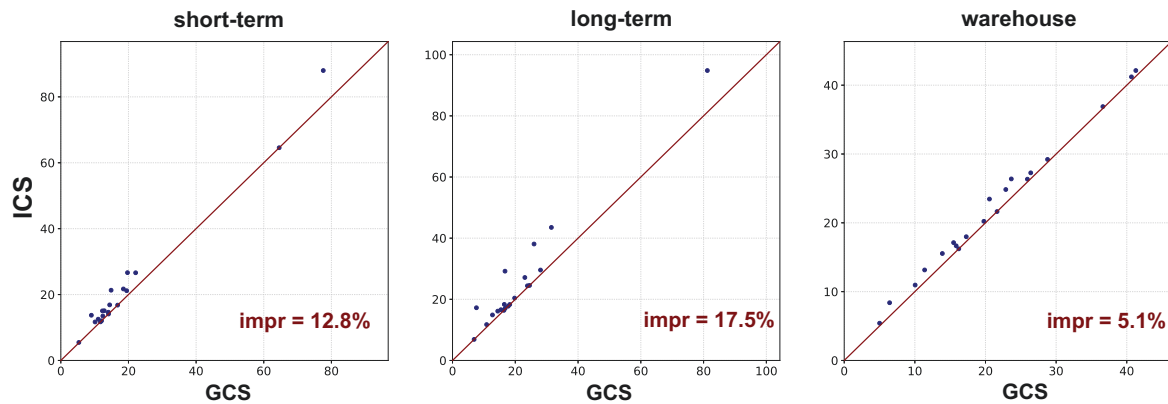


Figure 7 Scatter plots of the MAPE scores of 20 product categories under the ICS vs. GCS models. The left, middle, and right panels correspond to the cases of short-term, long-term, and warehouse forecast scenarios, respectively. Lower is better; therefore, GCS outperforms ICS for points above the 45° line.

our analysis by training all the models on offer sets obtained by applying the OSForecaster (as described above) to the training data². We observe that the GCS model’s performance increases slightly, but the improvement in the LC-MNL model’s performance is significant. However, despite the substantial improvement in the LC-MNL model’s performance, it is still unable to fully bridge the gap. We interpret this result as follows. *If* the underlying noise process is known, then it makes sense to incorporate it in the training process. But in practice, the noise process cannot be fully modeled because the factors affecting offer set noise (e.g., stock out and replenishment processes) vary over time, in which case the CTC models offer some protection against noise.

We conclude this section by comparing the predictive performance of the GCS model against a single-class ICS model; see Figure 7. We observe that the GCS model outperforms the ICS model on the MAPE metric by 12.8%, 17.5%, and 5.1% under the short-term, long-term, and warehouse forecast scenarios, respectively, on average, across 20 product categories. This finding suggests that a mixture of product-form consideration set distributions captures customer heterogeneity better, even when preferences are governed by a single rank list.

7. Case study on the car sharing dataset

The issue of noise in future offer sets is particularly acute for online platforms because product availability is determined by the market in real-time and is often hard to predict. In this section, we apply our consider-then-choose framework to a dataset from an industry partner, which runs an online peer-to-peer car-sharing platform. Our main finding is that consider-then-choose frameworks significantly outperform classical RUM models for predicting demand in these business environments.

² We thank an anonymous referee for suggesting this study.

In the rest of the section, we first provide some background information on our industry partner. Then, we describe the data and present our modeling assumptions. We incorporate the product feature information into the models in order to gain insights about consideration set formation. Then, we calibrate different variations of consider-then-choose models and a competitive, classical RUM benchmark model from the platform data and compare their predictive performance.

Our industry partner is an online, peer-to-peer car-sharing service that enables drivers to rent cars from private car owners, and owners to rent out their cars. The company offers its users a smartphone application to match car owners with renters on-demand. Car owners can use the application to list their vehicles by posting a picture of the vehicle and providing its detailed characteristics. In addition, car owners set the availability of their cars, hourly or daily prices, and potential conditions for sharing them. Every listed car has a device installed into it so that the renters are able to locate and unlock cars through the same application. As a car renter, the user of the platform can easily search for cars nearby and book the available alternative by entering the license number and credit card information.

For the empirical analysis in this section, we use a historical dataset including a sample of the rentals completed in a major US city over a period of two years. Each observation in the dataset is a rental (i.e., a renter who booked the listed car from a particular location given the set of available alternatives on a specific day/time). Our dataset includes 26.8K rentals from around five hundred car providers. For each rental, we have access to several observable features, such as car owner ID, hourly rental price, car access (i.e., open or closed), car location hours (i.e., 24 hours or restricted), car location type (i.e., garage, street, surface lot, or valet), car brand (e.g., BMW, Tesla, MINI), car type (i.e., economy, standard, full size, SUV, trucks, luxury), car age, and some other various binary car features such as transmission, premium wheels, power seats, bluetooth/wireless, leather interior, sunroof/moonroof, premium sound, power windows, GPS navigation system, roof rack, tinted windows. In Appendix A5.1 we examine the extent to which various features specified above (e.g., hourly rental price) impact the consideration set structure of renters. A detailed summary of the data is provided in Table A5 in Appendix A5. We split the dataset into two parts: the first 80%, in-sample, rental observations, and the remaining 20%, out-of-sample transactions.

The dataset consists of the rental request observations such that for every transaction we know which car was reserved and we can infer the set of available cars, listed in the online platform at the time of the request, with their characteristics. The offer sets are approximately built by

aggregating all listed and available cars within 0.3 miles distance from the location of the car which was in fact rented, defining tuples of the form (c, S) , where c is the chosen car and S is the set of cars available at the reservation time t .

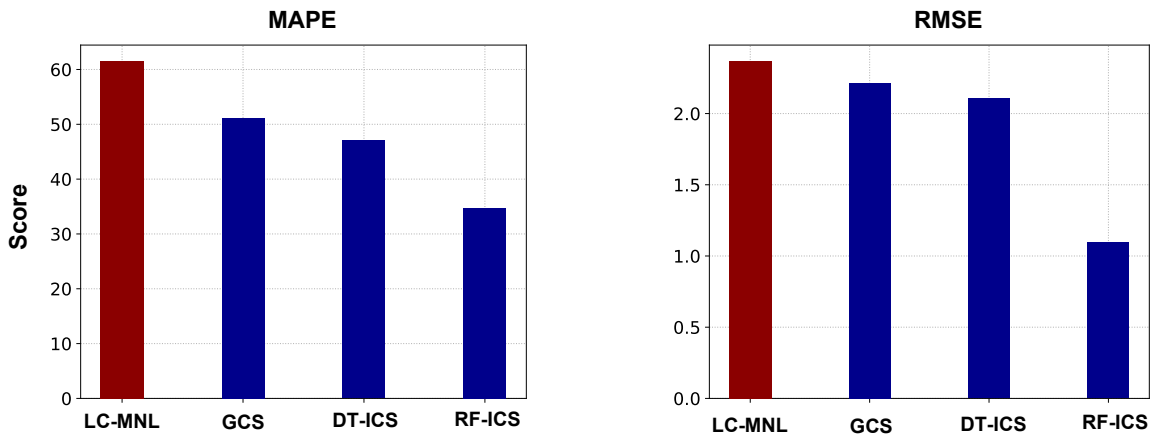
In general, in order to calibrate feature-based consider-then-choose models with our dataset, we need to estimate two types of parameters: the ranking π over all cars listed in the online platform (514 cars in total), and the parameters associated with the consideration set formation of renters. In order to simplify the estimation procedure, we assume that the ranking π is known a priori. Specifically, the cars are ranked according to their popularity among renters, defined as the number of times the vehicle was rented over the training dataset. Modeling the second stage choice process this way, we do not parameterize the ranking π which implies that the cars are assumed to have the same attributes over time, set at their average values. However, according to our dataset, this assumption is justified (see Appendix A5.2.1 for details).

Based on this popularity-based single ranking π , we estimate three variants of the CTC framework that account for product features with the purpose of characterizing the distribution over consideration sets: i) a GCS model defined as a mixture of logistic-based ICS (L-ICS) models, ii) a decision tree-based ICS (DT-ICS), and iii) a random forest-based ICS (RF-ICS). Details about the estimation of the L-ICS, DT-ICS, and RF-ICS models can be found in Appendix A3.1.

Our benchmark RUM model is the classic and competitive, feature-based LC-MNL model. For both GCS and LC-MNL, we tried up to $K = 5$ classes to find the optimal mixing distribution and report out-of-sample results for the number of mixtures that drove to best in-sample results.

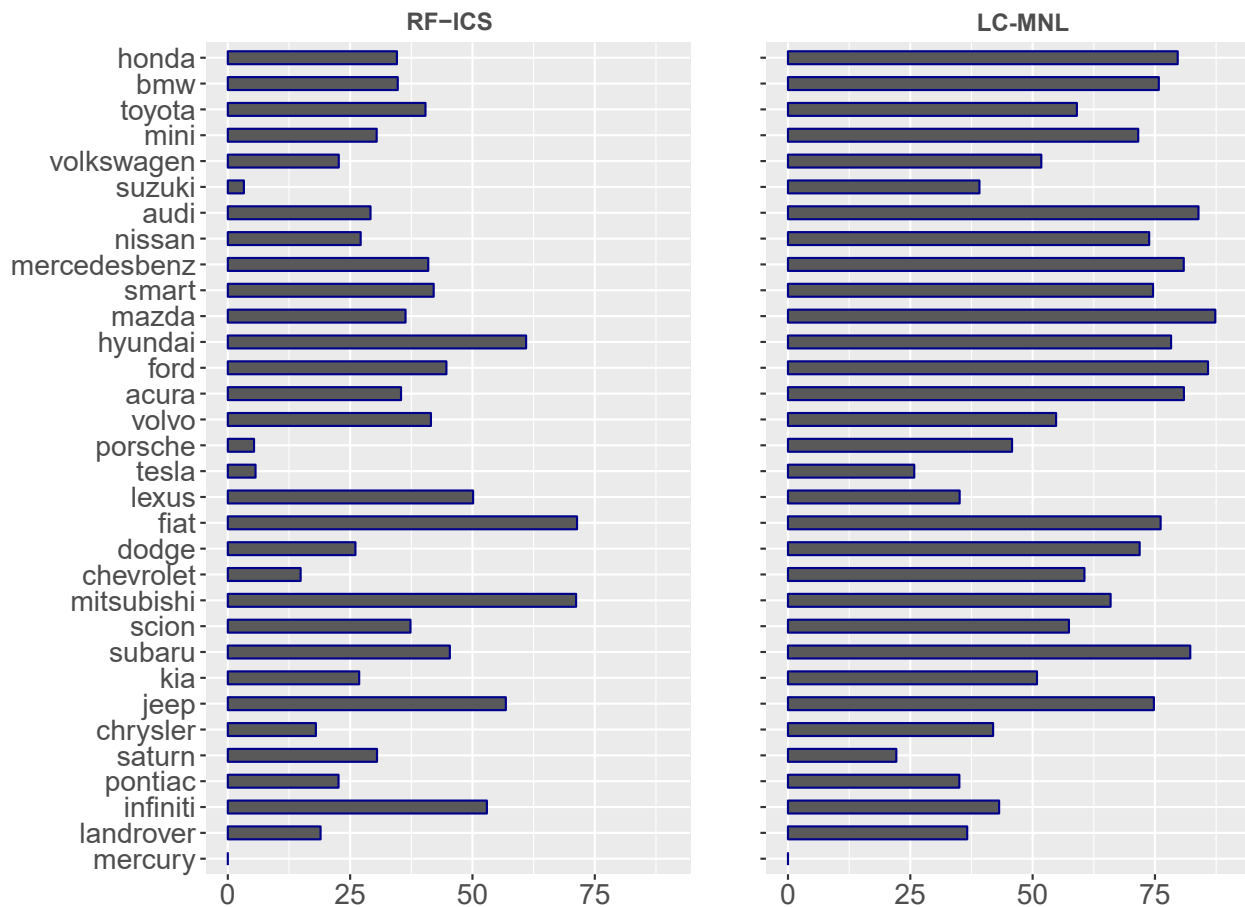
Next, we conduct an out-of-sample prediction testing of the analyzed models on the accuracy of two prediction measures: MAPE and RMSE (see Section 5.2), where lower scores stand for better prediction.

In order to optimize strategic and marketing decisions, the online platform needs to make long-term (or medium-term) demand forecasts for the cars listed on the online application. In real-world settings, the company can not rely on accurate data on car availabilities over time for the distant future, i.e., we can not test the prediction power of choice models by using the offer sets from the test dataset described above. Instead, the company might divide the city into several geographical areas and make predictions based on the aggregate assortment of cars listed in each area over the test time horizon. For our case study, we divide the city in 42 equal-spaced areas and estimate the assortments of cars by taking the superset of all the cars on offer over the entire horizon captured by the hold-out data for each area. Note that in this way we have 42 different offer sets (each corresponding to a particular area) while making predictions.



In Figure 8, we present the prediction performance results of the models based on MAPE (left panel) and RMSE (right panel) scores, averaged across all car brands. The MAPE score of consider-then-choose models exhibits an improvement of 16.7%, 23.4%, and 43.3% over LC-MNL for GCS, DT-ICS and RF-ICS, respectively. We also observe that consider-then-choose models obtain improvements of 6.2%, 10.9%, and 53.7% over LC-MNL for GCS, DT-ICS, and RF-ICS, respectively, based on RMSE metrics.

Figure 9 exhibits MAPE scores computed for every brand separately under the RF-ICS and LC-MNL models, where the brands are ordered according to their popularity (i.e., percentage of the total number of reservations in the training dataset coming from every brand), e.g., Honda is the most popular brand while Mercury is the least popular brand in the dataset. We note that the RF-ICS model outperforms the benchmark LC-MNL model almost consistently across all the brands. The panels also illustrate that MAPE scores vary significantly across brands both for RF-ICS and LC-MNL models. To further analyze this variation, in Figure 10 we regressed the improvement of RF-ICS over LC-MNL against the popularity of brands (left panel), and the improvement of RF-ICS over LC-MNL against the MAPE score of the LC-MNL model (right panel). We observe a clear positive correlation between MAPE score improvements and the popularity of brands, which indicates that we can better predict the demand for more popular brands. We can also see a clear positive correlation between the improvements and MAPE score under the LC-MNL model, allowing us to conclude that consider-then-choose type of models are especially relevant in prediction tasks, i.e., CTC models dominate LC-MNL, when the LC-MNL model provides a relatively poor prediction performance. Being robust to the noise, consider-then-choose models (and in particular, RF-ICS) provide significantly better predictive performance under these circumstances. Note that these insights are consistent with our numerical study based on the synthetic dataset in Section 5.



The results above indicate that consider-then-choose models forecast customer choices considerably better than the traditional LC-MNL model under both RMSE and MAPE scores. First of all, accounting for the consideration set formation with the linear-in-parameters GCS model with the logistically distributed error term, we can better predict the choices of customers. This improvement can be attributed to the effectiveness of consider-then-choose models to alleviate the noise impact on the outer set definition from sales transaction data. Moreover, we can further boost the predictive performance of the CTC models by modeling the consideration set formation in a nonlinear-in-parameters way, with decision trees or random forests. After calibrating DT-ICS and RF-ICS models we can get some insights of how customers form their consideration sets. In particular, Figure A17 in Appendix A5 illustrates an instance of the decision tree obtained after fitting the DT-ICS model.

In this paper, we analyze the importance of modeling customer choices by accounting for unobserved consideration sets. Even though consider-then-choose (CTC) models are gaining popularity in the

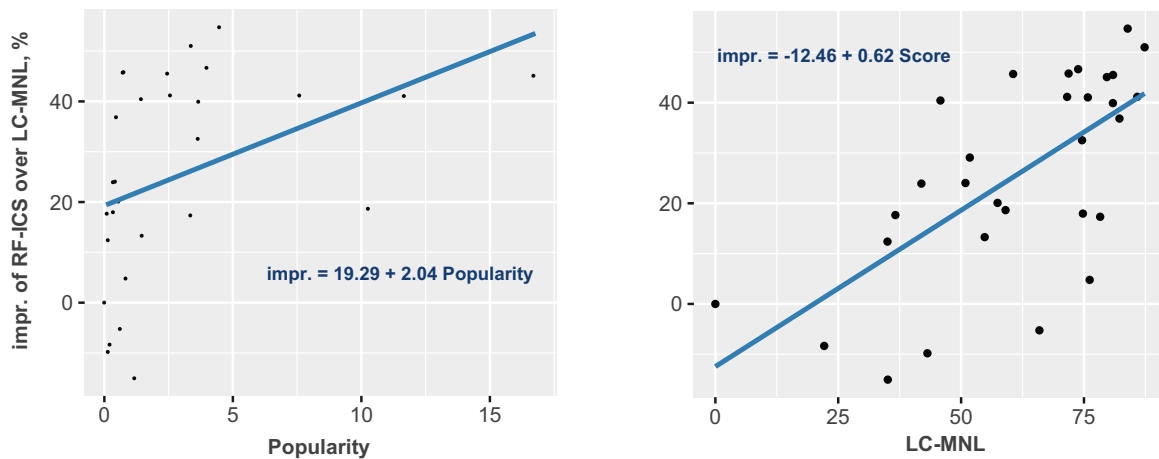


Figure 10 Left panel: scatter plot and linear regression of the percentage improvement of RF-ICS versus LC-MNL over brand popularity (i.e., percentage of total number of reservations in the training dataset), across 32 car brands. Right panel: scatter plot and linear regression of the percentage improvement of RF-ICS versus LC-MNL over LC-MNL prediction accuracy, across 32 car brands. In both panels, the prediction accuracy is measured by the MAPE score aggregated for every brand.

OM field, it is not clear from the existing literature when CTC models can outperform the classical models in the prediction performance if companies collect only transaction data. We show how one can effectively estimate a general class of CTC models and also address the problem of identifying such models when relying only on sales transaction data. We check the performance of the proposed methodology under a synthetic data setting where we control for different levels of noise and asymmetries of noise between the training and the testing data. Next, we apply CTC models to two real-world contexts: a retail operation and a car-sharing platform.

Our empirical results suggest that the predictive performance of CTC models significantly outperforms state-of-the-art RUM-based benchmarks widely used in marketing, economics, and more recently, in OM literature when there is noise in the data describing the offer sets. Moreover, we show that the relative improvement of consider-then-choose models in predictive performance becomes even more significant when there exist asymmetries between the accuracy of the description of the offer sets used to train the model, and the accuracy of the description of the offer sets used in the hold-out sample to derive forecasts. These results make our methodology promising for researchers interested in choice-based demand estimation, particularly for cases where the offer sets are not fully observable or whose definition is hard to anticipate looking forward.

References

Alba J, Chattopadhyay A (1985) Effects of context and postcategory cues on recall of competing brands. *Journal of Marketing Research* 22(3):340–349.

- Alptekinoglu A, Semple JH (2016) The exponential choice model: A new alternative for assortment and price optimization. *Operations Research* 64(1):79–93.
- Aouad A, Farias VF, Levi R (2020) Assortment optimization under consider-then-choose choice models. Forthcoming in *Management Science*.
- Blanchet J, Gallego G, Goyal V (2016) A markov chain approximation to choice modeling. *Operations Research* 64(4):886–905.
- Block H, Marschak J (1960) Random orderings and stochastic theories of responses. *Contributions to probability and statistics* 2:97–132.
- Bronnenberg BJ, Kruger MW, Mela CF (2008) Database paper the iri marketing data set. *Marketing science* 27(4):745–748.
- Campbell BM (1969) *The existence of evoked set and determinants of its magnitude in brand choice behavior*. Ph.D. thesis, Columbia University.
- Davis JM, Gallego G, Topaloglu H (2014) Assortment optimization under variants of the nested logit model. *Operations Research* 62(2):250–273.
- DeHoratius N, Raman A (2008) Inventory record inaccuracy: An empirical analysis. *Management Science* 54(4):627–641.
- Duran MA, Grossmann IE (1986) An outer-approximation algorithm for a class of mixed-integer nonlinear programs. *Mathematical programming* 36(3):307–339.
- Farias VF, Jagabathula S, Shah D (2013) A nonparametric approach to modeling choice with limited data. *Management Science* 59(2):305–322.
- Feldman J, Paul A, Topaloglu H (2019) Assortment optimization with small consideration sets. *Operations Research* 67(5):1283–1299.
- Feldman J, Topaloglu H (2015) Bounding optimal expected revenues for assortment optimization under mixtures of multinomial logits. *Production and Operations Management* 24(10):1598–1620.
- Feldman J, Topaloglu H (2018) Capacitated assortment optimization under the multinomial logit model with nested consideration sets. *Operations Research* 66(2):380–391.
- Gallego G, Li A (2017) Attention, consideration then selection choice model.
- Graham RL (1995) *Handbook of combinatorics* (Elsevier).
- Hauser JR (1978) Testing the accuracy, usefulness, and significance of probabilistic choice models: An information-theoretic approach. *Operations Research* 26(3):406–421.
- Hauser JR (2014) Consideration-set heuristics. *Journal of Business Research* 67(8):1688–1699.
- Hauser JR, Wernerfelt B (1990) An evaluation cost model of consideration sets. *Journal of consumer research* 16(4):393–408.

-
- Hausman JA (1996) Valuation of new goods under perfect and imperfect competition. *The economics of new goods*, 207–248 (University of Chicago Press).
- Hogarth RM, Karelaia N (2005) Simple models for multiattribute choice with many alternatives: When it does and does not pay to face trade-offs with binary attributes. *Management Science* 51(12):1860–1872.
- Howard JA, Sheth JN (1969) The theory of buyer behavior. *New York* 63:145.
- Hoyer WD (1984) An examination of consumer decision making for a common repeat purchase product. *Journal of consumer research* 11(3):822–829.
- Jagabathula S, Rusmevichientong P (2017) A nonparametric joint assortment and price choice model. *Management Science* 63(9):3128–3145.
- Jagabathula S, Subramanian L, Venkataraman A (2020) A conditional gradient approach for nonparametric estimation of mixing distributions. *Management Science* 66(8):3635–3656.
- Jagabathula S, Vulcano G (2018) A partial-order-based model to estimate individual preferences using panel data. *Management Science* 64(4):1609–1628.
- Kang Y, Gershwin SB (2005) Information inaccuracy in inventory systems: stock loss and stockout. *IIE transactions* 37(9):843–859.
- Lee E, Lee B (2012) Herding behavior in online p2p lending: An empirical investigation. *Electronic Commerce Research and Applications* 11(5):495–503.
- Lynch JG, Alba J, Hutchinson JW (1991) Memory and decision making. *Handbook of consumer behavior* 1–9.
- Manzini P, Mariotti M (2014) Stochastic choice and consideration sets. *Econometrica* 82(3):1153–1176.
- Montgomery H, Svenson O (1976) On decision rules and information processing strategies for choices among multiattribute alternatives. *Scandinavian Journal of Psychology* 17(1):283–291.
- Murphy KP (2012) *Machine learning: a probabilistic perspective* (MIT press).
- Newman JP, Ferguson ME, Garrow LA, Jacobs TL (2014) Estimation of choice-based models using sales data from a single firm. *Manufacturing & Service Operations Management* 16(2):184–197.
- Ratchford BT (1982) Cost-benefit models for explaining consumer choice and information seeking behavior. *Management Science* 28(2):197–212.
- Roberts JH, Lattin JM (1997) Consideration: Review of research and prospects for future insights. *Journal of Marketing Research* 406–410.
- Sher I, Fox JT, Bajari P, et al. (2011) Partial identification of heterogeneity in preference orderings over discrete choices. Technical report, National Bureau of Economic Research.
- Strauss D (1979) Some results on random utility models. *Journal of Mathematical Psychology* 20(1):35–52.
- Swait J, Ben-Akiva M (1987) Incorporating random constraints in discrete models of choice set generation. *Transportation Research Part B: Methodological* 21(2):91–102.

- Talluri K, Van Ryzin G (2004) Revenue management under a general discrete choice model of consumer behavior. *Management Science* 50(1):15–33.
- Tversky A (1972) Elimination by aspects: A theory of choice. *Psychological review* 79(4):281.
- van Ryzin G, Vulcano G (2017) An expectation-maximization method to estimate a rank-based choice model of demand. *Operations Research* 65(2):396–407.
- Vulcano G, Van Ryzin G, Ratliff R (2012) Estimating primary demand for substitutable products from sales transaction data. *Operations Research* 60(2):313–334.
- Wang R, Sahin O (2018) The impact of consumer search cost on assortment planning and pricing. *Management Science* 64(8):3649–3666.
- Westerlund T, Pettersson F (1995) An extended cutting plane method for solving convex minlp problems. *Computers & Chemical Engineering* 19:131–136.
- Wright P, Barbour F (1977) *Phased decision strategies: Sequels to an initial screening* (Graduate School of Business, Stanford University).

Demand estimation under uncertain consideration sets

APPENDIX

Srikanth Jagabathula

NYU Stern School of Business, New York, NY, sjagabat@stern.nyu.edu

Dmitry Mitrofanov

Carroll School of Management, Boston College, Chestnut Hill, MA, dmitry.mitrofanov@bc.edu

Gustavo Vulcano

School of Business, Universidad Torcuato Di Tella, and CONICET, Buenos Aires, Argentina,
gvulcano@utdt.edu

For completeness, we summarize the relevant notation from the main body and also introduce additional notation. We consider a universe \mathcal{I} of products $\{i_1, \dots, i_n\}$. We let \emptyset denote the no-purchase or the outside option. A customer is presented with a subset $S \subseteq \mathcal{I}$ of products and the customer chooses either one of the products in S or the outside option \emptyset . We let $\mathbb{P}(i)$ denote the probability that a customer chooses product $i \in \mathcal{I}$ and $\mathbb{P}(\emptyset)$ the probability that the customer chooses the outside option. We use \mathcal{S} to denote the set $\mathcal{S} = \mathcal{I} \cup \{\emptyset\}$. Let $\rho: \mathcal{S} \rightarrow [0, 1]$ define a distribution over consideration sets such that $\sum_{S \subseteq \mathcal{I}} \rho(S) = 1$. The preference relation \succ specifies a rank ordering over $n+1$ items which consist of the products in \mathcal{I} plus no-purchase option \emptyset with $\rho(i)$ denoting the preference rank of product i . The lower the rank of the product, the higher the preference, so that a customer's ranking ρ indicates that product i is preferred to product j if and only if $\rho(i) < \rho(j)$, or equivalently $i \succ j$. We assume that there is a distribution $\rho: \mathcal{S} \rightarrow [0, 1]$ over \mathcal{S} , which is the set of all full rankings or permutations of products in \mathcal{S} with cardinality $(n+1)!$.

To simplify the exposition, we also let $\mathcal{I} := \mathcal{I} \setminus \{\emptyset\}$, $\mathcal{S} := \mathcal{I} \cup \{\emptyset\}$, and $\mathbb{P}(i) = \Pr(i | \mathcal{S})$. Let $\langle \mathcal{I} \rangle$ denote the power set of \mathcal{I} , i.e., $|\langle \mathcal{I} \rangle| = 2^n$, and let \mathcal{U} denote $\{ \cup : i \in \mathcal{I} \}$ for any sets \mathcal{I} .

Proof of Proposition 1: First, we argue that if choice data are consistent with an underlying RUM model, then they are also consistent with a CTC model. To this end, let the distribution $\rho(\cdot)$ be a member of the RUM class. It defines a distribution over the $(n+1)!$ preference lists of products

in \mathcal{P} , which includes all the products in \mathcal{P} and the outside option. We map each ranking \succ in the support of (\cdot) to the tuple (\mathcal{P}, \succ) as follows: (a) \mathcal{P} consists of all the products in \mathcal{P} that are preferred over the outside option; that is, $\mathcal{P} = \{p \in \mathcal{P} : (p) \succ (o)\}$ and (b) \succ is the ranking obtained by moving the outside option to the last position (i.e., the $(|\mathcal{P}| + 1)$ th position) in ranking \succ . We then define the CTC model (\cdot, \cdot) such that $(\mathcal{P}, \succ) = (\mathcal{P}, \succ)$. It is straightforward to check that for any order set \mathcal{S} , the ranking \succ and the corresponding tuple (\mathcal{P}, \succ) result in the same choice. As a result, the choice probabilities under both (\cdot) and (\cdot, \cdot) should match.

We are now left to prove the other direction. Consider a CTC model (\cdot, \cdot) that defines a joint distribution over the preference lists of the products in \mathcal{P} that rank the outside option at the bottom, and the subsets of \mathcal{P} . We map each tuple (\mathcal{P}, \succ) in the support of (\cdot, \cdot) to the ranking \succ over the products in \mathcal{P} by repositioning the products not in \mathcal{P} to be below the outside option in the ranking \succ . More precisely, we have that $(\mathcal{P}, \succ) \succ (\mathcal{P}, \succ)$ if and only if $(\mathcal{P}, \succ) \succ (\mathcal{P}, \succ)$ whenever $p \in \mathcal{P}$ or $p \in \mathcal{P}$. In addition $(\mathcal{P}, \succ) \succ (\mathcal{P}, \succ)$ whenever $p \in \mathcal{P}$ and $(\mathcal{P}, \succ) \succ (\mathcal{P}, \succ)$ whenever $p \in \mathcal{P}$. We then define the RUM model (\cdot) such that $(\mathcal{P}, \succ) = (\mathcal{P}, \succ)$. Again, it is straightforward to check that both \succ and the tuple (\mathcal{P}, \succ) result in the same choice from each order set \mathcal{S} . It thus follows that the choice probabilities under both (\cdot, \cdot) and (\cdot) match.

The result of the proposition now follows. \square

We start this subsection by presenting Lemma A1 which shows how we can prove Proposition 3 by invoking a particular form of the inclusion-exclusion principle stated by Graham (1995). Then, we show how to prove Proposition 3 from first principles. Finally, we present Lemma A2 followed by the proof of Proposition 4 which relies on the proof of the combinatorial identity in Lemma A2.

LEMMA A1. For any sets $\mathcal{A} \subseteq \mathcal{B}$ and $\mathcal{C} \subseteq \mathcal{B}$, and the function $f: 2^{\mathcal{B}} \rightarrow \mathbb{R}$, we have

$$\sum_{\mathcal{D} \subseteq \mathcal{A}} \sum_{\mathcal{E} \subseteq \mathcal{C}} (-1)^{|\mathcal{D}| + |\mathcal{E}|} \cdot f(\mathcal{B} \setminus (\mathcal{D} \cup \mathcal{E})) = f(\mathcal{B} \setminus \mathcal{C}) \quad (\text{A1})$$

Proof: First, consider the inclusion-exclusion principle stated by Graham (1995) in the following form. Let \mathcal{B} be a finite set and $f: 2^{\mathcal{B}} \rightarrow \mathbb{R}$ be a real-valued function defined on the subsets of \mathcal{B} . Define the function $F: 2^{\mathcal{B}} \rightarrow \mathbb{R}$ by $F(\mathcal{A}) := \sum_{\mathcal{D} \subseteq \mathcal{A}} f(\mathcal{D})$, then $F(\mathcal{B} \setminus \mathcal{C}) := \sum_{\mathcal{D} \subseteq \mathcal{B} \setminus \mathcal{C}} (-1)^{|\mathcal{D}| + |\mathcal{C}|} f(\mathcal{D})$.

Then we show that the lemma follows from the stated above inclusion-exclusion principle. Let $F(\mathcal{A}) := \sum_{\mathcal{D} \subseteq \mathcal{A}} f(\mathcal{D})$, and $G(\mathcal{A}) := (-1)^{|\mathcal{A}|} \sum_{\mathcal{D} \subseteq \mathcal{A}} (-1)^{|\mathcal{D}|} \cdot f(\mathcal{D})$, which implies that

$$F(\mathcal{A}) \cdot (-1)^{|\mathcal{A}|} = \sum_{\mathcal{D} \subseteq \mathcal{A}} (-1)^{|\mathcal{D}|} \cdot f(\mathcal{D}) \quad \text{by invoking the inclusion-exclusion principle we obtain that}$$

$$(-1)^{-|\mathcal{A}|} \cdot f(\mathcal{A}) = \sum_{\mathcal{D} \subseteq \mathcal{A}} (-1)^{|\mathcal{D}| + |\mathcal{A}|} \cdot f(\mathcal{D}) \cdot (-1)^{|\mathcal{A}|} \quad \text{which implies that}$$

$$\begin{aligned} \binom{S \setminus T}{k} &= \binom{S}{k} = \sum_{\underline{C}} \binom{S}{k} = \sum_{\underline{C}} (-1)^{|\underline{C}|} \sum_{\underline{D}} (-1)^{|\underline{D}|} \cdot \binom{S}{k} = \sum_{\underline{C}} \sum_{\underline{D}} (-1)^{|\underline{C}| + |\underline{D}|} \cdot \binom{S}{k} \\ &= \sum_{\underline{C}} \sum_{\underline{D}} (-1)^{|\underline{C}| + |\underline{D}|} \cdot \binom{S \setminus T}{k} \end{aligned}$$

□

Proof of Proposition 3: For every $\underline{C} \subseteq \mathcal{C}$ we define boolean functions $f_{\underline{C}} : 2^{\mathcal{C}} \rightarrow \mathbb{R}$ and $g_{\underline{C}} : 2^{\mathcal{C}} \rightarrow \mathbb{R}$ by

$$\begin{aligned} f_{\underline{C}}(S) &= (-1)^{|\underline{C}|} \cdot \mathbf{I}[\underline{C} \subseteq S] \\ g_{\underline{C}}(S) &= (-1)^{|\underline{C}|} \cdot \mathbf{I}[\underline{C} \subseteq S] \end{aligned}$$

where $\mathbf{I}[\cdot]$ is an indicator function that is equal to 1, if condition \cdot is satisfied, and 0 otherwise. Then for all $\underline{C} \subseteq \mathcal{C}$ we claim that

$$\sum_{\underline{C}} f_{\underline{C}}(S) \cdot g_{\underline{C}}(S) = \begin{cases} 1 & \text{if } S = \mathcal{C} \\ 0 & \text{otherwise} \end{cases} \tag{A2}$$

First, we show that $\sum_{\underline{C}} f_{\underline{C}}(S) \cdot g_{\underline{C}}(S) = 1$ for every $\underline{C} \subseteq \mathcal{C}$:

$$\sum_{\underline{C}} f_{\underline{C}}(S) \cdot g_{\underline{C}}(S) = \sum_{\underline{C}} \mathbf{I}[\underline{C} \subseteq S] \cdot (-1)^{|\underline{C}|} \cdot \mathbf{I}[\underline{C} \subseteq S] = (-1)^{|\underline{C}|} \cdot \mathbf{I}[\underline{C} \subseteq S] = 1$$

Then we show that $\sum_{\underline{C}} f_{\underline{C}}(S) \cdot g_{\underline{C}}(S) = 0$ for all $\underline{C} \subseteq \mathcal{C}$ s.t. $\underline{C} \neq \mathcal{C}$:

$$\begin{aligned} \sum_{\underline{C}} f_{\underline{C}}(S) \cdot g_{\underline{C}}(S) &= \sum_{\underline{C}} \mathbf{I}[\underline{C} \subseteq S] \cdot (-1)^{|\underline{C}|} \cdot \mathbf{I}[\underline{C} \subseteq S] \\ &= (-1)^{|\underline{C}|} \cdot \sum_{\underline{C}} (-1)^{|\underline{C}|} \cdot \mathbf{I}[\underline{C} \subseteq S] \\ &= (-1)^{|\underline{C}|} \cdot (-1)^{|\underline{C}|} \cdot \sum_{\underline{C}} (-1)^{|\underline{C}|} \cdot \mathcal{C}^{|\underline{C}|} \cdot \mathbf{I}[\underline{C} \subseteq S] \text{ where } \mathcal{C} = \frac{1}{!(\mathcal{C} - \underline{C})!} \end{aligned}$$

[since the expression depends only on the cardinality of sets, the summation over the sets is reduced to the summation over the cardinality of sets]

$$= (-1)^{|\underline{C}|} \cdot (1 - 1)^{|\mathcal{C}| - |\underline{C}|} = 0$$

Consequently, the probability to choose the no purchase option from the offer set $\{S \setminus T\}$ is given by

$$\begin{aligned} \mathbb{P}(S \setminus T) &= \sum_{\underline{C}} \binom{S \setminus T}{k} = \sum_{\underline{C}} \binom{S}{k} \cdot (-1)^{|\underline{C}|} \cdot \mathbf{I}[\underline{C} \subseteq S] \\ &= \sum_{\underline{C}} \binom{S}{k} \cdot (-1)^{|\underline{C}|} \cdot \binom{S \setminus T}{k} \end{aligned} \tag{A3}$$

Then it follows that

$$\begin{aligned}
\sum_{\underline{c}} (-1)^{|\underline{c}|-1} \cdot \mathbb{P}(\underline{c} \setminus \underline{c}) &= \sum_{\underline{c}} \mathbb{P}(\underline{c} \setminus \underline{c}) \cdot (-1)^{|\underline{c}|-1} \mathbf{I}[\underline{c} \subseteq \underline{c}] \tag{A4} \\
&= (-1)^{|\underline{c}|-1} \cdot \sum_{\underline{c}} \mathbb{P}(\underline{c} \setminus \underline{c}) \cdot \mathbf{1}(\underline{c}) \\
&= (-1)^{|\underline{c}|-1} \cdot \sum_{\underline{c}} \sum_{\underline{1} \subseteq \underline{c}} \mathbf{1}(\underline{c}) \cdot (-1)^{|\underline{1}|-1} \cdot \mathbf{1}(\underline{1}) \cdot \mathbf{1}(\underline{c}) \\
&\quad \left[\text{by Equation (A3)} \right] \\
&= (-1)^{|\underline{c}|-1} \cdot \sum_{\underline{1} \subseteq \underline{c}} \mathbf{1}(\underline{c}) \cdot (-1)^{|\underline{1}|-1} \cdot \sum_{\underline{c}} \mathbf{1}(\underline{c}) \cdot \mathbf{1}(\underline{c}) \\
&= (-1)^{|\underline{c}|-1} \cdot \mathbf{1}(\underline{c}) \cdot (-1)^{|\underline{c}|-1} \left[\text{by Equation (A2)} \right] \\
&= \mathbf{1}(\underline{c})
\end{aligned}$$

Now it remains to prove the uniqueness of probability distribution function obtained from purchasing transactions data under the CTC choice model. Note that Equation (A3) relates probability distribution over consideration sets to the choice frequencies $\mathbb{P}(\underline{c} \setminus \underline{c})$ through the system of linear equations:

$$\mathbb{P}(\underline{c} \setminus \underline{c}) = \sum_{\underline{c}} \mathbf{1}(\underline{c}) \cdot (-1)^{|\underline{c}|-1} \cdot \mathbf{1}(\underline{c}) \quad \forall \underline{c} \iff \mathbf{c} = \mathbf{c} \tag{A5}$$

where $\mathbf{c} = (\mathbf{1}(\underline{c}))_{\underline{c}} \in \mathbb{R}^{2^{|I|} \times 1}$ denotes the $|2^{|I|}| \times 1$ vector of choice frequencies and $\mathbf{c} = (\mathbf{1}(\underline{c}))_{\underline{c}} \in \mathbb{R}^{2^{|I|} \times 1}$ denotes the $|2^{|I|}| \times 1$ vector that represents the probability distribution function over consideration sets. \mathbf{c} is the $|2^{|I|}| \times |2^{|I|}|$ matrix such that \mathbf{c} 's entry corresponding to the row \underline{c} and column \underline{c} is equal to $(-1)^{|\underline{c}|-1} \cdot \mathbf{1}(\underline{c})$. Therefore, the relation between the choice frequencies and the underlying model can be represented in a compact form as $\mathbf{c} = \mathbf{c} \cdot \mathbf{c}$. Then the proof of the uniqueness of \mathbf{c} reduces to showing that $\det(\mathbf{c}) \neq 0$. From Equation (A4) we have

$$\mathbf{c} = (-1)^{|\underline{c}|-1} \cdot \sum_{\underline{c}} \Pr(\underline{c} \setminus \underline{c}) \cdot \mathbf{1}(\underline{c}) \quad \forall \underline{c} \iff \mathbf{c} = \mathbf{c}$$

which establishes alternative linear relationship between choice frequencies $\Pr(\underline{c} \setminus \underline{c})$ and the model parameters \mathbf{c} in a compact form as $\mathbf{c} = \mathbf{c} \cdot \mathbf{c}$, where \mathbf{c} is the $|2^{|I|}| \times |2^{|I|}|$ matrix such that \mathbf{c} 's entry corresponding to the row \underline{c} and column \underline{c} is equal to $(-1)^{|\underline{c}|-1} \cdot \mathbf{1}(\underline{c})$. Therefore, we get

$$\begin{aligned}
\mathbf{c} &= \mathbf{c} \cdot \mathbf{c} \quad \left[\text{by Equation (A5)} \right] \\
\implies \mathbf{c} &= \mathbf{c} \cdot \mathbf{c} \implies \det(\mathbf{c}) = \det(\mathbf{c}) \cdot \det(\mathbf{c}) \\
\implies 1 &= \det(\mathbf{c}) \cdot \det(\mathbf{c}) \implies \det(\mathbf{c}) \neq 0
\end{aligned}$$

□

The next result will be invoked in the proof of upcoming Proposition 4.

LEMMA A2. *The combinatorial identity below is valid*

$$-\sum c \cdot \left[\sum_{-} (-1) c^{-} \right] = \begin{cases} 1 & \text{if } = \\ 0 & \text{if } \end{cases} \tag{A6}$$

where when = .

Proof: Let us consider two cases:

Case 1: = . In this case by invoking the assumptions of the lemma.

$$\begin{aligned} -\sum c \cdot \left[\sum_{-} (-1) c^{-} \right] &= -\sum c \cdot \left[\sum_{-} (-1) c^{-} \right] \\ &= -\sum \sum_{-} (-1) \cdot \frac{!}{! !(- -)!} = 1 \end{aligned}$$

where the last equality is proved by induction on = - :

Base case: = 1.

$$\begin{aligned} -\sum \sum_{-} (-1) \cdot \frac{!}{! !(- -)!} &= -\sum (-1)^{-} c = (-1) \cdot \sum (-1) c \\ &= (-1) \cdot ((1-1) - (-1)) = 1 \end{aligned}$$

Induction hypothesis: = .

Induction step: = + 1.

$$\begin{aligned} -\sum \sum_{-} (-1) \cdot \frac{!}{! !(- -)!} & \text{ [since } = - - 1] \\ &= -\sum \sum_{-} (-1) \cdot \frac{!}{! !(- -)!} + \sum (-1) \cdot \frac{!}{!(-)!(-)!} \\ &= -\sum \sum_{-} (-1) \cdot \frac{!}{! !(- -)!} + \frac{!}{!(-)!} \cdot \sum (-1) \cdot \frac{!}{!(-)!} \\ &= -\sum \sum_{-} (-1) \cdot \frac{!}{! !(- -)!} \\ &= -\sum \sum_{-} (-1) \cdot \frac{!}{! !(- -)!} - \sum (-1)^{-} \cdot \frac{!}{(- -)! ! !} \\ &= 1 - \sum (-1)^{-} \cdot \frac{!}{(- -)! ! !} \text{ [by induction hypothesis, } = -] \\ &= 1 + (-1)^{-} \cdot \frac{!}{(-)!} \cdot \sum (-1) \cdot \frac{(-)!}{(- -)! !} \\ &= 1 \end{aligned}$$

Case 2: _____ . the last equality is proved by induction on $n = m - 1$:

Base case: $n = 1$. Then $m = m - 1 \geq 0$, so that $\min(m) = 0$. And we have that

$$\begin{aligned} & - \sum_{c=0}^m c \cdot \left[\sum_{c=0}^m (-1)^c c^c \right] = - \sum_{c=0}^m c \cdot \left[\sum_{c=0}^m (-1)^c c^c \right] \\ & = (-1)^0 \sum_{c=0}^m (-1)^c c^c = 0 \end{aligned}$$

Induction hypothesis: $n = m$.

Induction step: $n = m + 1$.

Condition 1: $m \geq 0$. Then $\min(m) = 0$ and $\min(m - m - 1) = 0$. We have that

$$\begin{aligned} & - \sum_{c=0}^m c \cdot \left[\sum_{c=0}^m (-1)^c c^c \right] = - \sum_{c=0}^m c \cdot \left[\sum_{c=0}^m (-1)^c c^c \right] \text{ [since } m = m - m - 1] \\ & = - \sum_{c=0}^m c \cdot \left[\sum_{c=0}^m (-1)^c c^c \right] - \sum_{c=0}^m c \cdot (-1)^{m-m} \cdot c^{m-m} \\ & = - \sum_{c=0}^m c \cdot (-1)^{m-m} \cdot c^{m-m} \text{ [by induction hypothesis, } m = m - 1] \\ & = (-1)^{m-m} \cdot \sum_{c=0}^m (-1)^c \cdot c^c \cdot c^{m-m} = (-1)^{m-m} \cdot \frac{1}{1!} \cdot \sum_{c=0}^m (-1)^c \cdot \frac{(m-c)!}{1!(m-c)!(m-c)!} \end{aligned}$$

Now it is sufficient to show that $\sum_{c=0}^m (-1)^c \cdot \frac{(m-c)!}{1!(m-c)!(m-c)!} = 0$. We prove it by induction on m . For $m = 0$, it follows that $\sum_{c=0}^0 (-1)^c \cdot \frac{(0-c)!}{1!(0-c)!(0-c)!} = \sum_{c=0}^0 (-1)^c \cdot \frac{(0-c)!}{1!(0-c)!(0-c)!} = - \sum_{c=0}^0 (-1)^c \cdot c^c = 0$. Assuming that the result holds for $m = m$, we prove it for $m = m + 1$:

$$\begin{aligned} & \sum_{c=0}^m (-1)^c \cdot \frac{(m-c)!}{1!(m-c)!(m-c-1)!} = \sum_{c=0}^m (-1)^c \cdot \frac{(m-c)!(m-c)}{1!(m-c)!(m-c)!} \\ & = (m-c) \cdot \sum_{c=0}^m (-1)^c \cdot \frac{(m-c)!}{1!(m-c)!(m-c)!} - \sum_{c=0}^m (-1)^c \cdot \frac{(m-c)!}{1!(m-c)!(m-c)!} \\ & = - \sum_{c=0}^m (-1)^c \cdot \frac{(m-c)!}{1!(m-c)!(m-c)!} \text{ [by induction hypothesis, } m = m] \\ & = - \sum_{c=0}^m (-1)^c \cdot \frac{(m-c)!}{1!(m-c)!(m-c)!} \\ & = - \sum_{c=0}^m (-1)^c \cdot \frac{(m-c)!}{(m-1)!(m-c)!(m-c)!} \\ & = \sum_{c=0}^m (-1)^c \cdot \frac{(m-1-c)!}{1!(m-1-c)!(m-1-c)!} \\ & = \sum_{c=0}^m (-1)^c \cdot \frac{((m-1)-c)!}{1!((m-1)-c)!((m-1)-c)!} \\ & = 0 \text{ [by induction hypothesis, } m = m] \end{aligned}$$

Condition 2: $n \leq m$. Then $\min(n, m) = n$ and $\min(n, m-1) = n-1$. We have that

$$\begin{aligned}
 & - \sum_{k=0}^n c_k \cdot \left[\sum_{j=0}^{n-k} (-1)^j c_{n-k-j} \right] = - \sum_{k=0}^n c_k \cdot \left[\sum_{j=0}^{n-k} (-1)^j c_{n-k-j} \right] + \sum_{k=0}^n (-1)^k c_k \\
 & = - \sum_{k=0}^n c_k \cdot \left[\sum_{j=0}^{n-k} (-1)^j c_{n-k-j} \right] \\
 & = - \sum_{k=0}^n c_k \cdot \left[\sum_{j=0}^{n-k} (-1)^j c_{n-k-j} \right] - \sum_{k=0}^n c_k \cdot (-1)^{n-k} \cdot c_{n-k} \\
 & = - \sum_{k=0}^n c_k \cdot (-1)^{n-k} \cdot c_{n-k} \quad [\text{by induction hypothesis, } n = m-1] \\
 & = (-1)^n \cdot \sum_{k=0}^n (-1)^k \cdot c_k \cdot c_{n-k} = (-1)^n \cdot \frac{1}{1!} \cdot \sum_{k=0}^n (-1)^k \cdot \frac{(n-k)!}{!(n-k)!(n-k)!}
 \end{aligned}$$

Now it is sufficient to prove that $\sum_{k=0}^n (-1)^k \cdot \frac{(n-k)!}{!(n-k)!(n-k)!} = 0$. We prove it by induction on n . For $n=0$, it follows that $\sum_{k=0}^0 (-1)^k \cdot \frac{(n-k)!}{!(n-k)!(n-k)!} = \sum_{k=0}^0 (-1)^k \cdot \frac{(n-k)!}{!(n-k)!(n-k)!} = \frac{1}{1!} \cdot \sum_{k=0}^0 (-1)^k \cdot c_k = 0$. Assuming that the result holds for $n = m$, we prove it for $n = m+1$:

$$\begin{aligned}
 & \sum_{k=0}^n (-1)^k \cdot \frac{(n-k)!}{!(n-k-1-k)!(n-k)!} = \sum_{k=0}^n (-1)^k \cdot \frac{(n-k)!(n-k)}{!(n-k)!(n-k)!} \\
 & = (n-k) \cdot \sum_{k=0}^n (-1)^k \cdot \frac{(n-k)!}{!(n-k)!(n-k)!} - \sum_{k=0}^n (-1)^k \cdot \frac{(n-k)!}{!(n-k)!(n-k)!} \\
 & = - \sum_{k=0}^n (-1)^k \cdot \frac{(n-k)!}{!(n-k)!(n-k)!} \quad [\text{by induction hypothesis, } n = m] \\
 & = - \sum_{k=0}^n (-1)^k \cdot \frac{(n-k)!}{!(n-k)!(n-k)!} \\
 & = - \sum_{k=0}^n (-1)^k \cdot \frac{(n-k)!}{(n-1)!(n-k)!(n-k)!} \\
 & = \sum_{k=0}^n (-1)^k \cdot \frac{(n-1-k)!}{!(n-1-k)!(n-1-k)!} \\
 & = \sum_{k=0}^n (-1)^k \cdot \frac{((n-1)-k)!}{!((n-1)-k)!(n-1-k)!} \\
 & = 0 \quad [\text{by induction hypothesis}]
 \end{aligned}$$

□

Proof of Proposition 4: It follows from the proposition that

$$\begin{aligned} \mathbb{P}(\cdot) &= \sum_{\subseteq} \sum_{\supseteq \cup} (-1)^{|\cup| - |\cdot|} \cdot \mathbb{P}(\cdot) \cdot \mathbf{I}[|\cup| \leq |\cdot|] \\ &= \sum_{\subseteq} \sum_{\supseteq \cup} \mathbb{P}(\cdot) \cdot (-1)^{|\cap|} \cdot (-1)^{|\cup| - |\cdot|} \cdot \mathbf{I}[|\cup| \leq |\cdot|] \\ &= \sum_{\subseteq} \mathbb{P}(\cdot) \cdot (-1)^{|\cap|} \cdot \mathbf{I}[|\cup| \leq |\cdot|] \cdot \sum_{\supseteq \cup} (-1)^{|\cup| - |\cdot|} \cdot \mathbf{I}[|\cup| \leq |\cdot|] \end{aligned}$$

[since the expression depends only on the cardinality of sets \cup , the summation over the sets is reduced to the summation over the cardinality of sets \cup]

$$= \sum_{\subseteq} \mathbb{P}(\cdot) \cdot (-1)^{|\cap|} \cdot \mathbf{I}[|\cup| \leq |\cdot|] \cdot \sum_{-|\cup|}^{-|\cup|} (-1)^{\mathcal{C}^{-|\cup|}} \left[\text{where } \mathcal{C} = \frac{!}{!(\cdot -)!} \right]$$

For every \subseteq we define boolean functions $\mathbf{I}[\cdot] : 2 \rightarrow \mathbb{R}$ and $\mathcal{C}^{-|\cup|} : 2 \rightarrow \mathbb{R}$ by

$$\begin{aligned} \mathbf{I}[\subseteq] &= \mathbf{I}[|\subseteq| \leq |\cdot|] \\ \mathcal{C}^{-|\cup|} &= (-1)^{|\cap|} \cdot \mathbf{I}[|\cup| \leq |\cdot|] \cdot \sum_{-|\cup|}^{-|\cup|} (-1)^{\mathcal{C}^{-|\cup|}} \end{aligned}$$

Restricting consideration sets and order sets by the size of up to n (by the assumption of the proposition), we represent the probability to choose the no purchase option from the order set through a linear combination of boolean functions $\mathbf{I}[\cdot]$ as follows:

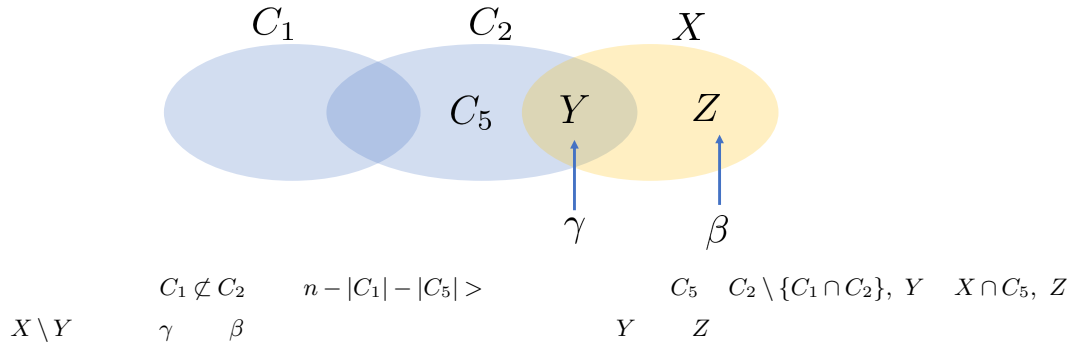
$$\mathbb{P}(\cdot) = \sum_{\subseteq} \mathbf{I}[\subseteq] \cdot \mathbf{I}[|\subseteq| \leq |\cdot|] = \sum_{\subseteq} \mathbf{I}[\subseteq] \cdot \mathcal{C}^{-|\cup|} \quad (\text{A7})$$

Then for all \subseteq such that $|\cup| \leq n$ we claim that

$$\sum_{\subseteq} \mathbf{I}[\subseteq] \cdot \mathcal{C}^{-|\cup|} = \begin{cases} 1 & \text{if } \cdot = \cup \\ 0 & \text{otherwise} \end{cases} \quad (\text{A8})$$

Consequently, it follows from the claim that

$$\begin{aligned} & \sum_{\subseteq} \mathbb{P}(\cdot) \cdot (-1)^{|\cap|} \cdot \mathbf{I}[|\cup| \leq |\cdot|] \cdot \sum_{-|\cup|}^{-|\cup|} (-1)^{\mathcal{C}^{-|\cup|}} \quad (\text{A9}) \\ &= \sum_{\subseteq} \mathbb{P}(\cdot) \cdot \mathbf{I}[\subseteq] = \sum_{\subseteq} \sum_{1 \subseteq} \mathbf{I}[\subseteq] \cdot \mathbf{I}[\subseteq] \cdot \mathbf{I}[\subseteq] \\ &= \sum_{1 \subseteq} \mathbf{I}[\subseteq] \cdot \sum_{\subseteq} \mathbf{I}[\subseteq] \cdot \mathbf{I}[\subseteq] = \mathbf{I}[\subseteq] \quad \left[\text{by Equation (A8)} \right] \end{aligned}$$



Now to complete the proof of the proposition, it is sufficient to prove the claim and show the uniqueness of the solution. We prove the claim by considering two different cases.

Case 1: $C_1 \not\subset C_2$.

$$\sum_{C \subseteq X} (-1)^{|C|} \cdot \mathbf{I}[|C| \leq n - |C_1| - |C_5|] \cdot \mathbf{I}[C \cap C_1 = \emptyset] \cdot \mathbf{I}[|C| \leq |C_2| - |C_5|]$$

$$\times \sum_{C \subseteq X} (-1)^{|C|} \mathcal{C}^{-|C|} \cdot \mathbf{I}[|C| \leq |C_2| - |C_5|]$$

[in this case, $C_1 \cap C_2 = \emptyset, |C_1 \cap C_2| = 0, |C_1 \cap C_5| = 0$, and $|C_1 \cup C_2| = |C_1| + |C_2|$]

$$= - \sum_{C \subseteq X} \mathbf{I}[|C| \leq n - |C_1| - |C_5|] \cdot \mathbf{I}[C \cap C_1 = \emptyset] \cdot \mathbf{I}[|C| \leq |C_2| - |C_5|] \cdot \sum_{C \subseteq X} (-1)^{|C|} \mathcal{C}^{-|C|} \cdot \mathbf{I}[|C| \leq |C_2| - |C_5|]$$

[since the expression depends only on the cardinality of sets, the summation over the sets is reduced to the summation over the cardinality of sets]

$$= - \sum_{k=0}^{|C_2| - |C_5|} \mathcal{C}^{-k} \cdot \left[\sum_{k=0}^{n - |C_1| - |C_5|} (-1)^k \mathcal{C}^{-k} \right] \left[\text{where } k = \text{cardinality of set } C \right]$$

$$= \begin{cases} 1 & \text{if } |C_2| - |C_5| = n - |C_1| - |C_5| \\ 0 & \text{if } |C_2| - |C_5| < n - |C_1| - |C_5| \end{cases}$$

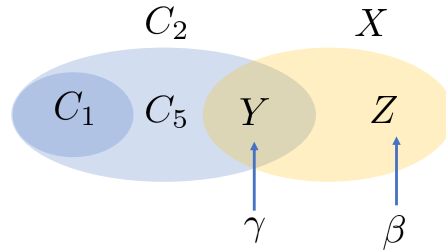
where the last equality follows by invoking Lemma A2, where $a = n - |C_1| - |C_5|$, $b = |C_2| - |C_5|$, and $c = |C_2| - |C_5|$.

Case 2: $C_1 \subset C_2$.

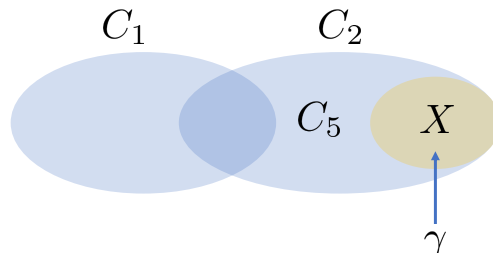
$$\sum_{C \subseteq X} (-1)^{|C|} \cdot \mathbf{I}[|C| \leq n - |C_1| - |C_5|] \cdot \mathbf{I}[C \cap C_1 = \emptyset] \cdot \mathbf{I}[|C| \leq |C_2| - |C_5|]$$

$$\times \sum_{C \subseteq X} (-1)^{|C|} \mathcal{C}^{-|C|} \cdot \mathbf{I}[|C| \leq |C_2| - |C_5|]$$

[since the expression depends only on the cardinality of sets, the summation over the sets is reduced to the summation over the cardinality of sets]



$$\begin{aligned}
 & C_1 \subset C_2 \quad k < n \quad n - |C_1| - |C_5| > \\
 & C_5 \subset C_2 \setminus \{C_1 \cap C_2\}, Y = X \cap C_5, Z = X \setminus Y \quad \gamma \quad \beta \\
 & Y \quad Z
 \end{aligned}$$



$$\begin{aligned}
 & C_1 \not\subset C_2 \quad n - |C_1| - |C_5| > \\
 & C_5 \subset C_2 \setminus \{C_1 \cap C_2\}, X \subseteq C_5 \quad \gamma \\
 & X
 \end{aligned}$$

$$= \begin{cases} \sum^{l_5} (-1)^{l_5} \cdot c^{l_5} \cdot \left[\sum_{l_2} (-1)^{l_2} (-1)^{l_1 - l_5} c^{-l_1 - l_5} \cdot \left[\sum_{l_1} (-1)^{l_1} (-1)^{l_2} c^{-l_1 - l_2} \right] \right] & \text{if } \not\subset \text{ and } -l_1 - l_2 - l_5 = 0 \text{ see Figure A1} \\ \sum^{l_5} (-1)^{l_5} \cdot c^{l_5} \cdot \left[\sum_{l_2} (-1)^{l_2} (-1)^{l_2} c^{-l_2} \cdot \left[\sum_{l_1} (-1)^{l_1} (-1)^{l_2} c^{-l_1 - l_2} \right] \right] & \text{if } \subset \text{ see Figure A2} \\ \sum^{l_5} (-1)^{l_5} \cdot c^{l_5} \cdot \left[\sum_{l_1} (-1)^{l_1} (-1)^{l_2} c^{-l_1 - l_2} \right] & \text{if } \not\subset \text{ and } -l_1 - l_2 - l_5 = 0 \text{ see Figure A3} \end{cases}$$

$$\begin{aligned}
 & \left[\text{where } = \setminus \{ \cap \} \right] \\
 & = \begin{cases} \left[-\sum^{l_5} (-1)^{l_5} \cdot c^{l_5} \right] \cdot \left[\sum_{l_2} (-1)^{l_2} (-1)^{l_1 - l_5} c^{-l_1 - l_5} \cdot \left[\sum_{l_1} (-1)^{l_1} (-1)^{l_2} c^{-l_1 - l_2} \right] \right] & \text{if } \not\subset \text{ and } -l_1 - l_2 - l_5 = 0 \\ \left[-\sum^{l_5} (-1)^{l_5} \cdot c^{l_5} \right] \cdot \left[\sum_{l_2} (-1)^{l_2} (-1)^{l_2} c^{-l_2} \cdot \left[\sum_{l_1} (-1)^{l_1} (-1)^{l_2} c^{-l_1 - l_2} \right] \right] & \text{if } \subset \\ \left[-\sum^{l_5} (-1)^{l_5} \cdot c^{l_5} \right] \cdot \left[\sum_{l_1} (-1)^{l_1} (-1)^{l_2} c^{-l_1 - l_2} \right] & \text{if } \not\subset \text{ and } -l_1 - l_2 - l_5 = 0 \end{cases} \\
 & = 0
 \end{aligned}$$

where the last equality follows since $\sum_{i \in \mathcal{C}} (-1)^{|\mathcal{C}| - |i|} \cdot \mathcal{C}^{|i|} = 0$.

In order to complete the proof, we show the uniqueness of the probability distribution function in our setting. First, note that Equation (A7) relates probability distribution $\mathbb{P}(\cdot)$ over consideration sets to the choice frequencies $\mathbb{P}(\cdot | \cdot)$ through the system of linear equations:

$$\mathbb{P}(\cdot) = \sum_{\subseteq} \mathbb{P}(\cdot | \cdot) \cdot \mathbb{1}(\cdot) \quad \forall \subseteq \iff \mathbb{P}(\cdot) = \cdot \tag{A10}$$

where $\mathbb{1}(\cdot)_{\subseteq}$ denotes the $|\mathcal{C}| \times 1$ vector of choice fractions and $\mathbb{P}(\cdot)_{\subseteq}$ denotes the $|\mathcal{C}| \times 1$ vector that represents the probability distribution function over consideration sets. $\mathbb{1}$ is the $|\mathcal{C}| \times |\mathcal{C}|$ matrix such that $\mathbb{1}_{ij}$ entry corresponding to the row i and column j is equal to $\mathbb{1}(i \subseteq j)$. As a result, the relation between the choice frequencies and the underlying model can be represented in a compact form as $\mathbb{P}(\cdot) = \mathbb{1} \cdot \mathbb{P}(\cdot | \cdot)$. Then the proof of the uniqueness of $\mathbb{P}(\cdot)$ reduces to showing that $\det(\mathbb{1}) \neq 0$. It follows from Equation (A9) that

$$\mathbb{P}(\cdot | \cdot) = \sum_{\subseteq} \mathbb{P}(\cdot) \cdot \mathbb{1}(\cdot) \quad \forall \subseteq \iff \mathbb{P}(\cdot | \cdot) = \cdot$$

which provides another relationship between choice frequencies $\mathbb{P}(\cdot | \cdot)$ and the model parameters in a linear form as $\mathbb{P}(\cdot | \cdot) = \cdot$, where \cdot is the $|\mathcal{C}| \times |\mathcal{C}|$ matrix such that \cdot_{ij} entry corresponding to the row i and column j is equal to $\mathbb{1}(i \subseteq j)$. Therefore, we get

$$\begin{aligned} \mathbb{P}(\cdot | \cdot) &= \cdot \cdot \quad \left[\text{by Equation (A10)} \right] \\ \implies \mathbb{P}(\cdot | \cdot) &= \cdot \implies \det(\cdot) = \det(\cdot) \cdot \det(\cdot) \\ \implies 1 &= \det(\cdot) \cdot \det(\cdot) \implies \det(\cdot) \neq 0 \end{aligned}$$

□

We start this subsection with the proof of Proposition 6. We then prove Lemmas A3, A4, and A5 followed by the proof of Proposition 7 which invokes these lemmas.

Proof of Proposition 6: It follows directly from the proof of Proposition 1 that $\mathcal{RUM} \subseteq \mathcal{GCS}$. Then, it remains to show that the RUM model class is not a specific case of the GCS model class. To this end, we provide a particular example of the RUM model class resulting in customers' choice frequencies that are inconsistent with the GCS choice rule.

Let $\mathcal{L} = \{ \cdot \}$, and recall the presence of the no purchase option \emptyset . Then let $\mathbb{P} : \mathcal{L} \rightarrow [0, 1]$ denote a specification of RUM class such that customers sample either preference list $\mathbb{P}(\cdot) = \{ \cdot \}$ with probability $\alpha \in (0, 1)$ or preference list $\mathbb{P}(\cdot) = \{ \cdot \}$ with probability $1 - \alpha$.

Consequently, probability distribution function over preference lists results in the following choice frequencies:

$$\begin{aligned} \mathbb{P}(\{ \quad \}) = \mathbb{P}(\{ \quad \}) = 1 &\Rightarrow \quad \text{is preferred to} \quad (\text{by Proposition 5}), \\ \mathbb{P}(\{ \quad \}) = 1 - \mathbb{P}(\{ \quad \}) = 1 &\Rightarrow \quad \text{is preferred to} \quad (\text{by Proposition 5}). \end{aligned}$$

These choice frequencies are inconsistent with GCS model class, which only allows a unique preference order of products, i.e., according to GCS choice rule either product is preferred to product or product is preferred to product. \square

LEMMA A3. Assume that for all consideration sets \subseteq we have that

$$\sum_{\subseteq} (-1)^{|\cdot|} \mathbb{P}(\cdot \setminus \cdot) \geq 0$$

with strict inequality for consideration sets of the size up to three, i.e., if $|\cdot| \leq 3$, then for all consideration sets \subseteq s.t. \subseteq it follows that

$$\sum_{\subseteq} (-1)^{|\cdot|} \mathbb{P}(\cdot \setminus \cdot) \geq 0$$

with strict inequality for consideration sets of the size up to three, i.e., if $|\cdot| \leq 3$.

Proof: Suppose that \subseteq and \subseteq . Let denote \setminus . We can now establish the following chain of equalities:

$$\begin{aligned} \sum_{\subseteq} (-1)^{|\cdot|} \mathbb{P}(\cdot \setminus \cdot) &= \sum_{\subseteq} (-1)^{|\cdot|} \cdot \mathbb{P}(\{\cdot \setminus \cdot\} \setminus \cdot) \\ &= \sum_{\subseteq} (-1)^{|\cdot|} \cdot \mathbb{P}(\{\cdot \setminus \cdot\} \setminus \cdot) \\ &= \sum_{\subseteq} \sum_{\subseteq} \sum_{\subseteq} (-1)^{|\cdot|} \cdot (-1)^{|\cdot|} \cdot \mathbb{P}(\{\cdot \setminus \cdot\} \setminus \cdot) \\ &\quad \left[\text{by invoking Lemma A1 for every } \subseteq, \text{ where } = \setminus = \right. \\ &\quad \left. = \text{ and } (\cdot \setminus \cdot) = (-1)^{|\cdot|} \cdot \mathbb{P}(\{\cdot \setminus \cdot\} \setminus \cdot) \right] \\ &= \sum_{\subseteq} \sum_{\subseteq} \sum_{\subseteq} (-1)^{|\cdot|} \cdot (-1)^{|\cdot|} \cdot \mathbb{P}(\{\cdot \setminus \cdot\} \setminus \cdot) \\ &= \sum_{\subseteq} \sum_{\in \langle \cup \rangle} (-1)^{|\cdot|} \cdot (-1)^{|\cdot|} \cdot \mathbb{P}(\cdot \setminus \cdot) \left[\text{where } = \cup, \text{ since } \cap = \emptyset \right] \\ &= \sum_{\subseteq} \sum_{\in \langle \cup \rangle} (-1)^{|\cup|} \cdot (-1)^{|\cdot|} \cdot \mathbb{P}(\cdot \setminus \cdot) \left[\text{since } \cap = \emptyset \right] \end{aligned}$$

$$\begin{aligned}
 &= \sum_{\subseteq} \sum_{\subseteq'} (-1)^{|'|-| |} \cdot \mathbb{P} (\setminus) \left[\text{where } ' = \cup \right] \\
 &\geq 0 \text{ with strict inequality when } | | \leq 3 \left[\text{by assumptions of the Lemma,} \right. \\
 &\left. \text{since } \sum_{\subseteq'} (-1)^{|'|-| |} \cdot \mathbb{P} (\setminus) \geq 0 \text{ with strict inequality when } | ' \leq 3 \right]
 \end{aligned}$$

□

The following three lemmas will be used in the proof of upcoming Proposition 7.

LEMMA A4. *If a sample of sales transaction data satisfies Conditions 1, 2, and 3, then for all $\in \cap$ where \subseteq and \subseteq we have that $\mathbb{P} () \geq \mathbb{P} ()$.*

Proof: Prove the result by induction on the $n=| | - | |$. We consider $\in \cap$ and \subseteq . For the base case $= 0$ we have that $=$ and $\mathbb{P} () = \mathbb{P} ()$. Assume that the result holds for $=$, i.e., $=$ and $| | - | | =$. Then we prove it for $= + 1$. Let us suppose w.l.o.g. that $= \cup \{ \}$ and \notin . Next, assume, by contradiction, that $\mathbb{P} (\setminus \{ \}) < \mathbb{P} ()$. Consequently, by Condition 2 it follows that $\mathbb{P} (\{ \}) < \mathbb{P} (\{ \})$. Then by Condition 1 we have that

$$\mathbb{P} (\{ \}) = \mathbb{P} (\{ \}) \tag{A11}$$

It now follows that

$$\begin{aligned}
 &\mathbb{P} (\{ \}) - \mathbb{P} (\{ \}) \\
 &= \left(1 - \mathbb{P} (\{ \}) \right) - \left(1 - \mathbb{P} (\{ \}) - \mathbb{P} (\{ \}) \right) \left[\text{by standard probability property} \right] \\
 &= \left(1 - \mathbb{P} (\{ \}) \right) - \left(1 - \mathbb{P} (\{ \}) - \mathbb{P} (\{ \}) \right) \left[\text{by Equation (A11)} \right] \\
 &= \left(1 - \mathbb{P} (\{ \}) \right) - \left(\mathbb{P} (\{ \}) - \mathbb{P} (\{ \}) \right) \left[\text{by standard probability property} \right] \\
 &= 1 - \mathbb{P} (\{ \}) - \mathbb{P} (\{ \}) + \mathbb{P} (\{ \}) < 0 \left[\text{by Condition 3 and Lemma A3, when } = = \{ \} \right]
 \end{aligned}$$

which contradicts to $\mathbb{P} (\{ \}) < \mathbb{P} (\{ \})$. Then we have

$$\begin{aligned}
 \mathbb{P} () &\leq \mathbb{P} (\setminus \{ \}) = \mathbb{P} () \left[\text{note that } | | - | | = \right] \\
 &\leq \mathbb{P} () \left[\text{by induction hypothesis} \right]
 \end{aligned}$$

Therefore, the result now follows by induction. □

LEMMA A5. *Consider $\in \subseteq \neq$. Then GCS choice model, with strict preference list and distribution over consideration sets where $() > 0$ if $| | \leq 3$, implies the following list of implications:*

$$\begin{aligned}
) \mathbb{P}(\setminus\{ \}) &\mathbb{P}(\setminus\{ \}) \implies \succ \text{ and } \forall' \subseteq \in': \mathbb{P}(\setminus\{ \}) > \mathbb{P}(\setminus\{ \}), \\
) \mathbb{P}(\setminus\{ \}) &= \mathbb{P}(\setminus\{ \}) \implies \succ \text{ and } \forall' \subseteq \in': \mathbb{P}(\setminus\{ \}) = \mathbb{P}(\setminus\{ \}), \\
) \mathbb{P}(\setminus\{ \}) &\neq \mathbb{P}(\setminus\{ \}) \implies \mathbb{P}(\setminus\{ \}) = \mathbb{P}(\setminus\{ \}).
\end{aligned}$$

Proof: a) Suppose that $\mathbb{P}(\setminus\{ \}) > \mathbb{P}(\setminus\{ \})$. Assume, by contradiction, that \succ . Then it can be inferred from purchase probability definition under the GCS, see Equation (1), that $\mathbb{P}(\setminus\{ \}) = \mathbb{P}(\setminus\{ \})$, which leads to a contradiction. As a result, we have that \succ since preferences are strict and asymmetric. Then $\forall' \subseteq \in'$ we establish that

$$\begin{aligned}
\mathbb{P}(\setminus\{ \}) - \mathbb{P}(\setminus\{ \}) &\geq (\setminus\{ \}) \left[\text{by Equation (1)} \right] \\
&0 \left[\text{by Assumption that } (\setminus\{ \}) = 0 \text{ if } |\setminus\{ \}| \leq 3 \right]
\end{aligned}$$

b) Suppose that $\mathbb{P}(\setminus\{ \}) = \mathbb{P}(\setminus\{ \})$. Assume, by contradiction, that \succ . Then it follows that

$$\begin{aligned}
\mathbb{P}(\setminus\{ \}) - \mathbb{P}(\setminus\{ \}) &\geq (\setminus\{ \}) \left[\text{by Equation (1)} \right] \\
&0 \left[\text{by Assumption that } (\setminus\{ \}) = 0 \text{ if } |\setminus\{ \}| \leq 3 \right]
\end{aligned}$$

which contradicts to the assumption above. As a result, we have that \succ , since preferences are strict and asymmetric. Then by Equation (1) we have that $\forall' \subseteq \in': \mathbb{P}(\setminus\{ \}) = \mathbb{P}(\setminus\{ \})$.

c) Suppose that $\mathbb{P}(\setminus\{ \}) \neq \mathbb{P}(\setminus\{ \})$. Then it is straightforward to verify that $\mathbb{P}(\setminus\{ \}) > \mathbb{P}(\setminus\{ \})$, since the following inequality holds from the Lemma A4: $\mathbb{P}(\setminus\{ \}) \geq \mathbb{P}(\setminus\{ \})$. Consequently, invoking the implication from part), we have \succ , and by Equation (1) we obtain that $\mathbb{P}(\setminus\{ \}) = \mathbb{P}(\setminus\{ \})$ \square

Proof of Proposition 7: Necessity: if purchasing transactions data is consistent with the GCS choice model with strict preference list and distribution over consideration sets where $(\setminus\{ \}) = 0$ if $|\setminus\{ \}| \leq 3$, then we claim that three axioms Condition 1, Condition 2, and Condition 3 are satisfied. First, it follows from Proposition 3 that Condition 3 is satisfied. Then Condition 1 and Condition 2 are satisfied by Lemma A5.

Sufficiency: we claim that the choice rule that satisfies Condition 1, Condition 2, and Condition 3 is a GCS choice model with the strict preference list where no purchase option is the least preferred item, and probability distribution function over consideration sets such that $(\setminus\{ \}) = 0$ if $|\setminus\{ \}| \leq 3$.

Define a binary relation between products $\subseteq \neq$, where $= 1$ if $\mathbb{P}(\setminus\{ \}) > \mathbb{P}(\setminus\{ \})$ for some \subseteq s.t. \in (note, by Condition 2 it implies that $\mathbb{P}(\setminus\{ \}) > \mathbb{P}(\setminus\{ \})$)

for all \subseteq s.t. \in), and zero otherwise. We claim that is complete, asymmetric, and transitive binary relation.

First, we prove that this binary relation is complete, i.e., either $= 1$ or $= 1$. Suppose that $\mathbb{P}(\setminus\{ \}) \leq \mathbb{P}()$ for some \subseteq , i.e., $= 0$. Then it follows from the Lemma A4 that $\mathbb{P}(\setminus\{ \}) = \mathbb{P}()$. Moreover, by Condition 2 we have that $\mathbb{P}(\{ \}) = \mathbb{P}(\{ \})$. We can now establish the following chain of equalities:

$$\begin{aligned} & \mathbb{P}(\{ \}) - \mathbb{P}(\{ \}) \\ &= \left(1 - \mathbb{P}(\{ \})\right) - \left(1 - \mathbb{P}(\{ \}) - \mathbb{P}(\{ \})\right) \quad [\text{by standard probability property}] \\ &= \left(1 - \mathbb{P}(\{ \})\right) - \left(1 - \mathbb{P}(\{ \}) - \mathbb{P}(\{ \})\right) \quad \left[\text{by Condition 2, see above}\right] \\ &= \left(1 - \mathbb{P}(\{ \})\right) - \left(\mathbb{P}(\{ \}) - \mathbb{P}(\{ \})\right) \quad [\text{by standard probability property}] \\ &= 1 - \mathbb{P}(\{ \}) - \mathbb{P}(\{ \}) + \mathbb{P}(\{ \}) = 0 \quad \left[\text{by Condition 3 and Lemma A3, where } = = \{ \}\right] \end{aligned}$$

which concludes that $= 1$. Therefore, completeness of binary relation now follows.

Second, we establish that the defined binary relation is asymmetric, i.e., if $= 1$ then $= 0$. Suppose that $\mathbb{P}(\setminus\{ \}) > \mathbb{P}()$ for some \subseteq , i.e., $= 1$. Then by Condition 1 we have that $\mathbb{P}(\setminus\{ \}) = \mathbb{P}()$ (note, by Condition 2 we have that for all $' \subseteq$ s.t. \in ': $\mathbb{P}(\setminus\{ \}) = \mathbb{P}(')$), which further implies that $= 0$. As a result, asymmetry of binary relation now follows.

Third, we show the transitivity of binary relation, i.e., if $= 1$ and $= 1$ then $= 1$ for all \in . Assume by contradiction that binary relation is not transitive. To this end, there exist \in such that $= 1$, $= 1$, $= 0$ with the following list of implications:

$$\begin{aligned} & = 1 \Rightarrow \mathbb{P}(\setminus\{ \}) > \mathbb{P}() \quad [\text{for some } \subseteq] \\ & \Rightarrow \mathbb{P}(\{ \}) > \mathbb{P}(\{ \}) \quad \left[\text{by Condition 2}\right] \\ & \Rightarrow \mathbb{P}(\{ \}) = \mathbb{P}(\{ \}) \quad \left[\text{by Condition 1}\right] \end{aligned} \tag{A12}$$

$$\Rightarrow \mathbb{P}(\{ \}) = \mathbb{P}(\{ \}) \quad \left[\text{by Condition 2}\right] \tag{A13}$$

$$\begin{aligned} & = 1 \Rightarrow \mathbb{P}(\setminus\{ \}) > \mathbb{P}() \quad [\text{for some } \subseteq] \\ & \Rightarrow \mathbb{P}(\{ \}) > \mathbb{P}(\{ \}) \quad \left[\text{by Condition 2}\right] \\ & \Rightarrow \mathbb{P}(\{ \}) = \mathbb{P}(\{ \}) \quad \left[\text{by Condition 1}\right] \end{aligned} \tag{A14}$$

$$\Rightarrow \mathbb{P}(\{ \}) = \mathbb{P}(\{ \}) \quad \left[\text{by Condition 2}\right] \tag{A15}$$

$$= 0 \Rightarrow \mathbb{P}(\setminus\{ \}) \leq \mathbb{P}() \quad [\text{for some } \subseteq]$$

$$\Rightarrow \mathbb{P}(A \setminus \{a\}) = \mathbb{P}(A) \quad \left[\text{by Lemma A4} \right] \quad (\text{A16})$$

$$\Rightarrow \mathbb{P}(\{a\}) = \mathbb{P}(A) \quad \left[\text{by Condition 2} \right] \quad (\text{A17})$$

Using the property of the choice rule, i.e., $\forall C \subseteq \mathcal{C} : \sum_{r \in C} \mathbb{P}(r) = 1$, for other sets $C = \{a, b\}$, $C = \{a, c\}$, and $C = \{a, b, c\}$ we further establish the following list of implications:

$$\begin{aligned} \text{For } C = \{a, b\} : \mathbb{P}(a) + \mathbb{P}(b) + \mathbb{P}(c) &= 1 \\ \Rightarrow \mathbb{P}(\{a\}) + \mathbb{P}(b) + \mathbb{P}(c) &= 1 \quad \left[\text{by Equation (A13)} \right] \\ \Rightarrow \mathbb{P}(\{a\}) + \mathbb{P}(b) + \mathbb{P}(c) &= 1 \quad \left[\text{by Equation (A14)} \right] \\ \Rightarrow \mathbb{P}(b) &= \mathbb{P}(\{a\}) - \mathbb{P}(c) \quad \left[\text{by standard probability property} \right] \end{aligned} \quad (\text{A18})$$

$$\begin{aligned} \text{For } C = \{a, c\} : \mathbb{P}(a) + \mathbb{P}(b) + \mathbb{P}(c) &= 1 \\ \Rightarrow \mathbb{P}(\{a\}) + \mathbb{P}(b) + \mathbb{P}(c) &= 1 \quad \left[\text{by Equation (A15)} \right] \\ \Rightarrow \mathbb{P}(\{a\}) + \mathbb{P}(b) + \mathbb{P}(c) &= 1 \quad \left[\text{by Equation (A16)} \right] \\ \Rightarrow \mathbb{P}(c) &= \mathbb{P}(\{a\}) - \mathbb{P}(b) \quad \left[\text{by standard probability property} \right] \end{aligned} \quad (\text{A19})$$

$$\begin{aligned} \text{For } C = \{a, b, c\} : \mathbb{P}(a) + \mathbb{P}(b) + \mathbb{P}(c) &= 1 \\ \Rightarrow \mathbb{P}(\{a\}) + \mathbb{P}(b) + \mathbb{P}(c) &= 1 \quad \left[\text{by Equation (A17)} \right] \\ \Rightarrow \mathbb{P}(\{a\}) + \mathbb{P}(b) + \mathbb{P}(c) &= 1 \quad \left[\text{by Equation (A12)} \right] \\ \Rightarrow \mathbb{P}(c) &= \mathbb{P}(\{a\}) - \mathbb{P}(b) \quad \left[\text{by standard probability property} \right] \end{aligned} \quad (\text{A20})$$

$$\begin{aligned} \text{For } C = \{a, b, c, d\} : \mathbb{P}(a) + \mathbb{P}(b) + \mathbb{P}(c) + \mathbb{P}(d) &= 1 \\ \Rightarrow 0 &= \mathbb{P}(\emptyset) - \mathbb{P}(\{a\}) - \mathbb{P}(\{b\}) - \mathbb{P}(\{c\}) + \mathbb{P}(a) + \mathbb{P}(b) \\ &+ \mathbb{P}(c) - \mathbb{P}(d) \quad \left[\text{since } \mathbb{P}(\emptyset) = 1, \text{ and by Equations (A18)-(A20)} \right] \\ &= 0 \quad \left[\text{by Condition 3 and Lemma A3, where } \mathbb{P}(a) = \mathbb{P}(b) = \mathbb{P}(c) \right] \end{aligned}$$

which leads to a contradiction. Therefore, the preference relation is transitive. Since we proved that binary relation is complete, asymmetric, and transitive, it specifies strict preference list \succ over products in \mathcal{C} , s.t. \succ_i $\mathbb{P}(i) = 1$. In addition, it immediately follows from the axioms that

is the least preferred item in the product universe according to the preference list \succ , i.e., for all $\in \mathcal{P}$ we have that $\mathbb{P}(\emptyset) = 0$:

$$\mathbb{P}(\emptyset) - \mathbb{P}(\{ \}) = 0 \left[\text{by Condition 3 and Lemma A3, where } = \{ \} \right]$$

which implies that $\mathbb{P}(\emptyset) = 0$ by definition.

Next, we prove that

$$\mathbb{P}(S) = \mathbb{P}(S' \setminus \{ \}) - \mathbb{P}(S') \quad \forall S \in \mathcal{S} \text{ s.t. } S \subseteq \mathcal{P}$$

where S' is the set of products that consists of product $\in S$ and all the items in \mathcal{P} that are preferred to item $\in S$, i.e., $S' = \{ \in \mathcal{P} : \succ \} \cup \{ \}$. The argument is proved by induction on the cardinality of the order set S , i.e., $|S|$. For the base case, $|S| = 1$, we have $\mathbb{P}(\{ \}) = 1 - \mathbb{P}(\{ \}) = \mathbb{P}(\emptyset) - \mathbb{P}(\{ \})$. Suppose the result follows for $|S| \leq k$, then we prove it for $|S| = k + 1$. We consider two cases.

Case 1: product $\in S$ is not the least preferred item in S' . In other words there exists $\in S'$ s.t. $\succ \in S'$. Then by definition of the binary relation \succ we have that $\mathbb{P}(S' \setminus \{ \}) = \mathbb{P}(S')$, and the result now follows:

$$\begin{aligned} \mathbb{P}(S) &= \mathbb{P}(S' \setminus \{ \}) \quad [\text{by Condition 1}] \\ &= \mathbb{P}(S' \setminus \{ \}) - \mathbb{P}(S') \quad [\text{by induction hypothesis,} \\ &\quad \text{and note that } \in S' \text{ since } \succ \in S'] \end{aligned}$$

Case 2: product $\in S$ is the least preferred item in S' . Consider order set $S' = \{ \in S' : \succ \} \cup \{ \}$ such that w.l.o.g. $\succ \succ \succ \succ \succ$. Assuming $\in S'$, we can now establish the following chain of equalities:

$$\begin{aligned} \mathbb{P}(S) &= 1 - \mathbb{P}(S') - \sum_{\in S'} \mathbb{P}(S') \\ &= -\mathbb{P}(S') + \mathbb{P}(\emptyset) - \sum_{\in S'} \mathbb{P}(\{ \in S' \}) \\ &= -\mathbb{P}(S') + \mathbb{P}(\emptyset) - \sum_{\in S'} \mathbb{P}(\{ \in S' \}) \quad [\text{by Condition 1}] \\ &= -\mathbb{P}(S') + \mathbb{P}(\emptyset) - \sum_{\in S'} \left(\mathbb{P}(\{ \in S' \}) - \mathbb{P}(\{ \in S' \}) \right) \\ &\hspace{15em} [\text{by induction hypothesis}] \\ &= -\mathbb{P}(S') + \mathbb{P}(\{ \in S' \}) \\ &= \mathbb{P}(\{ \in S' \}) - \mathbb{P}(\{ \in S' \}) = \mathbb{P}(S' \setminus \{ \}) - \mathbb{P}(S') \end{aligned}$$

Let us denote two particular sets \mathcal{C} and \mathcal{C}' as follows: $\mathcal{C} = \mathcal{C} \setminus \{r\}$, $\mathcal{C}' = \mathcal{C} \setminus \{r\}$. We can now establish the following chain of equalities:

$$\begin{aligned}
 \mathbb{P}(\mathcal{C}) &= \mathbb{P}(\mathcal{C} \setminus \{r\}) - \mathbb{P}(\mathcal{C}') \\
 &= \mathbb{P}(\mathcal{C} \setminus \{r\}) + \left(\sum_{\mathcal{C} \subseteq \mathcal{C}} \sum_{\mathcal{C} \subseteq \mathcal{C}} (-1)^{|\mathcal{C}'| - |\mathcal{C}|} \cdot \mathbb{P}(\mathcal{C} \setminus \{r\}) - \mathbb{P}(\mathcal{C} \setminus \{r\}) \right) - \mathbb{P}(\mathcal{C}') \\
 &\quad \left[\text{by invoking Lemma A1, where } \mathcal{C} = \mathcal{C} = \mathcal{C} \text{ and } (\mathcal{C} \setminus \{r\}) = \mathbb{P}(\mathcal{C} \setminus \{r\}) \right] \\
 &= \sum_{\mathcal{C} \subseteq \mathcal{C}} \sum_{\mathcal{C} \subseteq \mathcal{C}} (-1)^{|\mathcal{C}'| - |\mathcal{C}|} \cdot \mathbb{P}(\mathcal{C} \setminus \{r\}) - \mathbb{P}(\mathcal{C}') \left[\text{since } \mathcal{C} \setminus \{r\} = \mathcal{C} \setminus \{r\} \right] \\
 &= \sum_{\mathcal{C} \subseteq \mathcal{C}} \sum_{\mathcal{C} \subseteq \mathcal{C}} (-1)^{|\mathcal{C}'| - |\mathcal{C}|} \cdot \mathbb{P}(\mathcal{C} \setminus \{r\}) - \left(\sum_{\mathcal{C}' \subseteq \mathcal{C}'} \sum_{\mathcal{C} \subseteq \mathcal{C}} (-1)^{|\mathcal{C}'| - |\mathcal{C}|} \cdot \mathbb{P}(\mathcal{C} \setminus \{r\}) \right. \\
 &\quad \left. - \mathbb{P}(\mathcal{C} \setminus \{r\}) \right) - \mathbb{P}(\mathcal{C}') \\
 &\quad \left[\text{by invoking Lemma A1, where } \mathcal{C} = \mathcal{C} = \mathcal{C}' = \mathcal{C} \text{ and } (\mathcal{C} \setminus \{r\}) = \mathbb{P}(\mathcal{C} \setminus \{r\}) \right] \\
 &= \sum_{\mathcal{C} \subseteq \mathcal{C}} \sum_{\mathcal{C} \subseteq \mathcal{C}} (-1)^{|\mathcal{C}'| - |\mathcal{C}|} \cdot \mathbb{P}(\mathcal{C} \setminus \{r\}) - \sum_{\mathcal{C}' \subseteq \mathcal{C}'} \sum_{\mathcal{C} \subseteq \mathcal{C}} (-1)^{|\mathcal{C}'| - |\mathcal{C}|} \cdot \mathbb{P}(\mathcal{C} \setminus \{r\}) \left[\text{since } \mathcal{C} \setminus \{r\} = \mathcal{C}' \right] \\
 &= \sum_{r \in \langle \mathcal{C}' \cup \{r\} \rangle} \sum_{\mathcal{C} \subseteq \mathcal{C}} (-1)^{|\mathcal{C}'| - |\mathcal{C}|} \cdot \mathbb{P}(\mathcal{C} \setminus \{r\}) - \sum_{\mathcal{C}' \subseteq \mathcal{C}'} \sum_{\mathcal{C} \subseteq \mathcal{C}} (-1)^{|\mathcal{C}'| - |\mathcal{C}|} \cdot \mathbb{P}(\mathcal{C} \setminus \{r\}) \left[\text{since } \mathcal{C} = \mathcal{C}' \cup \{r\} \right] \\
 &= \sum_{r \in \langle \mathcal{C}' \cup \{r\} \rangle} \sum_{\mathcal{C} \subseteq \mathcal{C}} (-1)^{|\mathcal{C}'| - |\mathcal{C}|} \cdot \mathbb{P}(\mathcal{C} \setminus \{r\}) \\
 &= \sum_{r \in \langle \mathcal{C}' \cup \{r\} \rangle} \mathbb{P}(\mathcal{C}) \text{ where } \mathbb{P}(\mathcal{C}) = \sum_{\mathcal{C} \subseteq \mathcal{C}} (-1)^{|\mathcal{C}'| - |\mathcal{C}|} \cdot \mathbb{P}(\mathcal{C} \setminus \{r\}) \\
 &= \sum_{\mathcal{C} \subseteq \mathcal{C}} \mathbb{P}(\mathcal{C}) \cdot \mathbf{I}[\mathcal{C} \in \mathcal{C}] \cdot \mathbf{I}[\mathcal{C} \in \langle \mathcal{C}' \cup \{r\} \rangle] \\
 &= \sum_{\mathcal{C} \subseteq \mathcal{C}} \mathbb{P}(\mathcal{C}) \cdot \mathbf{I}[\mathcal{C} \in \mathcal{C}] \cdot \mathbf{I}[\succ \forall r \in \mathcal{C} \setminus \{r\}] \\
 &= \sum_{\mathcal{C} \subseteq \mathcal{C}} \mathbb{P}(\mathcal{C}) \cdot \mathbf{I}[\mathcal{C} \in \mathcal{C} \cap \mathcal{C}] \cdot \mathbf{I}[\succ \forall r \in \mathcal{C} \setminus \{r\}] \left[\text{since we assume that } r \in \mathcal{C}, \right. \\
 &\quad \left. \text{otherwise the choice probability is 0} \right]
 \end{aligned}$$

which is exactly the equation to compute the probability to purchase $r \in \mathcal{C}$ for the order set $\mathcal{C} \subseteq \mathcal{C}$ under the GCS choice model. As a result, we also have $\mathbb{P}(\mathcal{C}) = \sum_{\mathcal{C} \subseteq \mathcal{C}} \mathbb{P}(\mathcal{C}) \cdot \mathbf{I}[\mathcal{C} \cap \mathcal{C} = \emptyset]$ because of the standard probability law, i.e., $\mathbb{P}(\mathcal{C}) = 1 - \sum_{r \in \mathcal{C}} \mathbb{P}(\mathcal{C})$. Note that the above chain of equations specifies probability distribution function over consideration sets. Moreover, it follows from Proposition 3 that \succ is defined uniquely. In order to complete the proof, we show that the preference relation \succ is also defined uniquely. Suppose, by contradiction, there is another strict

preference order \succ' such that $\succ' \neq \succ$ and $\mathbb{P}(\cdot)_{\succ'} = \mathbb{P}(\cdot)_{\succ}$. Therefore there exist items $i \in \mathcal{N}$ s.t. $i \succ j$ and $j \succ' i$. By definition of the GCS choice rule, we have

$$\begin{aligned} \mathbb{P}(\{i, j\})_{\succ} &= \sum_{\subseteq} \mathbf{I}[i \in S] \cdot \mathbb{P}(S) \\ \mathbb{P}(\{i, j\})_{\succ'} &= \sum_{\subseteq} \mathbf{I}[i \in S] \cdot \mathbf{I}[j \in S] \cdot \mathbb{P}(S) \end{aligned}$$

As a result, we can establish now the following chain of inequalities:

$$\mathbb{P}(\{i, j\})_{\succ} - \mathbb{P}(\{i, j\})_{\succ'} \geq \mathbb{P}(S) \cdot \mathbf{I}[j \notin S] \geq 0 \quad \left[\text{by Condition 3} \right]$$

which contradicts to $\mathbb{P}(\cdot)_{\succ'} = \mathbb{P}(\cdot)_{\succ}$. \square

We start this section by showing how we can capture the impact of product features on consideration set formation and by providing the MINLP formulation for the estimation of the single class, logistic-based ICS (L-ICS) model. Then we describe the outer-approximation algorithm which is used to calibrate different variants of consider-then-choose models followed by the empirical validation of this algorithm. We finish this section by describing the EM algorithm to calibrate the GCS and CTC models.

To capture the impact of product features on consideration set formation, we use the following three ways to make propensity parameter α_i a function of the product features, where x_i is the observed i th feature of product i :

- *Logistic-based ICS model (L-ICS)*. We assume that customers have linear-in-parameters utility from considering product $i \in \mathcal{N}$, given by

$$u_i = \alpha_i + \sum_{k=1}^K \beta_k x_{ik} + \epsilon_i$$

where ϵ_i is a random variable distributed as a standard logistics, i.e., $\epsilon_i \sim \text{Logistic}(1)$. Therefore, product i is considered by an individual if and only if the utility from paying attention on it is non-negative, i.e.,

$$\alpha_i + \sum_{k=1}^K \beta_k x_{ik} + \epsilon_i \geq 0$$

Then the propensity of product i is given by

$$= \Pr[\alpha_i + \sum_{k=1}^K \beta_k x_{ik} + \epsilon_i \geq 0] = \frac{\exp(\alpha_i + \sum_{k=1}^K \beta_k x_{ik})}{1 + \exp(\alpha_i + \sum_{k=1}^K \beta_k x_{ik})}$$

- *Decision tree-based ICS model (DT-ICS)*. Here it is assumed that individuals decide which items to consider based on a tree with leaves $\in \{1, 2, \dots\}$ to which we can associate a mean probability of whether the item is going to be considered or not (see [Murphy \(2012\)](#)). Then we can write the probability to consider the item in the following way:

$$= \Pr[\text{item} \in \text{set}] = \sum_{\text{leaf}} \mathbb{I}[\text{item} \in \text{leaf}] = \sum_{\text{leaf}} \mathbb{I}(\text{leaf})$$

where leaf is the i -th region, i.e., the i -th leaf; features encodes the choice of features to split on and their threshold values, on the path from the root to the i -th leaf; and $\mathbb{I}(\text{leaf})$ is equal to 1 if item belongs to the i -th leaf, and equal to 0 otherwise.

- *Random forest-based ICS model (RF-ICS)*. In this case, we assume that individuals first randomly sample a tree and then decide which items to consider based on the sampled tree (see [Murphy \(2012\)](#)). Note that random forest avoids the overfitting problem of decision trees by adding more trees instead of building one big tree. We can write the probability of considering the item as follows:

$$= \Pr[\text{item} \in \text{set}] = \sum_{\text{tree}} \frac{1}{N} \mathbb{I}(\text{tree})$$

where $\mathbb{I}(\text{tree})$ is the probability of considering item according to the i -th decision tree.

Estimation methodology. In a similar spirit to the preliminaries in Section 4.1, we can formulate the maximum likelihood estimation problem for the single class, logistic-based ICS (L-ICS) model with product features in such a way so that we can apply the outer-approximation algorithm in Appendix A3.2 in order to calibrate it.

On the other hand, the calibration of the DT-ICS and RF-ICS models is more challenging. To this end, we need to estimate both the ranking and the decision tree (or random forest). Intuitively, both decision trees and random forests are not limited to the generalized linear model class. Instead, they can model the relationship between the features in a non-linear way by bisecting the space into smaller and smaller regions. Note that if the ranking is known, then the log-likelihood optimization is equivalent to calibrating a classification decision tree or random forest with the splitting criteria based on the entropy function. Then, given a decision tree or a random forest, the log-likelihood optimization problem reduces to solving a MILP to find θ . We could then heuristically iterate between these two steps of finding θ and finding a decision tree until finding a fixed-point or a time limit is reached. As a practical matter in Section 7, we assume that the ranking is known a priori (see Section A5.2 for the details).

In this part of the section we formulate the maximum likelihood estimation problem for the logistic-based ICS (L-ICS), and then simplify it in such a way so that we can apply the outer-approximation algorithm in Section A3.3 in order to calibrate it. Recall that $\delta_{j_t k_t} \in \{0, 1\}$, $\forall j_t, k_t \in \mathcal{I}_t$, $j_t \neq k_t$ is binary linear ordering variable such that $\delta_{j_t k_t} = 1$ if product j_t goes before product k_t in the preference list \succ (or, equivalently, $\delta_{k_t j_t} = 0$), and $\delta_{j_t k_t} = 0$ otherwise. The data log-likelihood function under this model is given by

$$\mathcal{L}(\beta, \delta) = \sum_t \left[\log \frac{\beta \mathbf{X}_{j_t}}{1 + \beta \mathbf{X}_{j_t}} + \sum_{\substack{k \in \mathcal{I}_t \\ k \neq j_t}} \left[\delta_{j_t k_t} \log \frac{1}{1 + \beta \mathbf{X}_{k_t}} \right] \right]$$

and the ML problem can be represented in the following way

$$\begin{aligned} \max_{\beta, \delta} \quad & \mathcal{L}(\beta, \delta) \\ \text{s.t.} \quad & \delta_{j_t k_t} + \delta_{k_t j_t} = 1 \quad \forall j_t, k_t \in \mathcal{I}_t, j_t \neq k_t \\ & \delta_{j_t k_t} + \delta_{k_t l_t} \leq 2 \quad \forall j_t, k_t, l_t \in \mathcal{I}_t, j_t \neq k_t \neq l_t \\ & \delta_{j_t k_t} \in \{0, 1\} \quad \forall j_t, k_t \in \mathcal{I}_t, j_t \neq k_t \end{aligned} \tag{A21}$$

where the constraints ensure that $\delta_{j_t k_t}$ indeed represents a total order. In particular, the first set of constraints ensures that either j_t is preferred over k_t or vice versa, and the second set of constraints imposes the total ordering among any three products. To simplify the likelihood function, we introduce a new variable $\gamma_{j_t k_t}$ defined as $\gamma_{j_t k_t} = \delta_{j_t k_t} - \delta_{k_t j_t}$, $\forall j_t, k_t \in \mathcal{I}_t, j_t \neq k_t$ and rewrite the likelihood function in the following way

$$\mathcal{L}(\beta, \gamma) = \sum_t \left[\log \frac{\beta \mathbf{X}_{j_t}}{1 + \beta \mathbf{X}_{j_t}} + \sum_{\substack{k \in \mathcal{I}_t \\ k \neq j_t}} (\delta_{j_t k_t} - 1) \log \left(\frac{1}{2} \right) + \sum_{\substack{k \in \mathcal{I}_t \\ k \neq j_t}} \left[\log \frac{1}{1 + \sum_i \gamma_{i k_t} \mathbf{X}_{i k_t}} \right] \right]$$

since if $\delta_{j_t k_t} = 1$ we have that $\delta_{k_t j_t} = 0 \quad \forall j_t, k_t \in \mathcal{I}_t, j_t \neq k_t$, and

$$\delta_{j_t k_t} \log \frac{1}{1 + \beta \mathbf{X}_{k_t}} = \log \frac{1}{1 + \beta \mathbf{X}_{k_t}} = \log \frac{1}{1 + \sum_i \gamma_{i k_t} \mathbf{X}_{i k_t}} = (\delta_{j_t k_t} - 1) \log \left(\frac{1}{2} \right) + \log \frac{1}{1 + \sum_i \gamma_{i k_t} \mathbf{X}_{i k_t}}$$

if $\delta_{j_t k_t} = 0$ we have that $\delta_{k_t j_t} = 1 \quad \forall j_t, k_t \in \mathcal{I}_t, j_t \neq k_t$, and

$$\delta_{k_t j_t} \log \frac{1}{1 + \beta \mathbf{X}_{k_t}} = 0 = -\log \left(\frac{1}{2} \right) + \log \frac{1}{1 + \beta \mathbf{X}_{k_t}} = (\delta_{j_t k_t} - 1) \log \left(\frac{1}{2} \right) + \log \frac{1}{1 + \sum_i \gamma_{i k_t} \mathbf{X}_{i k_t}}$$

Let β_{j_t} be the value of the largest component in vector β , i.e., $\beta_{j_t} = \max_k \beta_{k_t}$. And we also define $\mathcal{L}(\beta, \gamma)$ in the following way

$$\mathcal{L}(\beta, \gamma) = \sum_t \left[\log \frac{\beta \mathbf{X}_{j_t}}{1 + \beta \mathbf{X}_{j_t}} + \sum_{\substack{k \in \mathcal{I}_t \\ k \neq j_t}} \left[\log \frac{1}{1 + \sum_i \gamma_{i k_t} \mathbf{X}_{i k_t}} \right] \right]$$

We can then formulate the MLE problem in terms of the variables ():

$$\begin{aligned}
 \beta \tau \delta \quad & \mathcal{L}(\beta, \tau, \delta) + \sum_{\substack{k \in \mathcal{K} \\ t}} (n_k - 1) \log\left(\frac{1}{2}\right) & \text{(A22)} \\
 \mathbf{s.t.}: \quad & \beta_k \leq \tau_k \quad \forall k \\
 & \tau_k \leq \delta_k \quad \forall k \\
 & \beta_k \geq \alpha_k + \tau_k - \delta_k \quad \forall k \\
 & \beta_k \geq -\tau_k \quad \forall k \\
 & \beta_k + \tau_k = 1 \quad \forall k \quad \beta_k \leq \delta_k \\
 & \beta_k + \tau_k \leq 2 \quad \forall k \quad \beta_k \neq \delta_k \\
 & \beta_k \in \{0, 1\} \quad \forall k
 \end{aligned}$$

where the first four sets of linear constraints ensure that $\beta_k = \delta_k$, $\forall k$, given that β_k is a binary variable.

The high level idea behind outer-approximation algorithm is to approximate the convex constraint in MINLP by the set of linear constraints. This way we can solve the MINLP by solving the sequence of MILPs. In particular, we start with a feasible solution of MINLP. Then, we linearize the convex constraint at the previously obtained feasible solution. The next step is to solve MILP with a linearized convex constraint and obtain an additional point to linearize the convex constraint and continue iteratively. Note that, at each iteration, we add only one additional constraint to the optimization problem, solved at previous iteration. More formally, suppose that the problem is to solve the MINLP (P) defined below.

$$\begin{aligned}
 \theta \tau \delta \quad & \min_{\beta, \tau, \delta} \mathcal{L}(\beta, \tau, \delta) & \text{(P)} \\
 \mathbf{s.t.}: \quad & \beta_k \leq \mathcal{L}(\beta, \tau, \delta) \\
 & \beta_k + \tau_k \leq 0 \\
 & \beta_k = \alpha_k + \tau_k \\
 & \beta_k \leq \delta_k \\
 & \beta_k \in \{0, 1\} \quad \forall k
 \end{aligned}$$

where $\mathcal{L}(\beta, \tau, \delta)$ is a concave function; β , τ , δ , and θ are continuous decision variables; θ is a binary variable; α , β , τ , and δ are constant vectors; and β_k and δ_k are lower and upper bounds of β_k ,

respectively. Note that optimization problems (10) and (A22) have a similar structure and can be represented as the above mentioned MINLP (P) without loss of generality.

Before we provide the details of the outer-approximation algorithm, let us define the linear constraint () that is added to the optimization problem at the k th iteration for given (), i.e., we linearize the convex constraint at ():

$$() = \{ : \mathcal{L}() + \frac{\mathcal{L}() - \mathcal{L}()}{\delta} + \frac{\mathcal{L}() - \mathcal{L}()}{\tau} - \geq 0 \quad \in \mathcal{R} \}$$

The broad idea of the outer-approximation algorithm is that at the k th iteration, we substitute convex constraint $\leq \mathcal{L}()$ with a set of linear constraints $= \{ () \}$.

Next we define two subproblems () and () of the optimization problem (P) that are used to describe the Algorithm 1, which exploits the outer-approximation technique. First, let us define the concave subproblem () for given () (i.e., if () is known) in the following way

$$\begin{aligned} & \theta \tau \delta \quad (()) \\ \mathbf{s.t.}: & \leq \mathcal{L}() \\ & + + \leq 0 \\ & = + \end{aligned}$$

Note that solving MINLP (P) reduces to solving concave subproblem above if () is known. Second, let us define the MILP subproblem () for given () and () as follows

$$\begin{aligned} & \theta \tau \delta \quad () \\ \mathbf{s.t.}: & () \in \\ & + + \leq 0 \\ & = + \\ & \leq \leq \\ & \in \{0, 1\} \quad \forall \end{aligned}$$

Note that solving MINLP (P) reduces to solving MILP subproblem () if we approximate convex constraint $\leq \mathcal{L}()$ with a set of linear constraints ().

Now we can formally apply the outer-approximation method (Duran and Grossmann, 1986) to solve the optimization problem (P), see Algorithm 1. The proposed algorithm effectively exploits the structure of the optimization problem (P) where we have a linearity of the binary variables and convexity of the non-linear constraint which only depends on continuous variables. In order to

linearize the optimization problem, we use the outer-approximation of a convex set by the intersection of its collection of supporting half-spaces. To this end, the outer approximation defines the optimization problem () as MILP. Because of the potentially many continuous points required for outer-approximation, we solve a sequence of MILPs to build up increasingly tight relaxation of the original MINLP. Overall, the proposed Algorithm 1 consists of solving a finite sequence of convex problems (()) and relaxed versions of a MILP (()).

Algorithm 1 Outer-Approximation algorithm for optimization problem (P)

procedure OUTER-APPROXIMATION(P)

$\mathcal{R} = \mathcal{R} \times \mathcal{R}$, $\mathcal{R} = -\infty$, $\mathcal{R} = \infty$, $\mathcal{R} = 1$

Select arbitrary \mathcal{R} , i.e., it can be arbitrary full ranking

while $|\mathcal{R} - \mathcal{R}| > \epsilon$ **do**

Solve concave subproblem () such that $\mathcal{R} = \mathcal{R}^*$ (i.e., the optimal objective function of ()), and $(\mathcal{R}^*) = (\mathcal{R}^* \ \mathcal{R}^*)$ (i.e., the optimal solution of ())

Set $\mathcal{R} = \mathcal{R}^- \cap (\mathcal{R}^*)$

Solve MILP subproblem () such that $\mathcal{R} = \mathcal{R}^*$ (i.e., the optimal objective function of ()), and $(\mathcal{R}^*) = (\mathcal{R}^* \ \mathcal{R}^*)$ (i.e., the optimal solution of ())

$\mathcal{R} = \mathcal{R} + 1$

return (()).

A3.3.1. Outer-approximation algorithm vs. cutting plane algorithm Note that Algorithm 1 to solve the optimization problem (P) requires the solution of both convex optimization problem (()) and MILP (()). The solution of the convex optimization problem (()) in each iteration might be extremely computationally demanding, while in solving the MILP (()) the computational work, on the other hand, might be more moderate, because for every iteration we need to solve the MILP problem (()) which is the previous MILP problem (()) with only one additional linear constraint added. Therefore, we propose to use the cutting plane algorithm to solve the MINLP in this case (Westerlund and Pettersson, 1995), which would require the solution of only the finite sequence of MILP problem (()), see Algorithm 2. Note that Algorithm 2 is identical to the Algorithm 1 except that we skip Step 5 in the cutting plane Algorithm 2. Even though the main iteration loop of Algorithm 1 is, generally, more efficient, we have global convergence for both Algorithms 1 and 2.

Algorithm 2 Cutting plane algorithm for optimization problem (P)

procedure CUTTING PLANE(P)

 $\mathcal{R} = \mathcal{R} \times \mathcal{R}$, $\alpha = -\infty$, $\beta = \infty$, $\gamma = 1$

 Select arbitrary \mathcal{S} , i.e., it can be arbitrary full ranking

 Select arbitrary \mathcal{D} , i.e., it can be arbitrary distribution over consideration sets

 Set $\mathcal{S} = \mathcal{S} \cdot \mathcal{D}$, $\alpha = -\infty$, $\beta = \infty$
while $|\mathcal{S} - \mathcal{S}^*| > \epsilon$ **do**

 Set $\mathcal{S} = \mathcal{S} \cap \mathcal{S}^*$

 Solve MILP subproblem $\mathcal{P}(\mathcal{S})$ such that $\alpha = \alpha^*$, $\beta = \beta^*$ (i.e., the optimal objective function of $\mathcal{P}(\mathcal{S})$), and $(\mathcal{S}^*) = (\mathcal{S}^*, \mathcal{S}^*, \mathcal{S}^*)$ (i.e., the optimal solution of $\mathcal{P}(\mathcal{S})$)

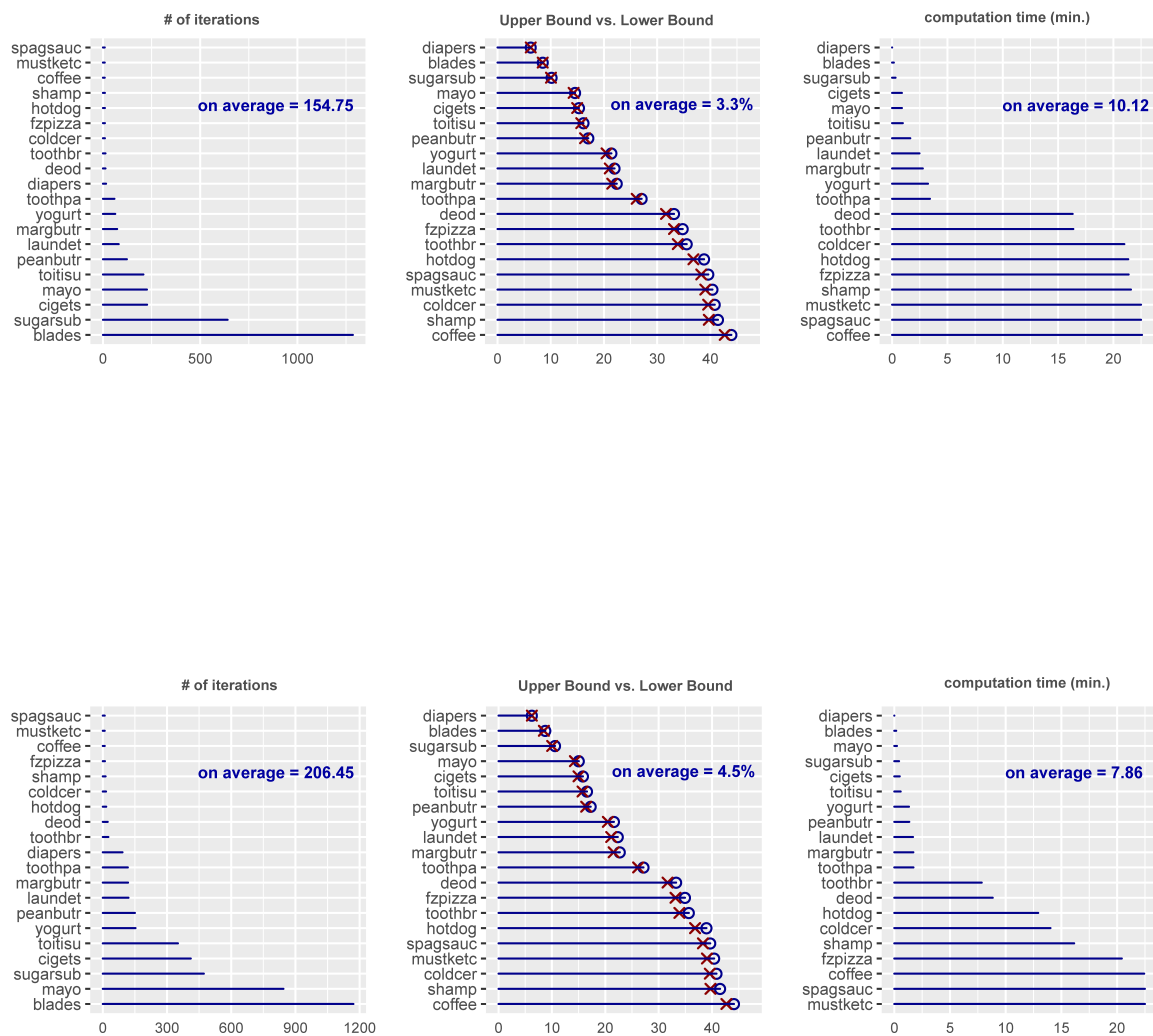
 $\gamma = \gamma + 1$
return (\mathcal{S}^*) .

A3.3.2. Empirical validation of the algorithms. In this section, we analyze the performance of outer-approximation algorithm (1) and cutting plane algorithm (2) to estimate ICS model with IRI Academic dataset (see data descriptive statistics in Table A1). We limited the running time of the algorithms by 3 hours, and the precision was set to 1e-6. It follows from Figure A4 that the optimality gap of the outer approximation algorithm (1) to calibrate the ICS model is 3.3% on average over 20 product categories. On the other hand, it is shown in Figure A5 that the optimality gap of the cutting plane algorithm (2) to calibrate the ICS model is 4.5% on average over 20 product categories. Following these findings, we apply outer-approximation algorithm (1) to calibrate the ICS model in our analysis as it provides significantly faster convergence to the optimal solution, which is consistent with previous studies.

In this section, we present the EM algorithm to calibrate the GCS model. We provide two versions of this algorithm that can be applied to the aggregate-level and individual-level sales transaction data.

A3.4.1. Estimation with aggregate level data. The log-likelihood function to calibrate the GCS model, after we reparametrize it by dividing all the transactions into σ segments, is given by

$$\log \mathcal{L}(\boldsymbol{\gamma}) = \sum_t \log \left(\sum_{j \in \sigma} \prod_{j \in \sigma} (1 - \gamma_j) \right) \quad (\text{A23})$$



where $w_i \geq 0$ is the weight of the class i a priori, s.t. $\sum_i w_i = 1$; \mathcal{O}_t denotes the set of ordered items at time t ; \mathbf{x}_t denotes the product purchased at time t ; T denotes the time horizon.

Unsurprisingly, the above likelihood function is nonconcave. In order to alleviate the complexity of solving the MLE problem directly, we use the Expectation Maximization (EM) algorithm. First, let us outline the main principles of EM procedure. We start with arbitrary initial parameter estimates $\mathbf{x}^{(0)}$. Then, we compute the conditional expected value of the log-likelihood function $E[\log \mathcal{L}(\mathbf{x}) | \mathbf{x}^{(0)}]$ (the E, expectation, step). Next, the resulting expected log-likelihood function is maximized to compute new estimates $\mathbf{x}^{(1)}$ (the M, maximization, step), and we repeat the

algorithm until convergence to get a sequence of estimates $\{\mathbf{x}^t = 1, 2, \dots\}$. We further describe the E-step and M-step of every iteration and how we start the algorithm in the context of our estimation problem.

Initialization: we initialize the EM with a random allocation of observations to one of the classes, resulting in an initial allocation $\mathcal{D} = \mathcal{D}_1 \cup \dots \cup \mathcal{D}_h$, which form a partition of the collection of all the transactions. Then we set $\theta_h = |\mathcal{D}_h| / (\sum_{j \in \mathcal{D}} |\mathcal{D}_j|)$. Then θ_h (i.e., θ_h) and θ_h for all $h \in \{1, \dots, h\}$, $\theta_h \in \mathcal{D}_h$ are obtained by solving the following optimization problem:

$$\theta_h \sum_{j \in \mathcal{D}_h} \left(\log \theta_h + \sum_{\substack{j \in \mathcal{D}_h \\ j \neq \theta_h}} \log(1 - \theta_h) \right)$$

which is solved by using the outer approximation algorithm in Section [A3.3](#).

E-step: we compute θ_h , which is the membership probability of every transaction at time t to belong to the segment \mathcal{D}_h (i.e., $\theta_h \in \mathcal{D}_h$, where \mathcal{D}_h is the set of transactions in class h) based on the parameter estimates $\{\theta_h, \theta_h, \theta_h\}$ and the purchasing transactions data $(\theta_h, \theta_h) | \theta_h$:

$$\begin{aligned} &= \Pr \left(\theta_h \in \mathcal{D}_h \mid \theta_h, \theta_h, \theta_h, (\theta_h, \theta_h) \right) \\ &= \Pr \left(\theta_h \in \mathcal{D}_h \mid \theta_h, \theta_h, \theta_h, (\theta_h, \theta_h) \right) \text{ [independence of purchases]} \end{aligned}$$

$$\begin{aligned}
&= \frac{\Pr\left(\left(\begin{matrix} t \\ \cdot \end{matrix}\right) \middle| \begin{matrix} - \\ - \\ - \\ \in \end{matrix}\right) \cdot \Pr\left(\begin{matrix} \in \\ \cdot \\ - \\ - \end{matrix} \middle| \begin{matrix} - \\ - \\ - \end{matrix}\right)}{\Pr\left(\left(\begin{matrix} t \\ \cdot \end{matrix}\right) \middle| \begin{matrix} - \\ - \\ - \end{matrix}\right)} \quad [\text{Bayes theorem}] \\
&= \frac{\Pr\left(\left(\begin{matrix} t \\ \cdot \end{matrix}\right) \middle| \begin{matrix} - \\ - \\ - \\ \in \end{matrix}\right) \cdot \Pr\left(\begin{matrix} \in \\ \cdot \\ - \\ - \end{matrix} \middle| \begin{matrix} - \\ - \\ - \end{matrix}\right)}{\sum \Pr\left(\left(\begin{matrix} t \\ \cdot \end{matrix}\right) \middle| \begin{matrix} - \\ - \\ - \\ \in \end{matrix}\right) \cdot \Pr\left(\begin{matrix} \in \\ \cdot \\ - \\ - \end{matrix} \middle| \begin{matrix} - \\ - \\ - \end{matrix}\right)} \quad [\text{Law of total probability}] \\
&= \frac{- \left[\begin{matrix} - \\ t \end{matrix} \prod_{\substack{j \in t \\ j > (q-1)}} (1 - \begin{matrix} - \\ j_t \end{matrix}) \right]}{\sum \left[\begin{matrix} - \\ \cdot \end{matrix} \left(\begin{matrix} - \\ t \end{matrix} \prod_{\substack{j \in t \\ j > (q-1)}} (1 - \begin{matrix} - \\ j_t \end{matrix}) \right) \right]}
\end{aligned}$$

As a result, the conditional expected value of the log-likelihood function is given by

$$\sum \sum \log \left(\begin{matrix} t \\ \cdot \end{matrix} \prod_{\substack{j \in t \\ j > \sigma}} (1 - \begin{matrix} - \\ j_t \end{matrix}) \right)$$

M-step: First, we update class membership probabilities for every segment $\in \{1, 2, \dots\}$:

$$= \sum \frac{\dots}{\dots}$$

and then optimize the conditional expected value of the log-likelihood function, obtained in the previous step, in terms of θ and σ :

$$\theta \sum \sum \log \left(\begin{matrix} t \\ \cdot \end{matrix} \prod_{\substack{j \in t \\ j > \sigma}} (1 - \begin{matrix} - \\ j_t \end{matrix}) \right)$$

which is solved using outer-approximation algorithm in Section A3.3.

A3.4.2. Estimation with panel data. In the EM algorithm above we assumed access to the aggregate level sales transaction data (i.e., sales transaction data without access to the customer tags). The EM algorithm is updated in the following way if we have access to the individual-level sales transaction data with n customers:

Initialization: we initialize the EM with a random allocation of individuals to one of the classes, resulting in an initial allocation $\mathcal{D} = \mathcal{D}_1 \cup \dots \cup \mathcal{D}_K$, which form a partition of the collection

of all the individuals. Then we set $\theta_h = |\mathcal{D}_h| / (\sum_{h \in \mathcal{H}} |\mathcal{D}_h|)$. Then θ_{it} (i.e., θ_{it}) and $\theta_{j \succ \sigma}^{it}$ for all $j \in \{1, \dots, q\}$, $\theta_{j \succ \sigma}^{it}$ are obtained by solving the following optimization problem:

$$\theta_h \sum_{i \in \mathcal{D}_h} \left(\log \theta_{it} + \sum_{\substack{j \in \{1, \dots, q\} \\ j \succ \sigma}} \log(1 - \theta_{j \succ \sigma}^{it}) \right)$$

which is solved by using the outer approximation algorithm in Section A3.3.

E-step: we compute θ_{it} , which is the membership probability of every individual i to belong to the segment h based on the parameter estimates $\{\theta_{it}, \theta_{j \succ \sigma}^{it}\}$ and the purchasing transactions data (θ_{it}) i :

$$\theta_{it} = \frac{\theta_{it} \prod_{j \in \{1, \dots, q\}} \theta_{j \succ \sigma}^{it} (1 - \theta_{j \succ \sigma}^{it})}{\sum_{h \in \mathcal{H}} \left[\theta_{it} \prod_{j \in \{1, \dots, q\}} \theta_{j \succ \sigma}^{it} (1 - \theta_{j \succ \sigma}^{it}) \right]}$$

M-step: next, we update class membership probabilities for every segment $h \in \{1, 2, \dots, q\}$:

$$\theta_h = \frac{\sum_{i \in \mathcal{D}_h} \theta_{it}}{\sum_{i \in \mathcal{D}} \theta_{it}}$$

and then optimize the conditional expected value of the log-likelihood function, obtained in the previous step, in terms of θ_{it} and $\theta_{j \succ \sigma}^{it}$:

$$\theta_h \sum_{i \in \mathcal{D}_h} \sum_{j \in \{1, \dots, q\}} \log \left(\theta_{it} \prod_{\substack{j \in \{1, \dots, q\} \\ j \succ \sigma}} (1 - \theta_{j \succ \sigma}^{it}) \right)$$

which is solved using the outer-approximation algorithm in Section A3.3. In order to simplify the optimization of the aforementioned conditional expected value of the log-likelihood function in Section 6, we assume that the ranking \succ is known a priori, i.e., the products are assumed to be ranked according to their sales.

A3.4.3. EM algorithm heuristics. Note that the proposed EM algorithm might become computationally challenging for large-scale problems as we need to run an outer-approximation algorithm for every h th iteration. Alternatively, we might further assume that the preference order \succ (i.e., \succ) over items in the product universe is known, e.g., we can rank the products according to their popularity in the sales transaction data or we can estimate the ranking from calibrating single class ICS model (see Section 4.1). In this case, the M step for h th iteration in the EM

algorithm reduces to solving a globally concave maximization problem with a unique, closed form solution (i.e., we don't need to apply an outer-approximation algorithm) given by:

$$= \frac{\sum \mathbb{I}[i_t = j]}{\sum \mathbb{I}[i_t = j] + \sum \mathbb{I}[i_t \in \mathcal{S}_t^j]}$$

which can be applied with aggregate level data (see Section A3.4.1), and

$$= \frac{\sum_i \sum_{it} \mathbb{I}[i_{it} = j]}{\sum_i \sum_{it} \mathbb{I}[i_{it} = j] + \sum_i \sum_{it} \mathbb{I}[i_{it} \in \mathcal{S}_{it}^j]}$$

which can be applied with panel data (see Section A3.4.2).

The CTC (i.e., general consideration - then - general choice) is the broadest class of consider-then-choose type of models where customers have heterogeneous preferences and consideration sets, i.e., before making a choice customers sample their preference order over the items in the product universe and the subset of items to consider from the general distributions over product rankings and consideration sets respectively.

A3.5.1. Estimation with aggregate level data. Similarly to the Section A3.4, we calibrate the CTC model by dividing transactions into segments such that customers in segment h sample their consideration sets based on the attention parameters θ_h and have their preferences characterized by the ranking \succ_h . Then the log-likelihood function can be represented in the following way

$$\log \mathcal{L}(\theta_h, \succ_h) = \sum_t \log \left(\sum_{j \in \mathcal{S}_t^h} \theta_h \prod_{j' \succ_h j_t} (1 - \theta_h) \right) \quad (\text{A24})$$

where $\theta_h \geq 0$ is the weight of the class h , s.t. $\sum_h \theta_h = 1$; \mathcal{S}_t^h denotes the set of ordered items at time t ; j_t denotes the product purchased at time t ; \mathcal{S}_t^h denotes the time horizon. Conceptually, we can obtain all the parameters of the CTC model (i.e., distributions over the preference lists and considerations sets) by maximizing the log-likelihood function above for a sufficiently large T .

Next we provide the initialization of the EM algorithm to calibrate the CTC model followed by the E and M steps of every iteration.

Initialization: we randomly allocate sales transaction to one of the H classes, resulting in an initial allocation $\mathcal{D} = \{\mathcal{D}^h\}_{h=1}^H$, which form a partition of the collection of all the transactions. Consequently, we set $\theta_h = |\mathcal{D}^h| / (\sum_h |\mathcal{D}^h|)$. Then \succ_h (i.e., \succ) and \mathcal{S}_t^h for all $h \in \{1, \dots, H\}$, $t \in \{1, \dots, T\}$ are obtained by solving the following optimization problem:

$$\max_{\theta_h, \succ_h} \sum_{t \in \mathcal{D}^h} \left(\log \theta_h + \sum_{j \in \mathcal{S}_t^h} \log(1 - \theta_h) \right)$$

which is solved by using the outer approximation algorithm for the ICS model in Section A3.3.

E-step: we compute π_{jt} , which is the membership probability of every transaction at time t to belong to the segment h based on the parameter estimates $\{\theta_h, \gamma_h\}$ and the purchasing transactions data (x_{jt}) :

$$\pi_{jt} = \frac{\theta_h \prod_{j \succ_h^{(q-1)} j_t} (1 - \theta_h)}{\sum_h \left[\theta_h \prod_{j \succ_h^{(q-1)} j_t} (1 - \theta_h) \right]}$$

M-step: first, we update class membership probabilities for every segment $h \in \{1, 2, \dots, H\}$:

$$\theta_h = \frac{\sum_t \pi_{jt} x_{jt}}{\sum_t \pi_{jt}}$$

and then optimize the conditional expected value of the log-likelihood function, obtained in the previous step, in terms of θ_h and γ_h for all $h \in \{1, 2, \dots, H\}$:

$$\sum_h \theta_h \sum_t \log \left(\prod_{\substack{j \in t \\ j \succ_h j_t}} (1 - \theta_h) \right)$$

which is solved by using the outer approximation algorithm for the ICS model in Section A3.3.

Note that in the proposed EM algorithm we need to apply the outer-approximation algorithm for every iteration. In order to reduce the computation time for the large-scale problems we might solve the optimization problem at M -step by ranking the products according to their popularity for each segment h . This way we can obtain the preference order γ_h for every segment h for m -th iteration. In this case, the M step in the EM algorithm reduces to solving a globally concave maximization problem with a unique and closed form solution given by:

$$\theta_h = \frac{\sum_t \mathbb{I}[x_{jt} = 1]}{\sum_t \mathbb{I}[x_{jt} = 1] + \sum_{j \succ_h j_t} \mathbb{I}[x_{jt} = 0]}$$

A3.5.2. Estimation with panel data. We update the EM algorithm above in the following way:

Initialization: we randomly allocate individuals to one of the H classes, resulting in an initial allocation $\mathcal{D} = \{\mathcal{D}_1, \dots, \mathcal{D}_H\}$. Consequently, we set $\theta_h = |\mathcal{D}_h| / (\sum_h |\mathcal{D}_h|)$. Then γ_h (i.e., \succ_h) and π_{jt} for all $h \in \{1, 2, \dots, H\}$, $j \in \mathcal{D}_h$ are obtained by solving the following optimization problem:

$$\sum_h \theta_h \sum_{j \in \mathcal{D}_h} \left(\log x_{jt} + \sum_{\substack{j \in \mathcal{D}_h \\ j \succ_h j_{it}}} \log(1 - \theta_h) \right)$$

which is solved by using the outer approximation algorithm for the ICS model in Section A3.3.

E-step: we compute π_{it}^i , which is the membership probability of every individual i to belong to segment i based on the parameter estimates $\{\theta_h, \beta_h, \gamma_h\}$ and the purchasing transactions data (x_{it}^i) :

$$\pi_{it}^i = \frac{\theta_h^i \prod_{j \in \mathcal{I}_h^{(q-1)}} (1 - \beta_{jit})}{\sum_{i=1}^I \theta_h^i \prod_{j \in \mathcal{I}_h^{(q-1)}} (1 - \beta_{jit})}$$

M-step: first, we update class membership probabilities for every segment $i \in \{1, 2, \dots, I\}$:

$$\theta_h^i = \frac{\sum_{i=1}^I \theta_h^i \prod_{j \in \mathcal{I}_h^{(q-1)}} (1 - \beta_{jit})}{\sum_{i=1}^I \theta_h^i \prod_{j \in \mathcal{I}_h^{(q-1)}} (1 - \beta_{jit})}$$

and then optimize the conditional expected value of the log-likelihood function, obtained in the previous step, in terms of θ_h^i and β_{jit} for all $i \in \{1, 2, \dots, I\}$:

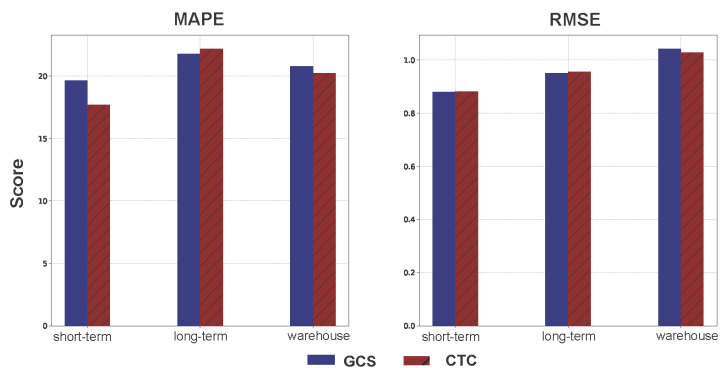
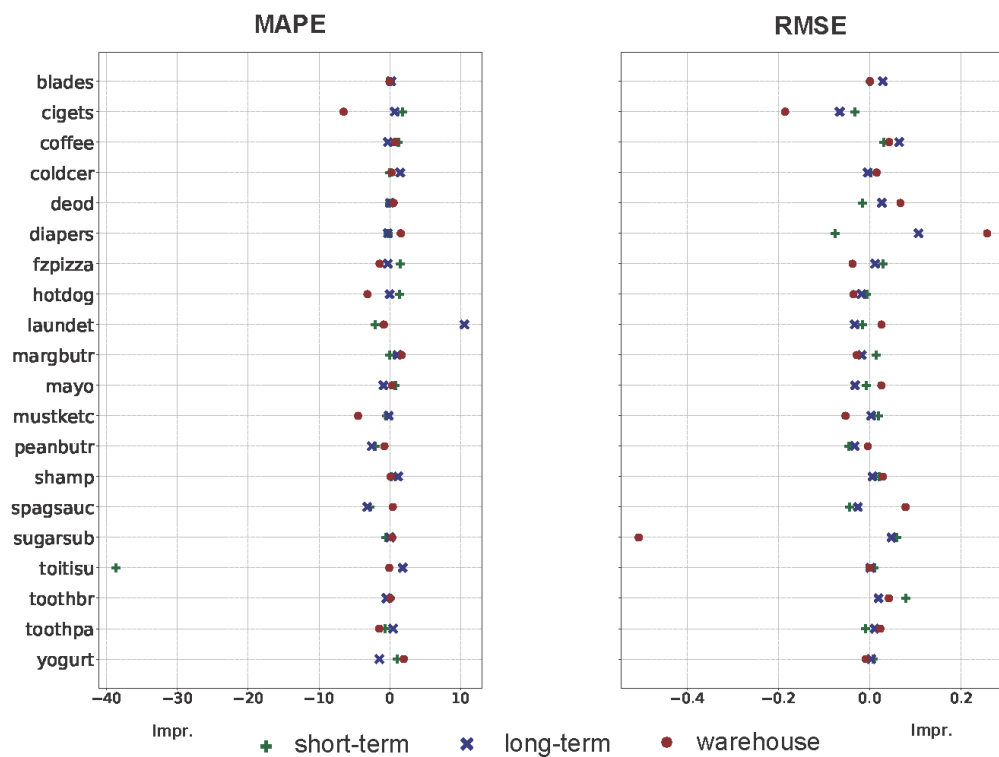
$$\theta_h^i \sum_{i=1}^I \sum_{j \in \mathcal{I}_h^{(q-1)}} \log \left(\prod_{j \in \mathcal{I}_h^{(q-1)}} (1 - \beta_{jit}) \right)$$

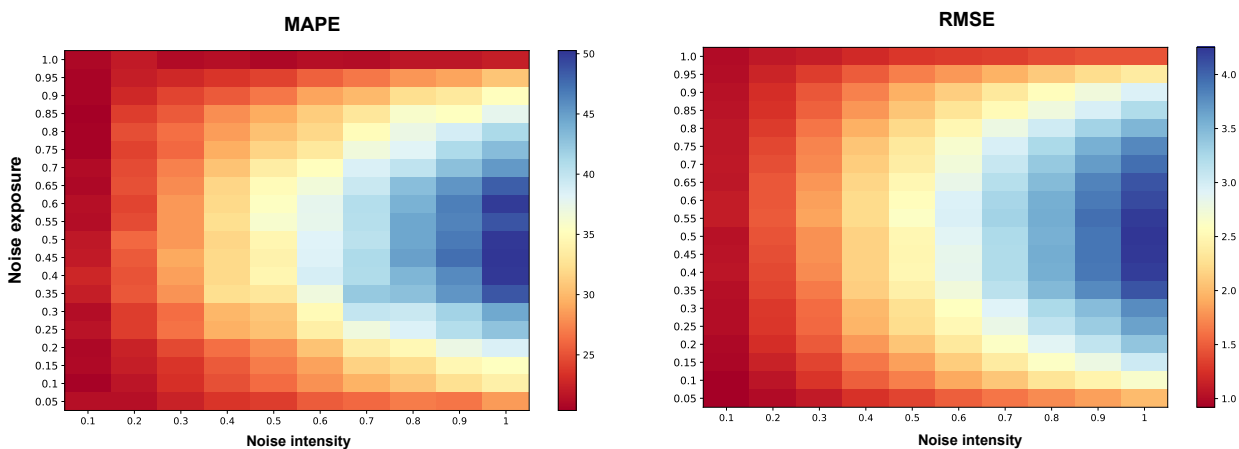
which is solved by using the outer approximation algorithm for the ICS model in Section A3.3.

In this section, we compare the prediction performance of the GCS with the CTC based on IRI dataset. It follows from Figures A6 and A7 that the GCS provides higher prediction accuracy than the CTC under the long-term forecast scenario whereas the CTC outperforms the GCS under the warehouse forecast scenario. As a result, we can not claim dominance of the GCS or the CTC.

In this section, we summarize the results of the extensive synthetic experiments conducted in order to check the robustness of the simulation results reported in Section 5, when we use the rank-based model as ground truth instead of the MNL. Recall that in Section 5, we compared the predictions of the MNL model against the ICS model to understand the conditions under which the MNL benchmark outperforms the consider-then-choose model in the presence of noise in the offer sets.

The setup of the experiments in this section is identical to the one in Section 5. In the new set of experiments we simulate sales transaction data according to the rank-based model with fifteen





customer types where type c customers are making purchases according to the ranking σ_c , i.e., when faced with a given choice set customers are assumed to purchase the available option that ranks highest in their preference list. To this end, we randomly sample the set of C fifteen rankings, each corresponding to a particular class, and also assume the equal probability of each class. In order to calibrate this rank-based model we exploit the EM algorithm proposed by [van Ryzin and Vulcano \(2017\)](#). This algorithm relies on the assumption that the set of rankings is known. Therefore, we initialize this algorithm with thirty preference orders such that fifteen of them are the rankings from the ground truth and the remaining fifteen rankings are sampled randomly. Moreover, we start the EM algorithm with an equal probability of sampling each preference list.

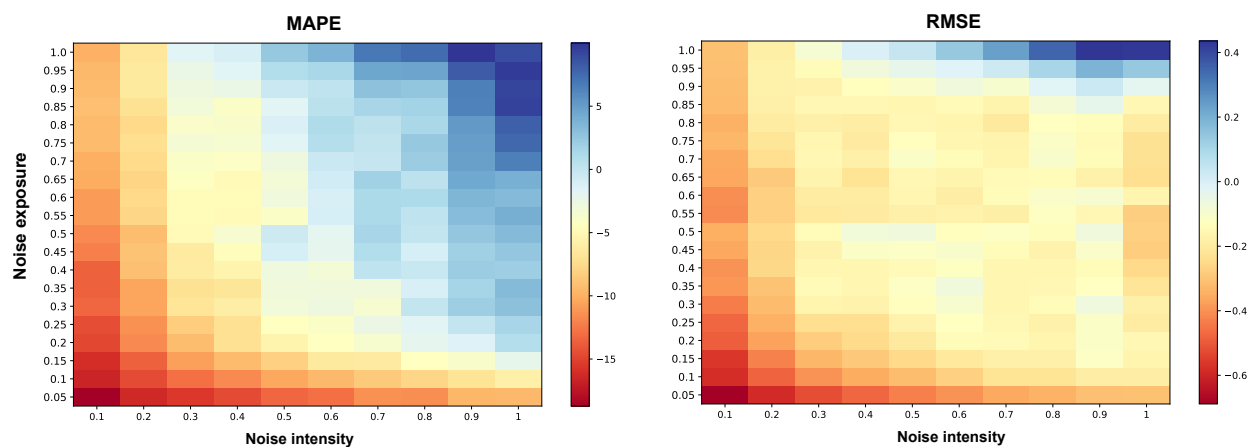
Similarly to Section 5, in Figures A8 and A9 we present the heatmaps of the prediction scores under the rank-based model, and the prediction scores improvements of the ICS model versus the rank-based model, respectively. And in Tables A2 and A3, we report the results for the regression (16), where the dependent variables are the prediction scores and the improvement scores, respectively. The main insights remain the same – the results of this extensive simulation study demonstrate that choice models based on the consider-then-choose framework are more robust to noise in order sets than their classical counterparts, i.e., the ICS model outperforms the ground truth rank-based model under noisier regimes.

With the objective of benchmarking the ICS model with a more competitive variant of the MNL, we resort to the LC-MNL. As it was mentioned above, we estimate the LC-MNL model for $C = 1, 2, \dots, 5$, classes and report the best performance measure from these five variants. Figure A11 illustrates the heatmaps of the MAPE and RMSE prediction score improvements of the ICS model

	Model (1) Score	Model (2) Score	Model (3) Score	Model (4) Score	Model (5) Score
	38.341*** (23.312)				48.526*** (15.904)
		32.639*** (20.358)			-9.259*** (-3.425)
			30.882*** (5.470)		42.516*** (29.661)
				4.827* (1.657)	13.450*** (19.378)
	7.072*** (6.930)	15.593*** (19.324)	17.891*** (8.679)	26.428*** (19.091)	-13.927*** (-14.501)
No. Observations:	200	200	200	200	200
R-squared:	0.733	0.677	0.131	0.014	0.955
Adj. R-squared:	0.732	0.675	0.127	0.009	0.954

t

* $p < .$ ** $p < .$ *** $p < .$

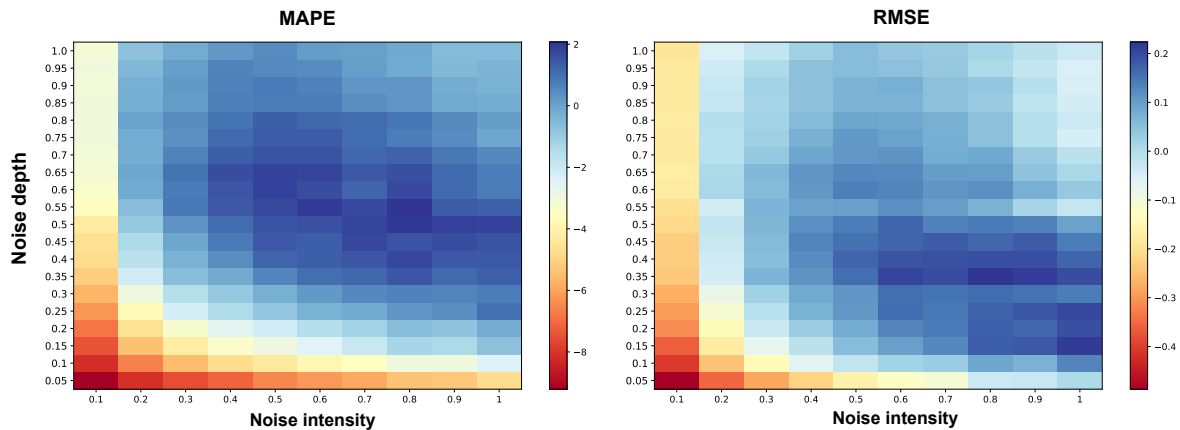


versus the LC-MNL when the ground truth is an MNL model as in Section 5. As expected, Figure A11 confirms that the LC-MNL is a more competitive benchmark. However, our qualitative results remain the same: the ICS model outperforms the LC-MNL model, which subsumes the ground truth MNL model, under sufficiently noisy regimes.

	Model (1) Impr.	Model (2) Impr.	Model (3) Impr.	Model (4) Impr.	Model (5) Impr.
	16.211*** (16.814)				23.908*** (13.640)
		13.646*** (14.968)			-6.997*** (-4.506)
			1.991 (0.697)		12.733*** (15.464)
				9.837*** (8.231)	12.419*** (31.147)
	-12.304*** (-20.567)	-8.641*** (-18.832)	-4.050*** (-3.884)	-6.917*** (-12.178)	-22.532*** (-40.841)
No. Observations:	200	200	200	200	200
R-squared:	0.588	0.531	0.002	0.255	0.933
Adj. R-squared:	0.586	0.528	-0.003	0.251	0.931

t

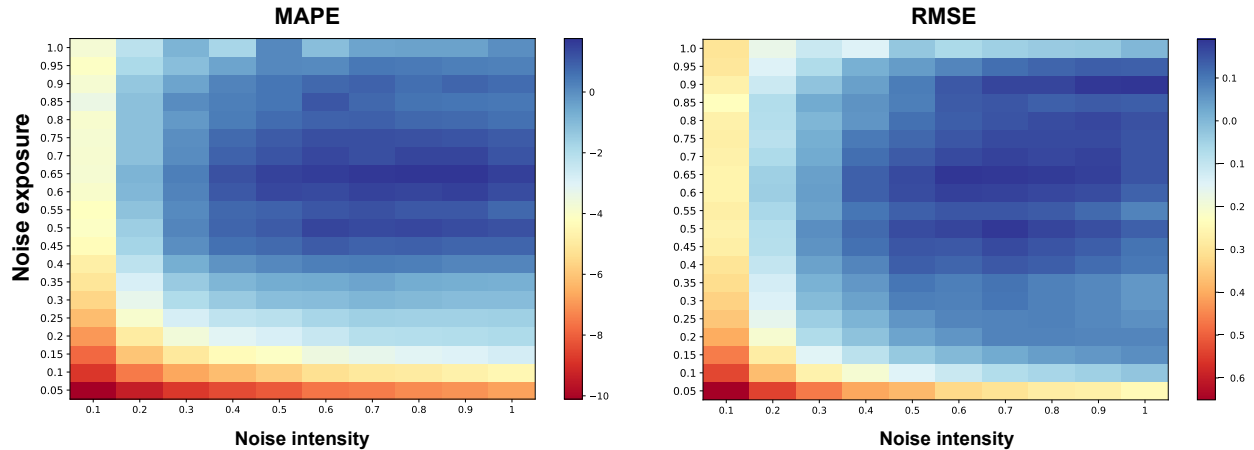
* $p < .$ ** $p < .$ *** $p < .$



η

γ

In this section, we provide more details on how the one-directional cannibalization property of the ICS model helps it to outperform the ground-truth MNL model in the simulation study in Section 5.



η

γ

This is just one potential mechanism that can explain the superior performance of the ICS model in noisy regimes. To streamline our analysis, we focus on the specific case of the simulation study in Section 5 where $\gamma = 0.5$ and $\eta = 1$, which corresponds to the case of maximum asymmetry between the training and testing datasets (i.e., the probability that a product is only offered in either the training or the testing dataset is 0.5, and every product in the exposure set will be added to either the training or the testing datasets). As it was mentioned in Section 5, the one-directional cannibalization property of the ICS model provides robustness to the noise in the offer sets because it alleviates its impact on the demand prediction for the higher ranked items (i.e., the presence or absence of lower-ranked products in the offer set does not affect the demand prediction for higher ranked products). We confirm this intuition by estimating the following regression specification:

$$= + \cdot + \tag{A25}$$

where i corresponds to item $i \in N$ and j corresponds to a specific instance, i.e., to each of the 1,000 generated instances. Δ_j is the improvement of item i demand forecasting error in instance j obtained by the ICS model over MNL, i.e., $\Delta_j = \text{MAPE}_{MNL} - \text{MAPE}_{ICS}$, where we rely on two ways to measure the item's demand forecasting accuracy: (1) absolute percentage error ($\text{MAPE} = \frac{|\hat{a}_j - a_j|}{a_j}$), where a_j denotes the observed sales for i in the test dataset and \hat{a}_j denotes our prediction and (2) demand prediction error ($\text{RMSE} = \frac{(\hat{a}_j - a_j)^2}{\sum_{a_j \in N} a_j}$). Note that the former and the latter prediction errors are factored into the computations of the MAPE and RMSE scores, respectively, and that they are expressed in percentage points so that the score improvement is also expressed in percentage points. r_j is the position of item i in the ranking inferred

from the ICS model calibration on instance i . rank_i is a categorical variable to control for the instance-level fixed effects. Table A4 presents the results from the regression specification (A25) where the improvement by the ICS model over MNL is measured using score_i in columns (1) and (3), while the improvement by the ICS model over MNL using the score_i score is presented in columns (2) and (4). In columns (1) and (2) we do not control for the instance-level fixed effects, while in columns (3) and (4) we do. It follows from Table A4 that the coefficient of the rank_i variable is negative and statistically significant, with a magnitude that implies a relevant decrease in score_i for the low-ranked products, which confirms the intuition that the ICS model is especially competitive when predicting the market shares of the top-ranked items because those items are immune to the presence or absence of lower-ranked products in the offer sets. Note that this finding is consistent across all four columns in Table A4 (with and without instance-level fixed effects, under both prediction scores). Note that the rank_i in all the columns in Table A4 is relatively small which is not unexpected given that we only have a single covariate in the regression models and there might be other factors (in addition to the rank_i variable) which could explain the variation in the outcome variable score_i .

	Model (1) Impr.	Model (2) Impr.	Model (3) Impr.	Model (4) Impr.
	-13.385*** (-3.274)	-1.428*** (-5.987)	-13.385*** (-3.177)	-1.428*** (-5.815)
	331.279*** (12.408)	27.625*** (14.699)		
Instance FE:	No	No	Yes	Yes
No. Observations:	15,000	15,000	15,000	15,000
R-squared:	0.001	0.003	0.001	0.003

t

* $p < .1$. ** $p < .05$. *** $p < .01$.

In this section, we show the robustness of our prediction results in Section 6 to different specifications of the prediction metrics. To this end, we focus on the variations of RMSE and MAPE metrics when aggregating predictions either over a one week intervals (i.e., RMSE_{1w} and MAPE_{1w}).

or over every order set (i.e., RMSE and MAPE). More formally, RMSE and MAPE metrics when aggregating predictions over one-week intervals can be computed in the following way

$$\text{MAPE} = \frac{100}{\sum_{t=1}^T \sum_{j \in I_t} | \frac{-}{10+} |} \quad (\text{A26})$$

$$\text{RMSE} = 100 \sqrt{\frac{1}{\sum_{t=1}^T \sum_{j \in I_t} \frac{(\frac{-}{10+})}{(\sum_{j \in I_t} \frac{-}{10+})}}}} \quad (\text{A27})$$

where I_t is the set of items that were ordered to customers in week t , T is the total number of weeks in the test dataset, $y_{j,t}$ is the observed number of times product j was purchased in week t , and $\hat{y}_{j,t}$ is the predicted number of times product j to be purchased in week t . In the same spirit, we can compute RMSE and MAPE metrics when aggregating predictions over each order set as follows

$$\text{MAPE} = \frac{100}{\sum_{m=1}^M \sum_{j \in I_m} | \frac{-}{10+} |} \quad (\text{A28})$$

$$\text{RMSE} = 100 \sqrt{\frac{1}{\sum_{m=1}^M \sum_{j \in I_m} \frac{(\frac{-}{10+})}{(\sum_{j \in I_m} \frac{-}{10+})}}}} \quad (\text{A29})$$

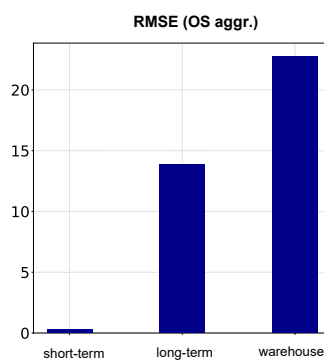
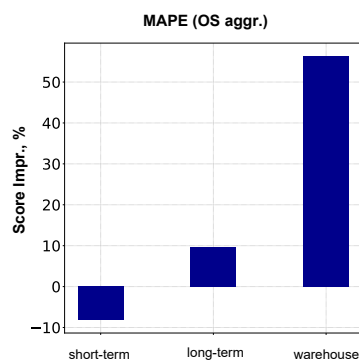
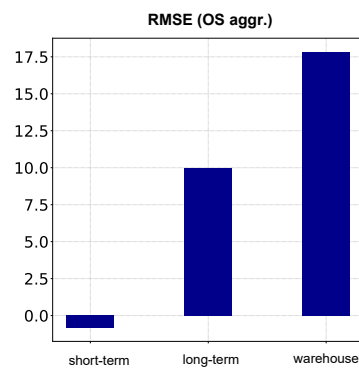
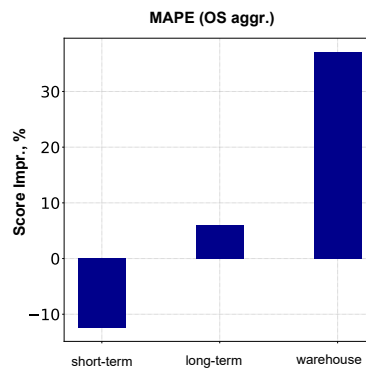
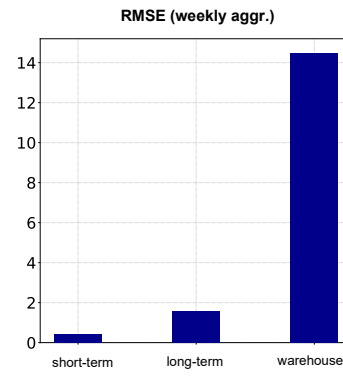
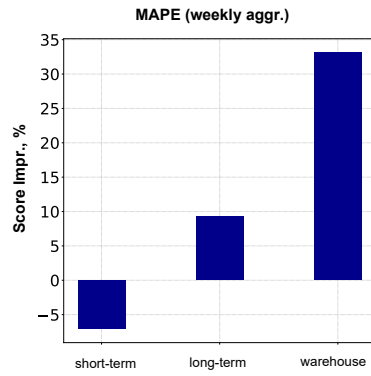
where I_m is the m th order set, M is the total number of different order sets in the test dataset, $y_{j,m}$ is the observed number of times product j was purchased under the order set m , and $\hat{y}_{j,m}$ is the predicted number of times product j to be purchased under the order set m . Alternatively, we could aggregate predictions based on the true order sets in the following way

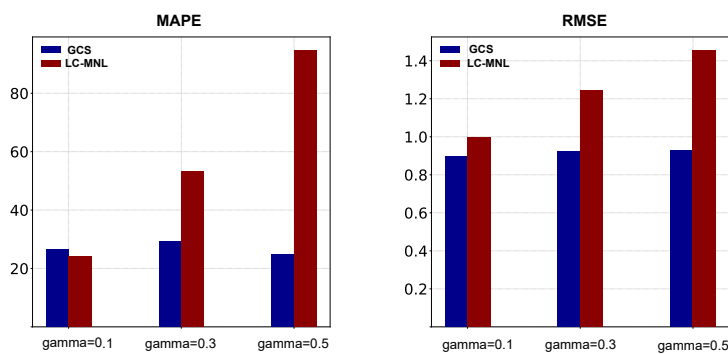
$$\text{MAPE} = \frac{100}{\sum_{m'=1}^{M'} \sum_{j \in I_{m'}} | \frac{-}{10+} |} \quad (\text{A30})$$

$$\text{RMSE} = 100 \sqrt{\frac{1}{\sum_{m'=1}^{M'} \sum_{j \in I_{m'}} \frac{(\frac{-}{10+})}{(\sum_{j \in I_{m'}} \frac{-}{10+})}}}} \quad (\text{A31})$$

where $I_{m'}$ is the m' th true order set, M' is the total number of different true order sets in the test dataset.

In Figure A12, we exhibit improvements of the GCS over the LC-MNL under MAPE (left panel) and RMSE (right panel), averaging across 20 product categories, for the three different scenarios. From these panels, we observe that the GCS significantly improves over the LC-MNL once we shift from short to long-term and from long-term to warehouse forecasts based on the MAPE and RMSE scores. Figure A13 (resp. Figure A14) presents qualitatively the same results when we compute prediction metrics based on the MAPE (resp. MAPE) and RMSE (resp. RMSE) scores. As a result, we conclude that the direction of our results in Section 6 is robust to the definition of the prediction metrics.

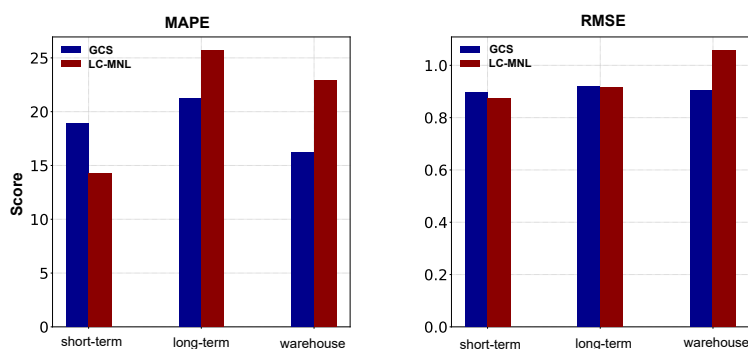


 η γ

In this section, we show the robustness of our prediction results in Section 6 to an alternative noise generation process. To this end, we add the noise only to the test dataset and we do it in a similar way to the simulation study presented in Section 5. Figure A15 illustrates the predictive performance of the GCS and the LC-MNL models when η is equal to 0.5 and γ takes three different values (0.1, 0.3, and 0.5). It follows from Figure A15 that our results are qualitatively the same under this alternative noise regime – the relative predictive performance of the GCS over the LC-MNL improves with the level of noise (measured by γ) in the other sets.

In this section, we show the robustness of our prediction results in Section 6 to the scenario when we use the same noise generation process to alter both the training and test datasets. Otherwise, the setup of the experiments is the same as in Section 6. Figure A16 illustrates that the findings are qualitatively the same under the updated setup – the relative predictive performance of the GCS over the LC-MNL improves with the level of noise in the other sets.

In this section, we calibrate the Logistic-based ICS (L-ICS) and MNL models accounting for car features and discuss the modeling assumptions based on the car sharing dataset (see data descriptive statistics in Table A5). We also provide explanatory analysis of choice models in order to gain insights about the consideration set formation of renters using the car feature information. In addition, we address the problem of a potential price endogeneity in our empirical explanatory



analysis. We argue that in our setting we are unlikely to have any price endogeneity problems while calibrating the models. We conclude this section by presenting Figure A17 which illustrates an instance of the decision tree obtained after fitting the DT-ICS model.

We start this section by calibrating the single-class L-ICS model with features to examine the extent to which various variables impact the consideration set structure. Assuming that the cars are ranked according to their popularity among renters (see Section 7.2), the problem of fitting the L-ICS model is the one of estimating the coefficients β . The car features available to the renters through the online platform are divided into three groups: (1) car brand; (2) car location type and accessibility, including car access (i.e., open or closed), car location hours (i.e., 24 hours or restricted), car location type (i.e., garage, street, surface lot, or valet); and the third group including (3a) car type (i.e., economy, standard, fullsize, SUV, trucks, luxury), and (3b) car features: hourly price, car age, and some other various binary car features such as transmission, premium wheels, power seats, bluetooth/wireless, leather interior, sunroof/moonroof, premium sound, power windows, GPS navigation system, roof rack, tinted windows. Assuming that the error terms ϵ_{it} are logistically distributed, we estimate the β vector using logistic regression analysis.

The results for the L-ICS model appear in the first column of Table A6. In the middle column, the table lists the average marginal effects (AME) of the L-ICS model when all the covariates are at their mean. Then we also calibrate the usual linear-in-parameters MNL model where the utility from reserving the car alternative i is represented with linear in parameters function $U_i = \beta_i + \beta'X_i$. On the right, Table A6 presents the estimates of the MNL model parameters. However, the interpretation of the β vector for the L-ICS and MNL models is different. The parameters of

	Mean	Std.	Min	Max
Brands				
<i>Acura</i>	2.52%	15.68%	0%	100%
<i>Audi</i>	4.54%	20.82%	0%	100%
<i>BMW</i>	11.73%	32.18%	0%	100%
<i>Buick</i>	0.21%	4.61%	0%	100%
<i>Chevrolet</i>	0.79%	8.84%	0%	100%
<i>Chrysler</i>	0.41%	6.37%	0%	100%
<i>Dodge</i>	0.82%	9%	0%	100%
<i>Fiat</i>	0.9%	9.42%	0%	100%
<i>Ford</i>	2.63%	16.01%	0%	100%
<i>Honda</i>	16.77%	37.36%	0%	100%
<i>Hyundai</i>	3.42%	18.16%	0%	100%
<i>Infiniti</i>	0.21%	4.61%	0%	100%
<i>Jeep</i>	0.41%	6.39%	0%	100%
<i>Kia</i>	0.49%	6.95%	0%	100%
<i>Land Rover</i>	0.17%	4.14%	0%	100%
<i>Lexus</i>	1.24%	11.06%	0%	100%
<i>Mazda</i>	3.44%	18.22%	0%	100%
<i>Mercedes Benz</i>	3.73%	18.95%	0%	100%
<i>Mercury</i>	0.07%	2.59%	0%	100%
<i>Mini</i>	7.66%	26.59%	0%	100%
<i>Mitsubishi</i>	0.68%	8.21%	0%	100%
<i>Nissan</i>	4.05%	19.71%	0%	100%
<i>Pontiac</i>	0.21%	4.57%	0%	100%
<i>Porsche</i>	1.5%	12.14%	0%	100%
<i>Saab</i>	0.03%	1.73%	0%	100%
<i>Saturn</i>	0.28%	5.25%	0%	100%
<i>Scion</i>	0.62%	7.85%	0%	100%
<i>Subaru</i>	3.71%	18.89%	0%	100%
<i>Suzuki</i>	0.53%	7.29%	0%	100%
<i>Smart</i>	5.44%	22.69%	0%	100%
<i>Tesla</i>	1.42%	11.83%	0%	100%
<i>Toyota</i>	10.33%	30.43%	0%	100%
<i>Volkswagen</i>	7.54%	26.41%	0%	100%
<i>Volvo</i>	1.53%	12.28%	0%	100%
Car types				
<i>Economy</i>	14.83%	35.54%	0%	100%
<i>Standard</i>	48.83%	49.99%	0%	100%
<i>Fullsize</i>	19.56%	39.67%	0%	100%
<i>SUV</i>	9.41%	29.2%	0%	100%
<i>Trucks</i>	3.31%	17.88%	0%	100%
<i>Luxury</i>	4.06%	19.74%	0%	100%
Car location type and accessibility				
<i>Car access [open]</i>	81.89%	38.51%	0%	100%
<i>Car access hours [all hours]</i>	93.52%	24.61%	0%	100%
<i>Car location type [garage]</i>	28.50%	45.14%	0%	100%
<i>Car location type [street]</i>	23.86%	42.62%	0%	100%
<i>Car location type [surface lot]</i>	43.85%	49.62%	0%	100%
<i>Car location type [valet]</i>	0.24%	4.84%	0%	100%
Car features				
<i>Price (per hour)</i>	8.63	4.61	2.0	300.0
<i>Car age</i>	5.32	3.18	-0.3	18.3
<i>Transmission [automatic]</i>	95.21%	21.35%	0%	100%
<i>Premium wheels</i>	29.38%	45.55%	0%	100%
<i>Power seats</i>	46.88%	49.90%	0%	100%
<i>Bluetooth/wireless</i>	33.74%	47.28%	0%	100%
<i>Leather interior</i>	53.56%	49.87%	0%	100%
<i>Sunroof/moonroof</i>	53.48%	49.88%	0%	100%
<i>Premium sound</i>	46.25%	49.86%	0%	100%
<i>Power windows</i>	92.90%	25.68%	0%	100%
<i>GPS navigation system</i>	23.05%	42.11%	0%	100%
<i>Roof rack</i>	6.98%	25.48%	0%	100%
<i>Tinted windows</i>	13.24%	33.89%	0%	100%

the L-ICS model listed in the left column of Table A6 show the estimated impact of exogenously imposed changes in car features on consideration set formation. Rather, the parameters of the MNL model in the right column of Table A6 show the influence of car features on the customer's choices, i.e., revealed preferences. Notably, quite a few coefficients (17 out of 53) estimated based on L-ICS and MNL models are not aligned, i.e., the covariates that increase (or decrease) the likelihood of considering the car under the L-ICS model might not necessarily increase (or decrease) the likelihood of booking the car under the MNL model, e.g., the utility of the renter from considering the car brand Jeep is higher by 0.98 ($t = 4.9$, $p < 0.01$) than the utility from considering the baseline brands while the utility of the renter from reserving the same brand under the MNL model is lower by 0.93 ($t = -6.2$, $p < 0.01$) than the utility from reserving baseline brands. Also, some of the covariates (7 out of 53) that are statistically significant in explaining the choice of renters under the MNL model might be statistically insignificant under the L-ICS model, e.g., the utility from choosing a car parked in the street is lower ($t = -8.33$, $p < 0.01$) than the utility from choosing the car located in the valet parking area while the discrepancy between these two parking location types is insignificant ($t = 1.36$, $p > 0.10$) under the L-ICS model. The price and car age coefficients are statistically significant and negative for both the L-ICS and MNL models. However, the impact of an additional \$1 increase in the car hourly rental price on the utility from considering the vehicle is equivalent to the car being 0.52 years older, while the impact of an additional \$1 increase in the car hourly rental price on the utility from booking the vehicle under the MNL model is equivalent to the car being 3.75 years older. According to these findings, car age plays a relatively more important role during the formation of the consideration set in the L-ICS model compared to its role in the choice process under the MNL model.

Next, we consider three types of car attributes (i.e., car brand, car location type and accessibility, and car type and features), with the objective of empirically verifying their impact on consideration set formation under the L-ICS model and on the choice probabilities under the MNL model. Models 1, 2, and 3 (both under the L-ICS and MNL) incorporate all the covariates except car types and features, car location type and accessibility, and brands, respectively, e.g., Model 1 excludes car types and features while including all the other covariates. According to Table A7, the car type and features attributes are less statistically significant than car brand attributes under the L-ICS model, whereas the opposite effect takes place under the MNL model. These findings are robust to the various measures of statistical significance and goodness-of-fit presented in Table A7 such as LL, AIC, BIC, Likelihood Ratio (LR) statistics, and Wald statistics. Overall, it is implied that car location type and accessibility play the least important role both for consideration set formation and for the final choice decision. The renters are likely to build their consideration sets based on car brands rather than on car properties, even though while evaluating alternatives towards choice customers are likely to pay more attention to car properties rather than to car brands.

	L-ICS		AME (L-ICS)		MNL	
	Coeff.	Std.err.	Coeff.	Std. err.	Coeff.	Std. err.
Brands						
<i>Acura</i>	0.29**	0.10	0.067**	0.024	-0.33***	0.094
<i>Audi</i>	-0.11	0.097	-0.025	0.023	0.22*	0.088
<i>BMW</i>	0.0061	0.089	0.0014	0.021	0.14	0.083
<i>Chrysler</i>	-0.19	0.15	-0.044	0.035	-1.24***	0.13
<i>Dodge</i>	0.21	0.12	0.048	0.029	-0.0095	0.11
<i>Fiat</i>	1.01***	0.15	0.24***	0.035	-0.030	0.11
<i>Ford</i>	-0.50***	0.095	-0.12***	0.022	0.055	0.087
<i>Honda</i>	0.13	0.086	0.029	0.020	0.0011	0.079
<i>Hyundai</i>	0.010	0.094	0.0024	0.022	-0.064	0.085
<i>Infiniti</i>	0.21	0.25	0.049	0.059	0.30	0.24
<i>Jeep</i>	0.98***	0.20	0.23***	0.047	-0.93***	0.15
<i>Kia</i>	0.36*	0.15	0.085*	0.035	-0.47***	0.13
<i>Land Rover</i>	1.01***	0.25	0.24***	0.058	0.74***	0.20
<i>Lexus</i>	0.043	0.11	0.0100	0.027	0.20	0.11
<i>Mazda</i>	0.27**	0.098	0.063**	0.023	0.070	0.089
<i>Mercedes Benz</i>	-0.32**	0.100	-0.074**	0.023	0.20*	0.090
<i>Mercury</i>	2.53***	0.75	0.59***	0.18	-0.064	0.29
<i>Mini</i>	0.37***	0.094	0.086***	0.022	0.23**	0.086
<i>Mitsubishi</i>	-0.37**	0.13	-0.087**	0.030	-0.51***	0.12
<i>Nissan</i>	0.39***	0.094	0.090***	0.022	0.11	0.085
<i>Pontiac</i>	0.87***	0.23	0.20***	0.054	0.54**	0.18
<i>Porsche</i>	-0.24*	0.11	-0.057*	0.027	0.45***	0.11
<i>Scion</i>	0.74***	0.17	0.17***	0.039	-0.67***	0.14
<i>Subaru</i>	0.32***	0.096	0.075***	0.023	0.43***	0.084
<i>Suzuki</i>	-0.071	0.15	-0.017	0.036	0.47**	0.16
<i>Smart</i>	-0.15	0.096	-0.036	0.022	0.090	0.084
<i>Tesla</i>	1.29***	0.16	0.30***	0.038	3.24***	0.15
<i>Toyota</i>	0.27**	0.089	0.063**	0.021	0.13	0.080
<i>Volkswagen</i>	0.26**	0.091	0.061**	0.021	0.46***	0.084
<i>Volvo</i>	1.26***	0.13	0.29***	0.029	0.26*	0.11
<i>Baseline brands</i>	Baseline		Baseline		Baseline	
Car types						
<i>Economy</i>	0.36***	0.068	0.084***	0.016	-0.28***	0.055
<i>Standard</i>	0.29***	0.060	0.068***	0.014	-0.13**	0.050
<i>Fullsize</i>	0.34***	0.066	0.079***	0.015	-0.36***	0.054
<i>SUV</i>	0.055	0.070	0.013	0.016	-0.43***	0.057
<i>Luxury</i>	0.37***	0.083	0.087***	0.019	-0.19**	0.070
<i>Trucks</i>	Baseline		Baseline		Baseline	
Car location type and accessibility						
<i>Car access [open]</i>	-0.36***	0.029	-0.084***	0.0067	0.029	0.026
<i>Car access hours [all hours]</i>	-0.089*	0.045	-0.021*	0.010	-0.097*	0.038
<i>Car location type [garage]</i>	-0.24***	0.061	-0.057***	0.014	0.13*	0.052
<i>Car location type [street]</i>	0.080	0.059	0.019	0.014	-0.40***	0.048
<i>Car location type [surface lot]</i>	-0.19***	0.057	-0.044***	0.013	0.27***	0.048
<i>Car location type [valet]</i>	Baseline		Baseline		Baseline	
Car features						
<i>Price (per hour)</i>	-0.022***	0.0033	-0.0051***	0.00077	-0.12***	0.0042
<i>Car age</i>	-0.042***	0.0039	-0.0099***	0.00091	-0.032***	0.0033
<i>Transmission [automatic]</i>	0.34***	0.046	0.080***	0.011	0.45***	0.037
<i>Premium wheels</i>	-0.0025	0.025	-0.00059	0.0058	-0.18***	0.021
<i>Power seats</i>	-0.21***	0.024	-0.048***	0.0055	0.043*	0.021
<i>Bluetooth/wireless</i>	-0.13***	0.025	-0.031***	0.0059	-0.29***	0.021
<i>Leather interior</i>	0.087**	0.030	0.020**	0.0070	0.12***	0.025
<i>Sunroof/moonroof</i>	0.0011	0.027	0.00026	0.0064	0.14***	0.024
<i>Premium sound</i>	0.25***	0.027	0.059***	0.0063	-0.14***	0.022
<i>Power windows</i>	-0.0058	0.042	-0.0013	0.0099	0.40***	0.036
<i>GPS navigation system</i>	-0.085**	0.029	-0.020**	0.0067	0.18***	0.023
<i>Roof rack</i>	0.16***	0.046	0.036***	0.011	-0.26***	0.037
<i>Tinted windows</i>	-0.087**	0.030	-0.020**	0.0070	-0.30***	0.026
<i>Constant</i>	-0.13	0.15				
No. of obs.	26791		26791		26791	
AIC	76980.3				69788.7	
BIC	77464.8				70309.2	
Log likelihood	-38436.1				-34841.4	
Pseudo R^2 square	0.024					

* $p < .$, ** $p < .$, *** $p < .$.

		Excluded groups	Log-like	AIC	BIC	LR	Wald
L-ICS	Model 1	Car types and features	-38664.7	77403.3	77735.3	457.04	449.39
	Model 2	Car location type and accessibility	-38552.1	77202.2	77641.9	231.94	231.67
	Model 3	Brands	-38840.7	77729.4	77944.7	809.07	764.90
MNL	Model 1	Car types and features	-35587.4	71246.8	71600.3	1492.05	1385.99
	Model 2	Car location type and accessibility	-35205.7	70507.5	70978.8	728.75	705.53
	Model 3	Brands	-35534.9	71115.8	71341.6	1387.03	1259.70

In this section, we further discuss the assumptions imposed by the CTC models with features that we take into account when calibrating the models. And then we also address the problem of a potential price endogeneity in our empirical explanatory analysis. We argue that in our setting we are unlikely to have any price endogeneity problems estimating the models.

A5.2.1. Semiparametric approach. Using the semiparametric approach in order to calibrate the two-stage CTC model, we assume that renters form their consideration set taking into account car features. Then we assume that during the second stage renters choose the most preferred car among the considered ones according to the preference order over the universe of car alternatives, which remains the same over time, i.e., the ranking is fixed over time. Modeling the second stage choice process this way, we do not parameterize the ranking which implies that the cars are assumed to have the same attributes over time. In this subsection we justify this assumption according to our dataset.

We start by analyzing the variation of the hourly price parameter over car alternatives. In Table A8, we report that the average coefficient of variation (CV) of the hourly price across all the car alternatives is around 5% while owners of cars listed on average around two different values of the price. Moreover, the most frequently used value of the hourly price corresponds to 78% of the car rentals and the second most frequently used value of the hourly price corresponds to 16% of the car rentals. The low variation of the rental price is explained by the policies of the online platform, for the time span of the dataset, which allows the owners to choose the price by themselves, i.e., the platform as a central agent did not dynamically adjust the listed rental price to efficiently match demand and supply as opposed to many ride-sharing platforms (e.g., Uber, Lyft) which optimize the price of the ride to match riders with drivers on-demand. Then in the same Table A8 we can also observe that more than 98% of car owners did not alter their car access (i.e., open or closed), access hours (i.e., 24 hours or restricted), and location type (i.e., garage, street, surface lot, or valet).

Additional descriptive statistics:	
<i>Number of rentals</i>	26791
<i>Number of car owners</i>	514
<i>Number of available alternatives (within 0.3 mile)</i>	5.7
<i>Rental duration (days)</i>	0.62
<i>Rental request in advance (days)</i>	1.24
<i>Price CV (averaging over car owners)</i>	0.053
<i>Average number of price modes</i>	2.33
<i>The most frequent price (percentage)</i>	0.78
<i>The second most frequent price (percentage)</i>	0.16
<i>Average number of car access modes</i>	1.07
<i>The most frequent car access (percentage)</i>	0.99
<i>Average number of car access hours modes</i>	1.03
<i>The most frequent car access hours (percentage)</i>	0.99
<i>Average number of car location type modes</i>	1.10
<i>The most frequent car location type (percentage)</i>	0.98

A5.2.2. Price endogeneity problem. Next, we want to address the concerns of potential price endogeneity in our empirical analysis. First of all, estimating the demand with personalized data significantly alleviates the price endogeneity problem since each renter has only a trivial influence on the number of cars supplied and the market rental price, while the empirical work with aggregate-level transaction data is more likely to face very severe endogeneity issues. Nevertheless, having access to individual consumer data is not always a big advantage because individuals demand could be correlated. For example, we might have unobservable demand or supply shocks if a local convention was organized on a particular day that might shift the demand curve. In this case, we need to use instrumental variables to address the endogeneity problem. The natural approach, in this case, would be to use the typical Hausman-style instrument (Hausman, 1996), i.e., the average rental price of similar cars in other geographical locations. However, in our dataset we are highly unlikely to have any price endogeneity issues because the rental price variation of the listed cars is very insignificant as it was discussed above, i.e., the price does not react to any unobservable shocks (see Table A8).

