



## Management Science

### MANAGEMENT SCIENCE



Publication details, including instructions for authors and subscription information:  
<http://pubsonline.informs.org>

## A Nonparametric Approach to Modeling Choice with Limited Data

Vivek F. Farias, Srikanth Jagabathula, Devavrat Shah,

To cite this article:

Vivek F. Farias, Srikanth Jagabathula, Devavrat Shah, (2013) A Nonparametric Approach to Modeling Choice with Limited Data. Management Science 59(2):305-322. <https://doi.org/10.1287/mnsc.1120.1610>

Full terms and conditions of use: <http://pubsonline.informs.org/page/terms-and-conditions>

This article may be used only for the purposes of research, teaching, and/or private study. Commercial use or systematic downloading (by robots or other automatic processes) is prohibited without explicit Publisher approval, unless otherwise noted. For more information, contact [permissions@informs.org](mailto:permissions@informs.org).

The Publisher does not warrant or guarantee the article's accuracy, completeness, merchantability, fitness for a particular purpose, or non-infringement. Descriptions of, or references to, products or publications, or inclusion of an advertisement in this article, neither constitutes nor implies a guarantee, endorsement, or support of claims made of that product, publication, or service.

Copyright © 2013, INFORMS

Please scroll down for article—it is on subsequent pages



INFORMS is the largest professional society in the world for professionals in the fields of operations research, management science, and analytics.

For more information on INFORMS, its publications, membership, or meetings visit <http://www.informs.org>

# A Nonparametric Approach to Modeling Choice with Limited Data

Vivek F. Farias

Sloan School of Management and Operations Research Center, Massachusetts Institute of Technology,  
Cambridge, Massachusetts 02142, vivekf@mit.edu

Srikanth Jagabathula

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

Devavrat Shah

Electrical Engineering and Computer Science, Massachusetts Institute of Technology,  
Cambridge, Massachusetts 02139, devavrat@mit.edu

Choice models today are ubiquitous across a range of applications in operations and marketing. Real-world implementations of many of these models face the formidable stumbling block of simply identifying the “right” model of choice to use. Because models of choice are inherently high-dimensional objects, the typical approach to dealing with this problem is positing, a priori, a parametric model that one believes adequately captures choice behavior. This approach can be substantially suboptimal in scenarios where one cares about using the choice model learned to make fine-grained predictions; one must contend with the risks of mis-specification and overfitting/underfitting. Thus motivated, we visit the following problem: For a “generic” model of consumer choice (namely, distributions over preference lists) and a limited amount of data on how consumers actually make decisions (such as marginal information about these distributions), how may one predict revenues from offering a particular assortment of choices? An outcome of our investigation is a *nonparametric* approach in which the data automatically select the right choice model for revenue predictions. The approach is practical. Using a data set consisting of automobile sales transaction data from a major U.S. automaker, our method demonstrates a 20% improvement in prediction accuracy over state-of-the-art benchmark models; this improvement can translate into a 10% increase in revenues from optimizing the offer set. We also address a number of theoretical issues, among them a qualitative examination of the choice models implicitly learned by the approach. We believe that this paper takes a step toward “automating” the crucial task of choice model selection.

*Key words:* nonparametric choice; choice models; revenue prediction; utility preference; preference list; marketing mix

*History:* Received December 14, 2009; accepted December 21, 2011, by Yossi Aviv, operations management.  
Published online in *Articles in Advance* January 8, 2013.

## 1. Introduction

A problem of central interest to operations managers is using historical sales data to predict the revenues or sales from offering a particular assortment of products to customers. As one can imagine, such predictions form crucial inputs to several important business decisions, both operational and otherwise. A classical example of such a decision problem is that of assortment planning: deciding the “optimal” assortment of products to offer customers with a view to maximizing expected revenues (or some related objective) subject to various constraints (e.g., limited display or shelf space). A number of variants of this problem, both static and dynamic, arise in essentially every facet of revenue management. Such problems are seen as crucial revenue management tasks and needless to say, accurate revenue or sales predictions fundamentally impact how well we can perform such tasks.

Why might these crucial predictions be difficult to make? Consider the task of predicting expected sales rates from offering a particular set of products to customers. In industry jargon, this is referred to as the “conversion rate,” and is defined as the probability of converting an arriving customer into a purchasing customer. Predicting the conversion rate for an offer set is difficult because the probability of purchase of each product depends on all the products on offer. This is due to substitution behavior, where an arriving customer potentially substitutes an unavailable product with an available one. Because of substitution, the sales observed for a product may be viewed as a combination of its “primary” demand and additional demand. *Customer choice models* have been used to model this behavior with success. At an abstract level, a choice model can be thought of as a conditional probability distribution that for any offer set yields the probability that an arriving customer purchases a given product in that set.

There is *vast* literature spanning marketing, economics, and psychology devoted to the construction of parametric choice models and their estimation from data. In the literature that studies the sorts of revenue management decision problems we alluded to above, such models are typically assumed *given*. The implicit understanding is that a complete prescription for these decision problems will require fitting the “right” parametric choice model to data, so as to make accurate revenue or sales predictions. This is a complex task. Apart from the fact that one can never be sure that the chosen parametric structure is a “good” representation of the underlying ground truth, parametric models are prone to overfitting and underfitting issues. Once a structure is *fixed*, one does not glean new structural information from data. This is a serious issue in practice because although a simple model (such as the multinomial logit (MNL) model) may make practically unreasonable assumptions (such as the so-called “IIA” (independent of irrelevant alternatives) assumption), fitting a more complex model can lead to worse performance because of overfitting—and one can never be sure.

In this paper, we propose a *nonparametric, data-driven* approach to making revenue or sales predictions that afford the revenue manager the opportunity to avoid the challenging task of fitting an appropriate parametric choice model to historical data. Our approach views choice models *generically*, namely, as distributions over rankings (or preference lists) of products. As will be seen subsequently, this view subsumes essentially all extant choice models. Further, this view yields a *nonparametric* approach to choice modeling where the revenue manager does not need to think about the appropriate parametric structure for his problem, or the trade-off between model parsimony and the risk of overfitting. Rather, through the use of a nonparametric approach, our goal is to offload as much of this burden as possible to the data itself.

### 1.1. Contributions

As previously mentioned, we consider entirely generic models of choice, specified as a distribution over all possible rankings (or preference lists) of products. Our view of data is aligned with what one typically has available in reality, sales rates of products in an assortment, for some set of product assortments. This is a general view of choice modeling. Our main contribution is to make this view operational, yielding a data-driven, nonparametric approach. Specifically, we make the following contributions in the context of this general setup:

- *Revenue Predictions.* Accurate revenue or sales predictions form core inputs for a number of important revenue/inventory management problems.

Available sales data will typically be insufficient to fully specify a generic model of choice of the type we consider. We therefore seek to identify the *set* of generic choice models consistent with available sales data. Given the need to make a revenue or sales prediction on a heretofore unseen assortment, we then offer the *worst-case* expected revenue possible for that assortment assuming that the true model lies in the set of models found to be consistent with observed sales data. Such an approach makes no a priori structural assumptions on the choice model, and has the appealing feature that as more data become available, the predictions will improve, by narrowing down the set of consistent models. This simple philosophy dictates challenging computational problems; for instance, the sets we compute are computationally unwieldy and, at first glance, highly intractable. Nonetheless, we successfully develop several simple algorithms of increasing sophistication to address these problems.

- *Empirical Evaluation.* We conducted an empirical study to gauge the practical value of our approach, both in terms of the absolute quality of the predictions produced, and also relative to using alternative parametric approaches. We describe the results of two such studies:

- (i) *Simulation study.* The purpose of our simulation study is to demonstrate that the robust approach can effectively capture model structure consistent with a number of different parametric models and produce good revenue predictions. The general setup in this study is as follows: We use a parametric model to generate synthetic transaction data. We then use these data in conjunction with our revenue prediction procedure to predict expected revenues over a swathe of offer sets. Our experimental design permits us to compare these predictions to the corresponding “ground truth.” The parametric families we considered included the MNL, nested logit (NL), and a mixture of multinomial logit (MMNL) models. To “stress-test” our approach, we conducted experiments over a wide range of parameter regimes for these generative parametric choice models, including some that were fit to DVD sales data from Amazon.com. The predictions produced are remarkably accurate.

- (ii) *Empirical study with sales data from a major U.S. automaker.* The purpose of our empirical study is twofold: (1) to demonstrate how our setup can be applied with real-world data, and (2) to pit the robust method in a “horse race” against the MNL and MMNL parametric families of models. For the case study, we used sales data collected daily at the dealership level over 2009 to 2010 for a range of small SUVs offered by a major U.S. automaker for a dealership zone in the Midwest. We used a portion of these sales data as “training” data. We made these

data available to our robust approach, as well as in the fitting of an MNL model and an MMNL model. We tested the quality of conversion-rate predictions (i.e., a prediction of the sales rate given the assortment of models on the lot) using the robust approach and the incumbent parametric approaches on the remainder of the data. We conducted a series of experiments by varying the amount of training data made available to the approaches. We conclude that (a) the robust method improves on the accuracy of either of the parametric methods by about 20% (this is large) in all cases, and (b) unlike the parametric models, the robust method is apparently not susceptible to underfitting and overfitting issues. In fact, we see that the performance of the MMNL model relative to the MNL model deteriorates as the amount of training data available decreases because of overfitting. Improved forecast accuracy improves the decisions made. For instance, a 20% improvement in forecast accuracy can result in a 10% increase in revenues from optimizing the offer set.

• *Descriptive Analysis.* In making revenue predictions, we did not need to concern ourselves with the choice model implicitly assumed by our prediction procedure. However, it is natural to consider criteria for selecting choice models consistent with observed data that are independent of any decision context. Thus motivated, we consider the natural task of finding the *simplest* choice model consistent with the observed data. As in much of contemporary high-dimensional statistics, we employ *sparsity*<sup>1</sup> as our measure of simplicity. First, we use the sparsest fit criterion to obtain a characterization of the choice models implicitly used by the robust revenue prediction approach. Loosely speaking, we show that the choice model implicitly used by the robust approach is essentially the sparsest model (Theorem 1), and the complexity of the model (as measured by its sparsity) scales with the “amount” of data. This provides an explanation for the immunity of the robust approach to overfitting/underfitting as observed in our case study. Second, we characterize the family of choice models that can be identified only from observed marginal data via the sparsest fit criterion (Theorems 2 and 3). Our characterization formalizes the notion that the complexity of the models that can be identified via the sparsest fit criterion scales with the amount of data at hand.

## 1.2. Relevant Literature

The study of choice models and their applications spans a vast literature across multiple fields including at least marketing, operations, and economics.

<sup>1</sup> By sparsity we refer to the number of rank lists or, in effect, customer types, assumed to occur with positive probability in the population.

In disciplines such as marketing, learning a choice model is an interesting goal unto itself given that it is frequently the case that a researcher wishes to uncover “why” a particular decision was made. Within operations, the goal is frequently more application oriented with the choice model being explicitly used as a predictive tool within some larger decision model. Because our goals are aligned with the latter direction, our literature review focuses predominantly on operations management (OM); we briefly touch on key work in marketing. We note that our consideration of sparsity as an appropriate nonparametric model selection criterion is closely related to the burgeoning statistical area of compressive sensing; we discuss those connections in §6.

The vast majority of decision models encountered in operations have traditionally ignored substitution behavior (and thereby choice modeling) altogether. Within airline revenue management (RM), this is referred to as the “independent demand” model (see Talluri and van Ryzin 2004b). Over the years, several studies have demonstrated the improvements that could be obtained by incorporating choice behavior into operations models. For example, within airline RM, the simulation studies conducted by Belobaba and Hopperstad (1999) on the well-known passenger origin and destination simulator (PODS) suggested the value of corrections to the independent demand model; more recently, Ratliff et al. (2008a) and Vulcano et al. (2010) demonstrated valuable average revenue improvements from using MNL choice-based RM approaches using real airline market data. Following such studies, there has been a significant amount of research in the areas of inventory management and RM attempting to incorporate choice behavior into operations models.

The bulk of the research on choice modeling in both the areas has been optimization related. That is to say, most of the work has focused on devising optimal decisions *given* a choice model. Talluri and van Ryzin (2004a), Gallego et al. (2004), van Ryzin and Vulcano (2008), Mahajan and van Ryzin (1999), and Goyal et al. (2009) are all papers in this vein. Kök et al. (2009) provide an excellent overview of the state-of-the-art in assortment optimization. Rusmevichientong et al. (2010) consider the MNL model and provide an efficient algorithm for the static assortment optimization problem and propose an efficient policy for the dynamic optimization problem. A follow-up paper, Rusmevichientong and Topaloglu (2012), considers the same optimization problem but where the mean utilities in the MNL model are allowed to lie in some arbitrary uncertainty set. Saure and Zeevi (2009) propose an alternative approach for the dynamic assortment optimization problem under a general random utility model.

The majority of the above-mentioned work focuses on optimization issues given a choice model. Papers such as Talluri and van Ryzin (2004a) discuss optimization problems with *general* choice models, and, as such, our revenue estimation procedure fits in perfectly there. In most cases, however, the choice model is assumed to be given and of the MNL type. Papers such as Saure and Zeevi (2009) and Rusmevichientong and Topaloglu (2012) loosen this requirement by allowing some amount of *parametric* uncertainty. In particular, Saure and Zeevi (2009) assume unknown mean utilities and learn these utilities, whereas the optimization schemes in Rusmevichientong and Topaloglu (2012) require knowledge of mean utilities only within an interval. In both cases, the structure of the model (effectively, MNL) is *fixed* up front.

The MNL model is by far the most popular choice model studied and applied in OM. The origins of the MNL model date all the way back to the Plackett–Luce model, proposed independently by Luce (1959) and Plackett (1975). Before becoming popular in the area of OM, the MNL model found widespread use in the areas of transportation (see seminal works of McFadden 1980 and Ben-Akiva and Lerman 1985) and marketing (starting with the seminal work of Guadagni and Little 1983, which paved the way for choice modeling using scanner panel data). See Wierenga (2008) and Chandukala et al. (2008) for a detailed overview of choice modeling in the area of marketing. The MNL model is popular because its structure makes it tractable, both in terms of estimating its parameters and solving decision problems. However, the tractability of the MNL model comes at a cost: It is incapable of capturing any heterogeneity in substitution patterns across products (see Debreu 1960) and suffers from independent of irrelevant alternatives property (see Ben-Akiva and Lerman 1985), both of which limit its practical applicability.

Of course, these issues with the MNL model are well recognized, and far more sophisticated models of choice have been suggested in the literature (see, for instance, Ben-Akiva and Lerman 1985, Anderson et al. 1992); the price one pays is that the more sophisticated models may not be easily identified from sales data and are prone to overfitting. It must be noted that an exception to this state of affairs is the paper by Rusmevichientong et al. (2006), which considers a general nonparametric model of choice similar to the one considered here in the context of an assortment pricing problem. The caveat is that the approach considered requires access to samples of entire customer preference lists that are unlikely to be available in many practical applications.

Our goal relative to all of the above-mentioned work is to *eliminate* the need for structural assumptions and thereby the associated risks as well.

We provide a means of going directly from raw sales transaction data to revenue or sales estimates for a given offer set. Although this does not represent the entirety of what can be done with a choice model, it represents a valuable application, at least within the operational problems discussed.

## 2. The Choice Model and Problem Formulations

We consider a universe of  $N$  products,  $\mathcal{N} = \{0, 1, 2, \dots, N - 1\}$ . We assume that the 0th product in  $\mathcal{N}$  corresponds to the “outside” or “no-purchase” option. A customer is associated with a permutation (or ranking)  $\sigma$  of the products in  $\mathcal{N}$ ; the customer prefers product  $i$  to product  $j$  if and only if  $\sigma(i) < \sigma(j)$ . A customer will be presented with a set of alternatives  $\mathcal{M} \subset \mathcal{N}$ ; any set of alternatives will, by convention, be understood to include the no-purchase alternative, i.e., the 0th product. The customer will subsequently choose to purchase her single most preferred product among those in  $\mathcal{M}$ . In particular, she purchases

$$\arg \min_{i \in \mathcal{M}} \sigma(i). \tag{1}$$

It is quickly seen that the above structural assumption is consistent with structural assumptions made in commonly encountered choice models including the multinomial logit, nested multinomial logit, or more general random utility models. Those models make many additional structural assumptions, which may or may not be reasonable for the application at hand. Viewed in a different light, basic results from the theory of social preferences dictate that the structural assumptions implicit in our model are no more restrictive than assuming that the customer in question is endowed with a utility function over alternatives and chooses an alternative that maximizes her utility from among those available. Our model of the customer is thus general.<sup>2</sup>

### 2.1. Choice Model

To make useful predictions on customer behavior that might, for instance, guide the selection of a set  $\mathcal{M}$  to offer for sale, one must specify a choice model. A general choice model is effectively a conditional probability distribution  $\mathbb{P}(\cdot | \cdot): \mathcal{N} \times 2^{\mathcal{N}} \rightarrow [0, 1]$  that yields the probability of purchase of a particular product in  $\mathcal{N}$  given the set of alternatives available to the customer.

We assume essentially the most general model for  $\mathbb{P}(\cdot | \cdot)$ . In particular, we assume that there exists a distribution  $\lambda: S_N \rightarrow [0, 1]$  over the set of all possible permutations  $S_N$ . Recall here that  $S_N$  is effectively the set

<sup>2</sup> As opposed to associating a customer with a fixed  $\sigma$ , one may also associate customers with distributions over permutations. This latter formalism is superfluous for our purposes.

of all possible customer types because every customer is associated with a permutation that uniquely determines her choice behavior. The distribution  $\lambda$  defines our choice model as follows: Define the set

$$\mathcal{S}_j(\mathcal{M}) = \{\sigma \in S_N: \sigma(j) < \sigma(i), \forall i \in \mathcal{M}, i \neq j\},$$

where  $\mathcal{S}_j(\mathcal{M})$  is simply the set of all customer types that would purchase product  $j$  when the offer set is  $\mathcal{M}$ . Our choice model is then given by

$$\mathbb{P}(j | \mathcal{M}) = \sum_{\sigma \in \mathcal{S}_j(\mathcal{M})} \lambda(\sigma) \triangleq \lambda^j(\mathcal{M}).$$

Not surprisingly, as previously mentioned, the above model subsumes essentially any model of choice one might concoct: In particular, all we have assumed is that at a *given* point in time a customer possess *rational* (transitive) (see Mas-Colell et al. 1995) preferences over all alternatives,<sup>3</sup> and that a particular customer will purchase her most preferred product from the offered set according to these preferences; a given customer sampled at different times may well have a distinct set of preferences.

## 2.2. Data

The class of choice models we work with is quite general and imposes a minimal number of behavioral assumptions on customers a priori. That said, the data available to calibrate such a model will typically be limited in the sense that a modeler will have sales rate information for a potentially small collection of assortments. Ignoring the difficulties of such a calibration problem for now, we posit a general notion of what we mean by observable data. The abstract notion we posit will quickly be seen as relevant to data one might obtain from sales information.

We assume that the data observed by the seller are given by an  $m$ -dimensional “partial information” vector  $y = A\lambda$ , where  $A \in \{0, 1\}^{m \times N!}$  makes precise the relationship between the observed data and the underlying choice model. Typically, we anticipate  $m \ll N!$  signifying, for example, the fact that we have sales information for only a limited number of assortments.

We now show how the type of data available in practice can be cast in the form of  $y = A\lambda$ . In the retail context, historical data about customer purchase behavior are available in the form of observed sales transactions from a set of displayed assortments. In particular, one typically has information about observed sales for a sequence of test assortments say  $\mathcal{M}_1, \mathcal{M}_2, \dots, \mathcal{M}_L$ . For each of the assortments  $\mathcal{M}_l$  and

products  $i$ , the sales data provide the fraction of purchasing customers who purchased product  $i$  when the displayed assortment was  $\mathcal{M}_l$ . Given these data, we claim that they can be written as a linear combination of  $\lambda$  for an appropriate choice of matrix  $A$ .

To see this, it is instructive to start with a simple special case that we call “comparison data.” Specifically, suppose that we have access to data that provide us with the information about the fraction of customers that prefer product  $i$  to product  $j$ , for all pairs of products  $i$  and  $j$ . For this case, the partial information vector  $y$  may be indexed by  $i, j$  with  $0 \leq i, j \leq N - 1; i \neq j$ . For each  $i, j$ ,  $y_{ij}$  denotes the fraction of customers that prefer product  $i$  to  $j$ . The matrix  $A$  is thus in  $\{0, 1\}^{N(N-1) \times N!}$ . A column of  $A$ ,  $A(\sigma)$ , will thus have  $A(\sigma)_{ij} = 1$  if and only if  $\sigma(i) < \sigma(j)$ . It is important to note here that we have introduced the comparison data for the simplicity of exposition and (as will become apparent later) for theoretical considerations. The way we have defined it, comparison data is in fact *not* readily available in practice: For starters, there is typically censoring because of which we can only observe the fraction of customers who prefer  $i$  to both  $j$  and 0 when we offer the pair of products  $i, j$ . In addition, practical applications do not typically provide sales information about all the  $\binom{N}{2}$  possible pairs of products. Nevertheless, comparison data provide a simple yet nontrivial example that makes our setup concrete.

More realistically, we have sales transaction data from a set of displayed assortments  $\mathcal{M}_1, \mathcal{M}_2, \dots, \mathcal{M}_L$ . In this case, denoting by  $y_{il}$  the fraction of customers purchasing product  $i$  when assortment  $\mathcal{M}_l$  is on offer, our partial information vector,  $y \in [0, 1]^{N \cdot L}$ , may thus be indexed by  $i, l$  with  $0 \leq i \leq N - 1, 1 \leq l \leq L$ . The matrix  $A$  is then in  $\{0, 1\}^{N \cdot L \times N!}$ . For a column of  $A$  corresponding to the permutation  $\sigma$ ,  $A(\sigma)$ , we will then have  $A(\sigma)_{il} = 1$  iff  $i \in \mathcal{M}_l$  and  $\sigma(i) < \sigma(j)$  for all products  $j$  in assortment  $\mathcal{M}_l \cup \{0\}$ .

Finally, we emphasize that the idea of viewing partial information as  $y = A\lambda$  is a very powerful one. It captures several different types of interesting partial information in addition to the transactional data as previously described. This becomes important both from a theoretical standpoint and from the standpoint of other applications. Although we will not explore any other application contexts, we discuss two other types of partial information, the “ranking data” and “top-set data,” for our theoretical analysis in §6.

## 2.3. Incorporating Choice in Decision Models: A Revenue Estimation Black Box

Although modeling choice is useful for a variety of reasons, we are largely motivated by decision models for OM problems that benefit from the incorporation of a choice model. In many of these models, the fundamental feature impacted by the choice

<sup>3</sup> Note, however, that the customer need not be aware of these preferences; from (1), it is evident that the customer need only be aware of his preferences for elements of the offer set.

model is a “revenue function” that measures revenue rates corresponding to a particular assortment of products offered to customers. Concrete examples include static assortment management, network revenue management under choice, and inventory management assuming substitution.

We formalize this revenue function. We associate every product in  $\mathcal{N}$  with a retail price  $p_j$ . Of course,  $p_0 = 0$ . The revenue function,  $R(\mathcal{M})$ , determines expected revenues to a retailer from offering a set of products  $\mathcal{M}$  to his customers. Under our choice model this is given by

$$R(\mathcal{M}) = \sum_{j \in \mathcal{M}} p_j \lambda^j(\mathcal{M}).$$

The function  $R(\cdot)$  is a fundamental building block for all of the OM problems previously described, so that we view the problem of estimating  $R(\cdot)$  as our central motivating problem. The above specification is general, and we refer to *any* linear functional of the type above as a revenue function. As another useful example of such a functional, consider setting  $p_j = 1$  for all  $j > 0$  (i.e., all products other than the no-purchase option). In this case, the revenue function  $R(\mathcal{M})$  yields the probability an arriving customer will purchase some product in  $\mathcal{M}$ ; i.e., the conversion rate under assortment  $\mathcal{M}$ .

Given a “black box” that is capable of producing estimates of  $R(\cdot)$  using some limited corpus of data, one may then hope to use such a black box for making assortment decisions over time in the context of the OM problems of the type discussed in the introduction.

### 2.4. Problem Formulations

Imagine we have a corpus of transaction data, summarized by an appropriate data vector  $y$  as described in §2.2. Our goal is to use *just* these data to make predictions about the revenue rate (i.e., the expected revenues garnered from a random customer) for some given assortment, say  $\mathcal{M}$ , that has never been encountered in past data. We propose accomplishing this by solving the following program:

$$\begin{aligned} & \underset{\lambda}{\text{minimize}} && R(\mathcal{M}) \\ & \text{subject to} && A\lambda = y, \\ & && \mathbf{1}^\top \lambda = 1, \\ & && \lambda \geq 0. \end{aligned} \tag{2}$$

In particular, the optimal value of this program will constitute our prediction for the revenue rate. In words, the feasible region of this program describes the set of all choice models consistent with the observed data  $y$ . The optimal objective value consequently corresponds to the *minimum* revenues possible for the assortment  $\mathcal{M}$  under any choice model

consistent with the observed data. Because the family of choice models we considered was *generic* this prediction relies on simply the data and basic economic assumptions on the customer that are tacitly assumed in essentially any choice model.

The philosophy underlying the above program can be put to other uses. For instance, one might seek to recover a choice model itself from the available data. In a parametric world, one would consider a suitably small, fixed family of models within which a unique model would best explain (but not necessarily be consistent with) the available data. It is highly unlikely that available data will determine a unique model in the *general* family of models we consider here. Our nonparametric setting thus requires an appropriate selection criterion. A natural criterion is to seek the “simplest” choice model that is consistent with the observed data. There are many notions of what one might consider simple. One criterion that enjoys widespread use in high-dimensional statistics is sparsity. In particular, we may consider finding a choice model  $\lambda$  consistent with the observed data, that has minimal support,  $\|\lambda\|_0 \triangleq |\{\lambda(\sigma) : \lambda(\sigma) \neq 0\}|$ , where  $|S|$  denotes the cardinality of set  $S$ . In other words, we might seek to explain observed purchasing behavior by presuming as small a number of modes of customer choice behavior as possible (where we associate a “mode” of choice with a ranking of products). More formally, we might seek to solve

$$\begin{aligned} & \underset{\lambda}{\text{minimize}} && \|\lambda\|_0 \\ & \text{subject to} && A\lambda = y, \\ & && \mathbf{1}^\top \lambda = 1, \\ & && \lambda \geq 0. \end{aligned} \tag{3}$$

Sections 3, 4, and 5 are focused on providing procedures to solve the program (2) and on examining the quality of the predictions produced on simulated data and actual transaction data, respectively. Section 6 will discuss algorithmic and interesting descriptive issues pertaining to (3).

### 3. Revenue Predictions: Computation

In the previous section, we formulated the task of computing revenue predictions via a nonparametric model of choice and any available data as the mathematical program (2), which we repeat below, in a slightly different form for clarity:

$$\begin{aligned} & \underset{\lambda}{\text{minimize}} && \sum_{j \in \mathcal{M}} p_j \lambda_j(\mathcal{M}) \\ & \text{subject to} && A\lambda = y, \\ & && \mathbf{1}^\top \lambda = 1, \\ & && \lambda \geq 0. \end{aligned}$$

This mathematical program is a linear program (LP) in the variables  $\lambda$ . Interpreting the program in words, the constraints  $A\lambda = y$  ensure that any  $\lambda$  assumed in making a revenue estimate is *consistent* with the observed data. Other than this consistency requirement, writing the probability that a customer purchases  $j \in \mathcal{M}$ ,  $\mathbb{P}(j | \mathcal{M})$ , as the quantity  $\lambda_j(\mathcal{M}) \triangleq \sum_{\sigma \in \mathcal{S}_j(\mathcal{M})} \lambda(\sigma)$  assumes that the choice model satisfies the basic structure laid out in §2.1. We make no other assumptions outside of these, and ask for the lowest expected revenues possible for  $\mathcal{M}$  under *any* choice model satisfying these requirements.

Thus, while the assumptions implicit in making a revenue estimate are something that the user need not think about, the two natural questions that arise are the following:

1. How does one solve this conceptually simple program in practice given that the program involves an intractable number of variables?
2. Even if one did succeed in solving such a program, are the revenue predictions produced useful or are they too loose to be of practical value?

This section focuses on the first question. In practical applications, such a procedure would need to be integrated into a larger decision problem; therefore, it is useful to understand the computational details that we present at a high level in this section. The second “so what” question will be the subject of §§4 and 5, where we examine the performance of the scheme on simulated transaction data and finally on a real-world sales prediction problem using real data. Finally, in §6, we examine an interesting property enjoyed by the choice models implicitly assumed in making the predictions in this scheme.

### 3.1. The Dual to the Robust Problem

At a high level our approach to solving (2) will be to consider the dual of that program and then derive efficient exact or approximate descriptions to the feasible regions of these programs. We begin by considering the dual program to (2). In preparation for taking the dual, let us define

$$\mathcal{A}_j(\mathcal{M}) \triangleq \{A(\sigma) : \sigma \in \mathcal{S}_j(\mathcal{M})\},$$

where we recall that  $\mathcal{S}_j(\mathcal{M}) = \{\sigma \in S_N : \sigma(j) < \sigma(i), \forall i \in \mathcal{M}, i \neq j\}$  denotes the set of all permutations that result in the purchase of  $j \in \mathcal{M}$  when the offered assortment is  $\mathcal{M}$ . Because  $S_N = \bigcup_{j \in \mathcal{M}} \mathcal{S}_j(\mathcal{M})$  and  $\mathcal{S}_j(\mathcal{M}) \cap \mathcal{S}_i(\mathcal{M}) = \emptyset$  for  $i \neq j$ , we have implicitly specified a partition of the columns of the matrix  $A$ . Armed with this notation, the dual of (2) is

$$\begin{aligned} & \text{maximize}_{\alpha, \nu} \quad (\alpha^\top y + \nu) \\ & \text{subject to} \quad \max_{x^j \in \mathcal{A}_j(\mathcal{M})} (\alpha^\top x^j + \nu) \leq p_j, \quad \text{for each } j \in \mathcal{M}, \end{aligned} \quad (4)$$

where  $\alpha$  and  $\nu$  are dual variables corresponding, respectively, to the data consistency constraints  $A\lambda = y$  and the requirement that  $\lambda$  is a probability distribution (i.e.,  $\mathbf{1}^\top \lambda = 1$ ), respectively. Of course, this program has a potentially intractable number of constraints. We explore two approaches to solving the dual:

1. An extremely simple to implement approach that relies on sampling constraints in the dual that will, in general, produce approximate solutions that are upper bounds to the optimal solution of our robust estimation problem.
2. An approach that relies on producing effective representations of the sets  $\mathcal{A}_j(\mathcal{M})$ , so that each of the constraints  $\max_{x^j \in \mathcal{A}_j(\mathcal{M})} (\alpha^\top x^j + \nu) \leq p_j$  can be expressed efficiently. This approach is slightly more complex to implement, but in return can be used to sequentially produce tighter approximations to the robust estimation problem. In certain special cases, this approach is provably efficient and optimal.

### 3.2. The First Approach: Constraint Sampling

The following is an extremely simple to implement approach to approximately solve the problem (4):

1. Select a distribution over permutations,  $\psi$ .
2. Sample  $n$  permutations according to the distribution. Call this set of permutations  $\hat{\mathcal{S}}$ .
3. Solve the program:

$$\begin{aligned} & \text{maximize}_{\alpha, \nu} \quad (\alpha^\top y + \nu) \\ & \text{subject to} \quad \alpha^\top A(\sigma) + \nu \leq p_j, \quad \text{for each } j \in \mathcal{M}, \sigma \in \hat{\mathcal{S}}. \end{aligned} \quad (5)$$

Observe that (5) is essentially a “sampled” version of the problem (4), wherein constraints of that problem have been sampled according to the distribution  $\psi$  and are consequently a relaxation of that problem. A solution to (5) is consequently an upper bound to the optimal solution to (4).

The question of whether the solutions thus obtained provide meaningful approximations to (4) is partially addressed by recent theory developed by Calafiore and Campi (2005). In particular, it has been shown that for a problem with  $m$  variables and given  $n = O((1/\epsilon)(m \ln(1/\epsilon) + \ln(1/\delta)))$  samples, we must have that with probability at least  $1 - \delta$  the following holds: An optimal solution to (5) violates at most an  $\epsilon$  fraction of constraints of the problem (4) under the measure  $\psi$ . Hence, given a number of samples that scales only with the number of variables (and is independent of the number of constraints in (4)), one can produce a solution to (4) that satisfies all but a small fraction of constraints. The theory does not provide any guarantees on how far the optimal cost of the relaxed problem is from the optimal cost of the original problem.



The heuristic nature of this approach notwithstanding, it is extremely simple to implement, and in the experiments conducted in the next section, provided close to optimal solutions.

### 3.3. The Second Approach: Efficient Representations of $\mathcal{A}_j(\mathcal{M})$

We describe here one notion of an efficient representation of the sets  $\mathcal{A}_j(\mathcal{M})$ , and assuming we have such a representation, we describe how one may solve (4) efficiently. We deal with the issue of actually coming up with these efficient representations in Online Appendix B (online appendices available at [http://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=2187779](http://papers.ssrn.com/sol3/papers.cfm?abstract_id=2187779)), where we develop an efficient representation for ranking data and demonstrate a generic procedure to sequentially produce such representations.

Let us assume that every set  $\mathcal{S}_j(\mathcal{M})$  can be expressed as a disjoint union of  $D_j$  sets. We denote the  $d$ th such set by  $\mathcal{S}_{jd}(\mathcal{M})$  and let  $\mathcal{A}_{jd}(\mathcal{M})$  be the corresponding set of columns of  $A$ . Consider the convex hull of the set  $\mathcal{A}_{jd}(\mathcal{M})$ ,  $\text{conv}\{\mathcal{A}_{jd}(\mathcal{M})\} \triangleq \bar{\mathcal{A}}_{jd}(\mathcal{M})$ . Recalling that  $A \in \{0, 1\}^{m \times N}$ ,  $\mathcal{A}_{jd}(\mathcal{M}) \subset \{0, 1\}^m$ .  $\bar{\mathcal{A}}_{jd}(\mathcal{M})$  is thus a polytope contained in the  $m$ -dimensional unit cube,  $[0, 1]^m$ . In other words,

$$\bar{\mathcal{A}}_{jd}(\mathcal{M}) = \{x^{jd}: A_1^{jd} x^{jd} \geq b_1^{jd}, A_2^{jd} x^{jd} = b_2^{jd}, A_3^{jd} x^{jd} \leq b_3^{jd}, x^{jd} \in \mathbb{R}_+^m\} \quad (6)$$

for some matrices  $A^{jd}$  and vectors  $b^{jd}$ . By a canonical representation of  $\mathcal{A}_j(\mathcal{M})$ , we thus understand a partition of  $\mathcal{S}_j(\mathcal{M})$  and a polyhedral representation of the columns corresponding to every set in the partition as given by (6). If the number of partitions as well as the polyhedral description of each set of the partition given by (6) is polynomial in the input size, we will regard the canonical representation as efficient. Of course, there is no guarantee that an efficient representation of this type exists; clearly, this must rely on the nature of our partial information, i.e., the structure of the matrix  $A$ . Even if an efficient representation did exist, it remains unclear whether we can identify it. Ignoring these issues for now, in the remainder of this section, we demonstrate how given a representation of the type (6), one may solve (4) with the time complexity that is polynomial in the size of the representation.

For simplicity of notation, in what follows we assume that each polytope  $\bar{\mathcal{A}}_{jd}(\mathcal{M})$  is in standard form,

$$\bar{\mathcal{A}}_{jd}(\mathcal{M}) = \{x^{jd}: A^{jd} x^{jd} = b^{jd}, x^{jd} \geq 0\}.$$

Now because an affine function is always optimized at the vertices of a polytope, we know that

$$\max_{x^{jd} \in \bar{\mathcal{A}}_j(\mathcal{M})} (\alpha^\top x^{jd} + \nu) = \max_{d, x^{jd} \in \bar{\mathcal{A}}_{jd}(\mathcal{M})} (\alpha^\top x^{jd} + \nu).$$

We have thus reduced (4) to a “robust” LP. Now, by strong duality, we have the following:

$$\begin{aligned} & \underset{x^{jd}}{\text{maximize}} \quad \alpha^\top x^{jd} + \nu & \underset{\gamma^{jd}}{\text{minimize}} \quad (b^{jd})^\top \gamma^{jd} + \nu \\ & \text{subject to} \quad A^{jd} x^{jd} = b^{jd} \equiv & \text{subject to} \quad (\gamma^{jd})^\top A^{jd} \geq \alpha. \\ & \quad \quad \quad x^{jd} \geq 0. & \end{aligned} \quad (7)$$

We have thus established the following useful equality:

$$\begin{aligned} & \left\{ \alpha, \nu: \max_{x^j \in \bar{\mathcal{A}}_j(\mathcal{M})} (\alpha^\top x^j + \nu) \leq p_j \right\} \\ & = \{ \alpha, \nu: (b^{jd})^\top \gamma^{jd} + \nu \leq p_j, (\gamma^{jd})^\top A^{jd} \geq \alpha, d = 1, 2, \dots, D_j \}. \end{aligned}$$

It follows that solving (2) is equivalent to the following LP whose complexity is polynomial in the description of our canonical representation:

$$\begin{aligned} & \underset{\alpha, \nu}{\text{maximize}} \quad \alpha^\top y + \nu \\ & \text{subject to} \quad (b^{jd})^\top \gamma^{jd} + \nu \leq p_j \\ & \quad \quad \quad \text{for all } j \in \mathcal{M}, d = 1, 2, \dots, D_j \quad (8) \\ & \quad \quad \quad (\gamma^{jd})^\top A^{jd} \geq \alpha \\ & \quad \quad \quad \text{for all } j \in \mathcal{M}, d = 1, 2, \dots, D_j. \end{aligned}$$

As discussed, our ability to solve (8) relies on our ability to produce an efficient canonical representation of  $\mathcal{S}_j(\mathcal{M})$  of the type (6). In Online Appendix B, we first consider the case of ranking data, where such an efficient representation may be produced. We then illustrate a method that produces a sequence of “outer approximations” to (6) for general types of data, and thereby allows us to produce a sequence of improving lower bounding approximations to our robust revenue estimation problem, (2). This provides a general procedure to address the task of solving (4) or, equivalently, (2).

We end this section with a brief note on noise. So far we have assumed that the choice probabilities  $y_i$  obtained from historical data are known exactly and fit the model exactly so that there exists a choice model  $\lambda$  such that  $y = A\lambda$ . This is, of course, hardly the case in practice. Specifically, there are two sources of errors. First is the finite sample error caused due to the fact that the choice probabilities  $y_i$  can only be estimated through a sample average of finitely many samples; depending on the number of samples available, there is uncertainty in the estimate of  $y_i$ . Second, there could be model misfit errors. That is, even if the choice probabilities  $y_i$  are known exactly, our choice model may not be an exact fit, making the set of equalities  $y = A\lambda$  and  $\lambda$  a distribution infeasible. To overcome these issues, we incorporate these errors

in making our predictions. Specifically, we assume that we are given an uncertainty region  $\mathcal{E}$  constructed from the data such that there exists a choice model  $\lambda$  with  $A\lambda \in \mathcal{E}$ . The uncertainty region  $\mathcal{E}$  may either be an uncertainty ellipsoid or a “box” derived from sample averages of the associated choice probabilities and the corresponding confidence intervals. Given such an uncertainty region  $\mathcal{E}$ , we predict revenues by solving the following LP:

$$\begin{aligned} & \underset{\lambda, y \in \mathcal{E}}{\text{minimize}} && \sum_{j \in \mathcal{M}} p_j \lambda_j(\mathcal{M}) \\ & \text{subject to} && A\lambda = y, \\ & && \mathbf{1}^\top \lambda = 1, \\ & && \lambda \geq 0. \end{aligned}$$

Provided  $\mathcal{E}$  is convex, this program is essentially no harder to solve than the variant of the problem we have discussed, and similar methods to those developed in this section apply. It is clear from the convex program that if the uncertainty region  $\mathcal{E}$  is “too small,” then the constraint set will become infeasible. On the other hand, if  $\mathcal{E}$  is “too large,” then we would have a poor fit to the data and conservative revenue predictions. To balance the extremes, in our empirical analyses, we choose the “smallest” uncertainty region such that the constraint set becomes feasible. Precise details of how we do that are provided in §5.

#### 4. Revenue Predictions: Data-Driven Computational Study

In this section, we describe the results of an extensive simulation study, the main purpose of which is to demonstrate that the robust approach can capture various underlying parametric structures and produce good revenue predictions. For this study, we pick a range of random utility parametric structures used extensively in current modeling practice.

The broad experimental procedure we followed is the following:

1. Pick a structural model. This may be a model derived from real-world data or a purely synthetic model.
2. Use this structural model to simulate sales for a set of test assortments. This simulates a data set that a practitioner likely has access to.
3. Use this transaction data to estimate marginal information  $y$ , and use  $y$  to implement the robust approach.
4. Use the implemented robust approach to predict revenues for a distinct set of assortments, and compare the predictions to the *true* revenues computed using the ground-truth structural model chosen for benchmarking in step 1.

Note that the above experimental procedure lets us isolate the impact of structural errors from that of finite sample errors. Specifically, our goal is to understand how well the robust approach captures the underlying choice structure. For this purpose, we ignore any estimation errors in data by using the ground-truth parametric model to compute the *exact* values of any choice probabilities and revenues required for comparison. Therefore, if the robust approach has good performance across an interesting spectrum of structural models that are believed to be good fits to data observed in practice, we can conclude that the robust approach is likely to offer accurate revenue predictions with no additional information about structure across a wide-range of problems encountered in practice.

##### 4.1. Benchmark Models and Nature of Synthetic Data

The above procedure generates data sets using a variety of ground-truth structural models. We pick the following “random utility” models as benchmarks. A self-contained and compact exposition on the foundations of each of the benchmark models described next can be found in the online appendix.

*Multinomial Logit Family (MNL)*. For this family, we have

$$\mathbb{P}(j | \mathcal{M}) = w_j / \sum_{i \in \mathcal{M}} w_i,$$

where the  $w_i$  are the parameters specifying the models. See Online Appendix C.1 for more details.

*Nested Logit Family (NL)*. This model is a first attempt at overcoming the independence of irrelevant alternatives effect, a shortcoming of the MNL model. For this family, the universe of products is partitioned into  $L$  mutually exclusive subsets, or “nests,” denoted by  $\mathcal{N}_1, \mathcal{N}_2, \dots, \mathcal{N}_L$  such that

$$\mathcal{N} = \bigcup_{l=1}^L \mathcal{N}_l \quad \text{and} \quad \mathcal{N}_l \cap \mathcal{N}_m = \emptyset, \quad \text{for } m \neq l.$$

This model takes the form

$$\begin{aligned} \mathbb{P}(j | \mathcal{M}) &= \mathbb{P}(\mathcal{N}_l | \mathcal{M}) \mathbb{P}(j | \mathcal{N}_l, \mathcal{M}) \\ &= \frac{(w(l, \mathcal{M}))^\rho}{\sum_{m=1}^L (w(m, \mathcal{M}))^\rho} \frac{w_j}{w(l, \mathcal{M})}, \end{aligned} \tag{9}$$

where  $\rho < 1$  is a certain scale parameter, and

$$w(l, \mathcal{M}) \stackrel{\text{def}}{=} \alpha_l w_0 + \sum_{i \in (\mathcal{N}_l \cap \mathcal{M}) \setminus \{0\}} w_i.$$

Here,  $\alpha_l$  is the parameter capturing the level of membership of the no-purchase option in nest  $l$  and satisfies,  $\sum_{l=1}^L \alpha_l^\rho = 1$ ,  $\alpha_l \geq 0$ , for  $l = 1, 2, \dots, L$ . In cases when  $\alpha_l < 1$  for all  $l$ , the family is called the *cross nested logit (CNL)* family. For a more detailed

description, including the corresponding random utility function and bibliographic details, see Online Appendix C.2.

*Mixed Multinomial Logit Family (MMNL).* This model accounts specifically for customer heterogeneity. In its most common form, the model reduces to

$$\mathbb{P}(j | \mathcal{M}) = \int \frac{\exp(\beta^T x_j)}{\sum_{i \in \mathcal{M}} \exp(\beta^T x_i)} G(d\beta; \theta),$$

where  $x_j$  is a vector of observed attributes for the  $j$ th product, and  $G(\cdot, \theta)$  is a distribution parameterized by  $\theta$  selected by the econometrician that describes heterogeneity in taste. Because the coefficients  $\beta$  are assumed to be random (unlike for the MNL model), this model is often also termed in the literature the random coefficients multinomial logit (RC-MNL) model. For this paper, we restrict ourselves to Gaussian MMNL models in which we assume that  $\beta$  has a multivariate Gaussian distribution. For a more detailed description, including the corresponding random utility function and bibliographic details, see Online Appendix C.3.

**4.1.1. Transaction Data Generated.** Having selected (and specified) a structural model from the above list, we generated sales transactions as follows:

1. Fix an assortment of two products,  $i, j$ .
2. Compute the values of  $P(i | \{i, j, 0\})$ ,  $P(j | \{i, j, 0\})$  using the chosen parametric model.
3. Repeat the above procedure for all pairs,  $\{i, j\}$ , and single item sets,  $\{i\}$ .

The above data are succinctly summarized as an  $N^2 - N$  dimensional data vector  $y$ , where  $y_{i,j} = P(i | \{i, j, 0\})$  for  $0 \leq i, j \leq N - 1$ ,  $i \neq j$ . Given the above data, the precise specialization of the robust estimation problem (2) that we solve can be found in Online Appendix B.3.

**4.2. Experiments Conducted**

With the above setup we conducted two broad sets of experiments. In the first set of experiments, we picked specific models from the MNL, CNL, and MMNL model classes; the MNL model was constructed using DVD shopping cart data from Amazon.com, and the CNL and MMNL models were obtained through slight “perturbations” of the MNL model. To avoid any artifacts associated with specific models, in the second set of experiments, we conducted stress tests by generating a number of instances of models from each of the MNL, CNL, and MMNL model classes. We next present the details of the two sets of experiments.

**4.2.1. The Amazon Model.** We considered an MNL model fit to Amazon.com DVD sales data

collected between July 1, 2005, to September 30, 2005,<sup>4</sup> where an individual customer’s utility for a given DVD,  $j$  is given by

$$U_j = \theta_0 + \theta_1 x_{j,1} + \theta_2 x_{j,2} + \xi_j.^5$$

Here,  $x_{j,1}$  is the price of the package  $j$  divided by the number of physical discs it contains, and  $x_{j,2}$  is the total number of helpful votes received by product  $j$  and  $\xi_j$  is a standard Gumbel. The model fit to the data has  $\theta_0 = -4.31$ ,  $\theta_1 = -0.038$  and  $\theta_2 = 3.54 \times 10^{-5}$ . See Table 2 in the online appendix for the attribute values taken by the 15 products we used for our experiments. We abbreviate this model AMZN for future reference.

We also considered the following synthetic perturbations of the AMZN model:

1. *AMZN-CNL.* We derived a CNL model from the original AMZN model by partitioning the products into four nests with the first nest containing products 1 to 5, the second nest containing products 6 to 9, the third nest containing products 10 to 13, and the last nest containing products 14 and 15. We chose  $\rho = 0.5$ . We assigned the no-purchase option to every nest with nest membership parameter  $\alpha_i = (1/4)^{1/\rho} = 1/16$ .

2. *AMZN-MMNL.* We derived an MMNL model from the original AMZN model by replacing each  $\theta_i$  parameter with the random quantity  $\beta_i = (1 + \eta_{i,j})\theta_i$ , for  $i = 0, 1, 2$  with  $\eta_{i,j}$  a customer specific random variable distributed as a zero mean normal random variable with standard deviation  $s = 0.25$ . Put differently, we assumed that the random coefficients  $\beta_i$  are independent and have a Gaussian distribution with mean  $\theta_i$  and standard deviation  $s\theta_i$ .

Figure 1 shows the results of the generic experiment for each of the three models. Each experiment queries the robust estimate on sixty randomly drawn assortments of sizes between one and seven and compares these estimates to those under the respective true model for each case.

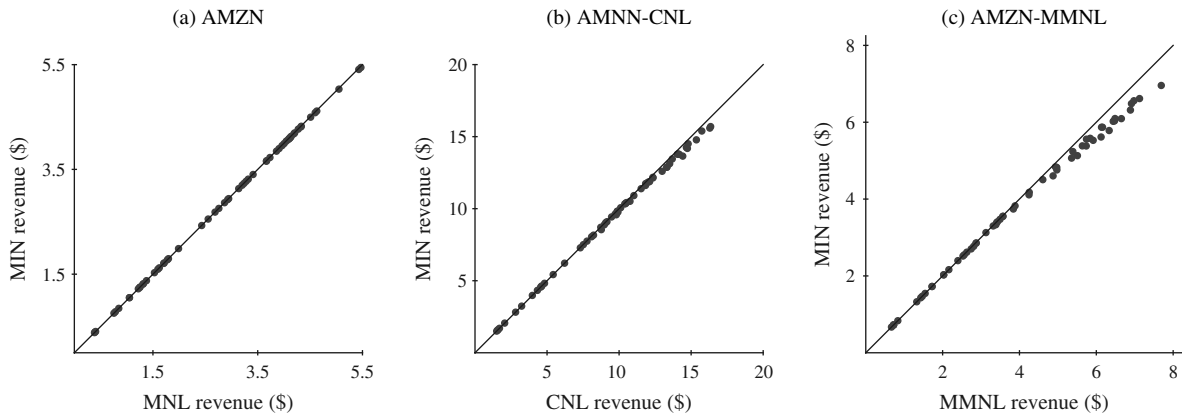
**4.2.2. Synthetic Model Experiments.** The above experiments considered structurally diverse models, each for a *specific* set of parameters. Are the conclusions suggested by Figure 1 artifacts of the set of parameters? To assuage this concern, we performed stress tests by considering each structural model in turn, and for each model generating a number of instances of the model by drawing the relevant parameters from a generative family. For each structural model, we considered the following generative families of parameters:

1. *MNL Random Family.* Twenty randomly generated models on 15 products, each generated by drawing mean utilities,  $\ln w_j$ , uniformly between  $-5$  and  $5$ .

<sup>4</sup> The specifics of this model were shared with us by the authors of Rusmevichientong et al. (2010).

<sup>5</sup> The corresponding weights  $w_j$  are given by  $w_j = \exp(\theta_0 + \theta_1 x_{j,1} + \theta_2 x_{j,2})$ .

Figure 1 Robust Revenue Predictions (MIN) vs. True Revenues for the AMZN, AMZN-CNL, and AMZN-MMNL Models



Note. Each of the 60 points in a plot represents the coordinate (true revenue, MIN revenue) for a randomly drawn assortment.

2. *CNL Random Family.* We maintained the nests, selection of  $\rho$  and  $\alpha_j$  as in the AMZN-CNL model. We generated 20 distinct CNL models, each generated by drawing  $\ln w_j$  uniformly between  $-5$  and  $5$ .

3. *MMNL Random Family.* We preserved the basic nature of the AMZN-MMNL model. We considered 20 randomly generated MMNL models. Each model differs in the distribution of the parameter vector  $\beta$ . The random coefficients  $\beta_j$  in each case are defined as follows:  $\beta_j = (1 + \eta_{i,j})\theta_j$ , where  $\eta_{i,j}$  is a  $N(\mu_j, 0.25)$  random variable. Each of the 20 models corresponds to a single draw of  $\mu_j$  for  $j = 0, 1, 2$  from the uniform distribution on  $[-1, 1]$ .

For each of the 60 structural model instances previously described, we randomly generated 20 offer sets of sizes between one and seven. For a given offer set  $\mathcal{M}$ , we queried the robust procedure and compared the revenue estimate produced to the true revenue for that offer set; we can compute the latter quantity theoretically. In particular, we measured the relative error,  $\varepsilon(\mathcal{M}) \stackrel{\text{def}}{=} (R^{\text{true}}(\mathcal{M}) - R^{\text{MIN}}(\mathcal{M})) / R^{\text{MIN}}(\mathcal{M})$ . Figure 2 represents distributions of relative error for the three generative families previously described. Each histogram consists of 400 test points; a given test point

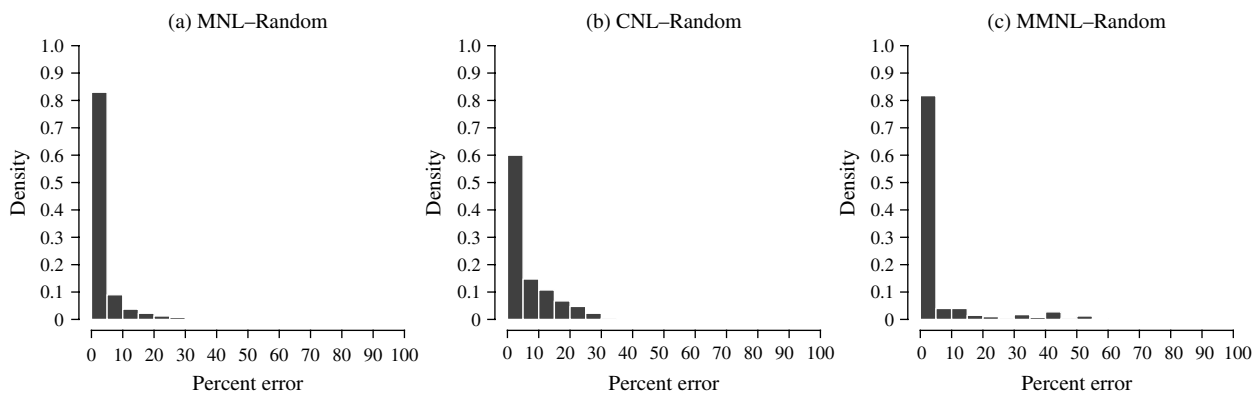
corresponds to one of the 20 randomly generated structural models in the relevant family, and a random assortment.

In the above stress tests, we kept the standard deviation of the coefficients  $\beta_i$  in the MMNL models as  $s\theta_i$  with the multiplier  $s$  fixed at 0.25. The standard deviation of the coefficients in the MMNL model can be treated as a measure of the heterogeneity or the “complexity” of the model. Naturally, if we keep the amount of transaction data fixed and increase the standard deviation—and hence the complexity of the underlying model—we expect the accuracy of robust estimates to deteriorate. To give a sense of the sensitivity of the accuracy of robust revenue predictions to changes in the standard deviation of coefficients, we repeated the stress tests with the MMNL model class for three different values of the multiplier  $s$ : 0.1, 0.25, and 0.4. Figure 3 shows the comparison of the density plots of relative errors for the three cases.

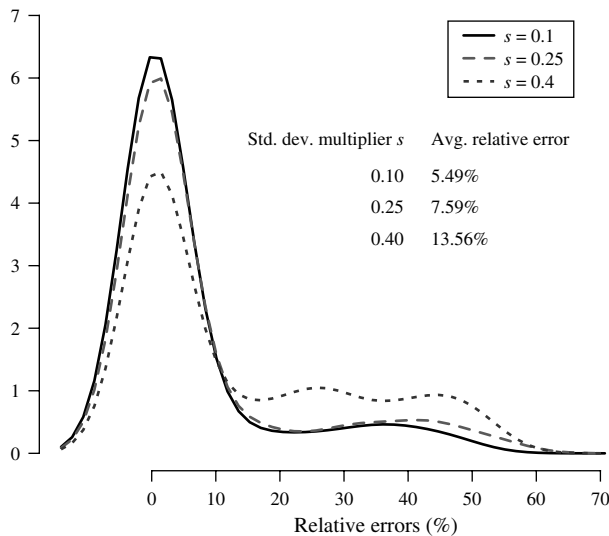
We draw the following broad conclusion from the above experiments:

1. Given limited marginal information for distributions over permutations  $\lambda$  arising from a number

Figure 2 Relative Error Across Multiple Instances of the MNL, CNL, and MMNL Structural Models



**Figure 3 Accuracy of Robust Revenue Predictions Deteriorates with Increase in Model Complexity, Measured in Terms of the Standard Deviation  $s\theta_i$  of the Normally Distributed Coefficients  $\beta_i$  in the MMNL Model**



*Notes.* The standard deviation multiplier takes three different values: 0.1, 0.25, and 0.4. The densities were estimated through kernel density estimation. The density estimates go below zero as a result of smoothing.

of commonly used structural models of choice, the robust approach effectively captures diverse parametric structures and provides close revenue predictions under a range of practically relevant parametric models.

2. With the type of marginal information  $y$  fixed, the accuracy of robust revenue predictions deteriorates (albeit mildly) as the complexity of the underlying model increases; this is evidenced by the deterioration of robust performance as we go from the MNL to the MMNL model class, and similarly as we increase the standard deviation of the coefficients for the MMNL model while keeping the amount of data fixed.

3. The design of our experiments allows us to conclude that in the event that a given structural model among the types used in our experiments predicts revenue rates accurately, the robust approach is likely to be just as good *without* the knowledge of the relevant structure. In the event that the structural model used is a poor fit, the robust approach will continue to provide meaningful guarantees on revenues under the mild condition that it is tested in an environment where the distribution generating sales is no different from the distribution used to collect marginal information.

## 5. Revenue Predictions: Case Study with a Major U.S. Automaker

In this section, we present the results of a case study conducted using sales transaction data from the dealer network of a major U.S. automaker. Our goal in this study is to use historical transaction data to

predict the sales rate or conversion rate for any given offer set of automobiles on a dealer lot. This conversion rate is defined as the probability of converting an arriving customer into a purchasing customer. The purpose of the case study is twofold: (1) to demonstrate how the prediction methods developed in this paper can be applied in the real world and the quality of the predictions they offer in an absolute sense, and (2) to pit the robust method for revenue predictions in a horse race against parametric approaches based on the MNL and MMNL families of choice models. To test the performance of these approaches in different regimes of calibration data, we carried out cross-validations with varying amounts of training/calibration data. The results of the experiments conducted as part of the case study provide us with the evidence to draw two main conclusions:

1. The robust method predicts conversion rates more accurately than either of the parametric methods. In our case study, the improvement in accuracy was about 20% across all regimes of calibration data.

2. Unlike the parametric methods we study, the robust approach is apparently *not* susceptible to overfitting and underfitting.

The 20% improvement in accuracy is substantial. The second conclusion has important implications as well: In practice, it is often difficult to ascertain whether the data available are “sufficient” to fit the model at hand. As a result, parametric structures are prone to overfitting or underfitting. The robust approach, on the other hand, *automatically* scales the complexity of the underlying model class with data available, so in principle one should be able to avoid these issues. This is borne out by the case study. In the remainder of this section, we describe the experimental setup and then present the evidence to support the above conclusions.

### 5.1. Setup

We collect data comprising purchase transactions of a specific range of small SUVs offered by a major U.S. automaker over 16 months. The data are collected at the dealership level (i.e., the finest level possible) for a network of dealers in the Midwest. Each transaction contains information about the date of sale, the identity of the SUV sold, and the identity of the other SUVs on the dealership lot at the time of sale. Here, by “identity” we mean a unique model identifier that collectively identifies a package of features, color, and invoice price point. We make the assumption that purchase behavior within the zone can be described by a single distribution over preference lists. To ensure the validity of this assumption, we restrict attention to a specific dealership zone, defined as the collection of dealerships within an appropriately defined geographical area with relatively homogeneous demographic features. Our data consisted of

sales information on 14 distinct SUV identities (as previously described), which we term products.

**5.1.1. Data.** To describe the data and our methods precisely, we introduce some notation. We let  $\mathcal{M}^{\text{training}}$  and  $\mathcal{M}^{\text{test}}$ , respectively, denote the set of assortments used as part of training and test data. For any given assortment,  $\mathcal{M}$  and product  $i \in \mathcal{M}$ , define  $C_{i,\mathcal{M}}$  as the total number of sales of product  $i$  across all dealerships over the data collection period, such that the assortment on offer at the time of sale was  $\mathcal{M}$ . Similarly,  $C_{0,\mathcal{M}}$  denotes the number of customers who purchased nothing when  $\mathcal{M}$  was on offer. Note that we do not have access to this latter quantity; we describe how it is estimated momentarily. Finally, let  $T^{\text{training}}$  denote the set of tuples  $(i, \mathcal{M})$  such that  $\mathcal{M} \in \mathcal{M}^{\text{training}}$ . We observed a total of  $M \triangleq |\mathcal{M}^{\text{training}}| + |\mathcal{M}^{\text{test}}| = 203$  distinct assortments (or subsets) of the 14 products in the data set, where each assortment  $\mathcal{M}_i$ ,  $i = 1, 2, \dots, M$ , was on offer at some point at some dealership in the dealership zone.

**5.1.2. Demand Untruncation.** As previously discussed,  $C_{0,\mathcal{M}}$  is unavailable because we do not observe arriving customers that do not purchase. This issue impacts choice modeling irrespective of whether one chooses the nonparametric approach adopted here of any of the extant parametric approaches. We follow a strategy that is common in practice when one has access to the rich data we do here. In more data-limited scenarios, more sophisticated techniques can be applied; see, for instance, Talluri and van Ryzin (2004a), Vulcano et al. (2010, 2012) and Ratliff et al. (2008b). Importantly, our estimates of  $C_{0,\mathcal{M}}$  will be common to our robust revenue prediction approach and the incumbent parametric approaches we study. Central to the above task is estimating the number of customers that considered assortment  $\mathcal{M}$ . In particular, given this estimate, we can compute  $C_{0,\mathcal{M}}$  as

$$C_{0,\mathcal{M}} = (\text{number of customer arrivals when } \mathcal{M} \text{ is on offer}) - \sum_{j \in \mathcal{M}} C_{j,\mathcal{M}}.$$

For a given dealership this is the number of arriving customers over days when  $\mathcal{M}$  was on offer at that dealership; we then simply sum this figure over all dealerships. In particular,

$$\begin{aligned} & \text{number of customer arrivals when } \mathcal{M} \text{ was on offer} \\ &= \sum_d \alpha_d \text{ days}_d(\mathcal{M}), \end{aligned} \quad (10)$$

where  $\alpha_d$  denotes the average number of customers arriving daily at dealership  $d$ , and  $\text{days}_d(\mathcal{M})$  denotes the number of days for which  $\mathcal{M}$  was on offer at dealership  $d$ . There are a number of ways of estimating  $\alpha_d$ ;

in fact, this is the focus of the untruncation literature alluded to previously. As mentioned, we adopt a strategy that is common in practice when one has access to historical sales at a retailer: Assume that  $\alpha_d$  is proportional to total sales at the dealer over the year preceding the year in which the data was collected. We estimate the proportionality constant using cross-validation on our training data. It is worth noting that this scheme of estimating  $\alpha_d$  implicitly assumes a coarse relationship between arrivals and sales in the preceding year. This relationship is assumed *only* for the purposes of estimating  $\alpha_d$ .

**5.1.3. Robust Method.** Given  $\mathcal{M}^{\text{training}}$  and an assortment  $\mathcal{M} \in \mathcal{M}^{\text{test}}$ , the conversion rate of  $\mathcal{M}$  is predicted by the robust approach by solving the following LP:

$$\begin{aligned} & \text{minimize} \quad \sum_{j \in \mathcal{M}} \mathbb{P}_\lambda(j | \mathcal{M}) \\ & \text{subject to} \quad a_{i,\mathcal{M}} \leq \mathbb{P}_\lambda(i | \mathcal{M}) \leq b_{i,\mathcal{M}}, \\ & \quad \quad \quad \forall (i, \mathcal{M}) \in T^{\text{training}} \quad (11) \\ & \quad \quad \quad \mathbf{1}^\top \lambda = 1, \\ & \quad \quad \quad \lambda \geq 0, \end{aligned}$$

where we recall that  $\mathbb{P}_\lambda(i | \mathcal{M}) = \sum_{\sigma \in \mathcal{S}_i(\mathcal{M})} \lambda(\sigma)$  with  $\mathcal{S}_i(\mathcal{M})$ , denoting the set  $\{\sigma: \sigma(i) < \sigma(j) \forall j \in \mathcal{M}, i \neq j\}$ , and  $[a_{i,\mathcal{M}}, b_{i,\mathcal{M}}]$  denotes the interval to which  $\mathbb{P}_\lambda(i | \mathcal{M})$  belongs. We obtained an approximate solution to the LP in (11) by taking its dual and using the approach of constraint sampling as described in §3.2. The LP (11) is a slight modification of the LP in (2) in that the prices  $p_i$  are all set to 1 and the equalities  $y_t = \mathbb{P}_\lambda(\mathcal{M})$  are changed to inequalities to account for finite sample errors in the data. Setting all the prices to 1 has the effect of computing the conversion rate for the assortment. For each tuple  $(i, \mathcal{M}) \in T^{\text{training}}$ , we computed the left and right end points, respectively, as  $a_{i,\mathcal{M}} = y_{i,\mathcal{M}}(1 - z\varepsilon_{i,\mathcal{M}})$  and  $b_{i,\mathcal{M}} = y_{i,\mathcal{M}}(1 + z\varepsilon_{i,\mathcal{M}})$ , where

$$\varepsilon_{i,\mathcal{M}} = \sqrt{\frac{1 - y_{i,\mathcal{M}}}{C_{i,\mathcal{M}}}} \quad \text{and} \quad y_{i,\mathcal{M}} = \frac{C_{i,\mathcal{M}}}{C_{0,\mathcal{M}} + \sum_{i \in \mathcal{M}} C_{i,\mathcal{M}}}.$$

Here,  $\hat{y}_{i,\mathcal{M}}\varepsilon_{i,\mathcal{M}}$  is the standard error, and  $z$  is a constant multiplier that determines the width of the confidence interval. Different values of  $z$  give us approximate confidence intervals for  $\mathbb{P}_\lambda(i | \mathcal{M})$ . For our experiments, we had set  $z$  to be 3.15, which corresponded to the smallest value of  $z$  for which (11) was feasible; incidentally, this value of  $z$  also corresponds to approximate 99.8% confidence interval for  $\mathbb{P}_\lambda(i | \mathcal{M})$ .

**5.1.4. Parametric Methods.** As benchmarks, we fit an MNL as well as an MMNL model given the  $C_{i,\mathcal{M}}$  and  $C_{0,\mathcal{M}}$  estimates previously described. For the MNL model, we assumed the following specific random

utility structure  $U_i = V_i + \xi_i$ ,  $i = 0, 1, 2, \dots, N$ , where  $V_i$  is the mean utility and  $\xi_i$  are independent and identically distributed Gumbel distributed with location parameter 0 and scale parameter 1, and  $N = 14$  is the number of products. For the MMNL model we assumed the following specific random utility structure:  $U_i = V_i + \beta x_i + \xi_i$ ,  $i = 0, 1, 2, \dots, 14$ , where as before  $V_i$  denotes the mean utility and  $N = 14$  the number of products,  $\xi_i$  are independent and identically distributed Gumbel with location parameter 0 and scale parameter 1,  $x_i$  are dummy features with  $x_0 = 0$  and  $x_i = 1$  for  $i > 0$ , and  $\beta$  is Gaussian with mean 0 and variance  $s^2$ .

For both models, we used the training data  $\hat{y}_{i,\mathcal{M}}$ , for all  $(i, \mathcal{M}) \in T^{\text{training}}$  to determine the maximum-likelihood (ML) estimates of the parameters. Specifically, fixing  $V_0$  to 0, we used BIOGEME (Bierlaire 2003, 2008) to estimate  $V_i$ ,  $i > 0$ , and  $s$ .

### 5.2. Experiments and Results

We now describe the experiments we conducted and present the results we obtained. To test the predictive performance of the robust, the MNL, and the MMNL methods, we carried out  $k$ -fold cross-validations with  $k = 2, 5, 10$ . In  $k$ -fold cross-validation (see Mosteller and Tukey 1968), we arbitrarily partition the collection of assortments  $\mathcal{M}_1, \mathcal{M}_2, \dots, \mathcal{M}_M$  into  $k$  partitions of about equal size, except maybe the last partition. Then, using  $k - 1$  partitions as training data to calibrate the methods, we test their performance on the  $k$ th partition. We repeat this process  $k$  times with each of the  $k$  partitions used as test data exactly once. This repetition ensures that each assortment is tested at least once. Note that as  $k$  decreases, the number of training assortments decreases resulting in more limited data scenarios. Such limited data scenarios are of course of great practical interest.

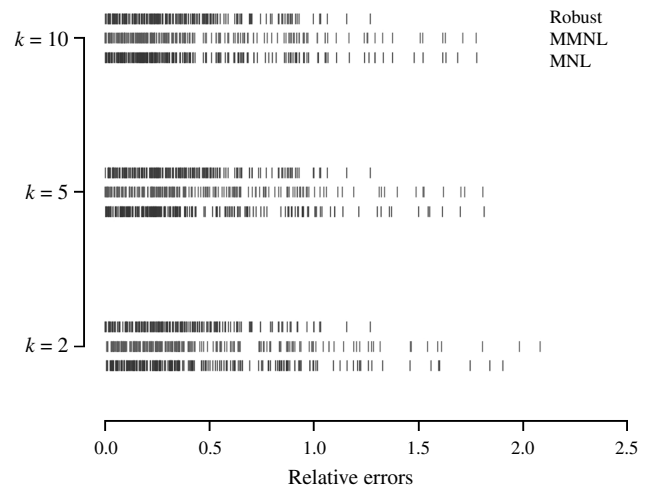
We measure the prediction accuracy of the methods using the relative error metric. In particular, letting  $\hat{y}(\mathcal{M})$  denote the conversion-rate prediction for test assortment  $\mathcal{M}$ , the incurred relative error is defined as  $|\hat{y}(\mathcal{M}) - y(\mathcal{M})|/y(\mathcal{M})$ , where

$$y(\mathcal{M}) := (\text{number of customers who purchase product when } \mathcal{M} \text{ is on offer}) \cdot (\text{number of customer arrivals when } \mathcal{M} \text{ was on offer})^{-1}.$$

In the case of the parametric approaches,  $\hat{y}(\mathcal{M})$  is computed using the choice model fit to the training data. In the case of the robust approach, we solve an appropriate mathematical program. A detailed description of how  $\hat{y}(\mathcal{M})$  is determined by each method is provided in the online appendix.

We now present the results of the experiments. Figure 4 shows the comparison of the relative errors of the three methods from  $k$ -fold cross-validations for  $k = 10, 5, 2$ . Table 1 shows the mean relative error

**Figure 4 Robust Method Outperforms Both MNL and MMNL Methods in Conversion-Rate Predictions Across Various Calibration Data Regimes**



*Notes.* This figure compares relative errors of the three methods in  $k$ -fold cross-validations for  $k = 10, 5, 2$ . Each point corresponds to the relative error for a particular test assortment.

percentages of the three methods and the percentage of improvement in mean relative error achieved by the robust method over the MNL and MMNL methods for the three calibration data regimes of  $k = 10, 5, 2$ . It is clear from the definition of  $k$ -fold cross-validation that as  $k$  decreases, the amount of calibration data decreases, or equivalently calibration data sparsity increases. Such sparse calibration data regimes are of course of great practical interest.

The immediate conclusion we draw from the results is that the prediction accuracy of the robust method is better than those of both MNL and MMNL methods in all calibration data regimes. In particular, using the robust method results in close to 20% improvement in prediction accuracy over the MNL and MMNL methods. We also note that although the prediction accuracy of the more complex MMNL method is marginally better than that of the MNL method in the high calibration-data regime of  $k = 10$ , it quickly becomes worse as the amount of calibration data available decreases. This behavior is a consequence of overfitting caused by the complexity of the MMNL model. The performance of the robust method, on the other hand, remains stable across the different regimes of calibration data.

**Table 1 Mean Relative Errors in Percentages of Different Methods**

$k$	MNL	MMNL	Robust	Percentage of improvement over	
				MNL	MMNL
10	43.43	43.39	34.79	19.89	19.80
5	43.25	45.73	35.79	17.23	21.62
2	45.65	46.61	36.83	19.33	20.99

## 6. The Sparsest Choice Model Consistent with Data

In making revenue predictions, we did not need to concern ourselves with the choice model implicitly assumed by our prediction procedure. However, it is natural to consider criteria for selecting choice models consistent with the observed data that are independent of any decision context. Thus motivated, we consider the natural task of finding the *simplest* choice model consistent with the observed data. As in much of contemporary high-dimensional statistics (see, for example, Candes et al. 2006, Cormode and Muthukrishnan 2006), we employ sparsity as our measure of simplicity. Sparse models use as few preference lists as possible to explain observed substitutions and have provided a great deal of tractability in multiple applications (see, for example, van Ryzin and Vulcano 2008). Our goal in this section is to first understand the choice models implicitly assumed by the robust procedure through the lens of the sparsity criterion, and second, to understand the discriminative power of this criterion.

Toward this goal, we begin by characterizing choice models implicitly used by the robust approach in terms of their sparsity. Loosely speaking, we establish that the choice model implicitly used by the robust approach is indeed simple or sparse. In particular, such choice models have sparsity within at most one of the sparsity of the sparsest model consistent with the data. As such, we see that the choice model implicitly selected by our robust revenue prediction procedure is, in essence, the sparsest choice model consistent with the data. From a descriptive perspective, this establishes the appealing fact that simplicity or sparsity is a natural property possessed by all choice models used in making robust revenue predictions. We also establish that the sparsity of the choice model used by the robust approach scales with the dimension of the data vector  $y$  thereby establishing that the complexity of the model used by the robust approach scales with the amount of data available. This provides a potential explanation for the immunity of the robust approach to overfitting/underfitting issues, as evidenced in our case study.

Next, we turn to understanding the discriminative power of the sparsest fit criterion. Toward this end, we describe a family of choice models that can be uniquely identified from the given marginal data using the sparsest fit criterion. We intuitively expect the complexity of identifiable models to scale with the amount of data that is available. We formalize this intuition by presenting for various types of data, conditions on the model generating the data under which identification is possible. These conditions characterize families of choice models that can be identified in terms of their sparsity and formalize the scaling

between the complexity of a model class and the amount of data needed to identify it.

### 6.1. Revenue Prediction and Sparse Models

We now provide a characterization of the choice models implicitly used by the robust procedure through the lens of model sparsity. As previously mentioned, loosely speaking, we can establish that the choice models selected implicitly via our revenue estimation procedure are, in essence, close to the sparsest model consistent with the observed data. In other words, the robust approach implicitly uses the simplest models consistent with observed data to predict revenues.

To state our result formally, let us define the set  $\mathcal{Y}$  as the set of all possible data vectors, namely, the convex hull of the columns of the matrix  $A$ . For some  $y \in \mathcal{Y}$  and an arbitrary offer set,  $\mathcal{M}$ , let  $\lambda^{\min}(y)$  be an optimal *basic feasible* solution to the program used in our revenue estimation procedure, namely, (2). Moreover, let,  $\lambda^{\text{sparse}}(y)$  be the *sparsest* choice model consistent with the data vector  $y$ ; i.e.,  $\lambda^{\text{sparse}}(y)$  is an optimal solution to (3). We then have that with probability one, the sparsity (i.e., the number of rank lists with positive mass) under  $\lambda^{\min}(y)$  is close to that of  $\lambda^{\text{sparse}}(y)$ . In particular, we have the following theorem:

**THEOREM 1.** *For any distribution over  $\mathcal{Y}$  that is absolutely continuous with respect to Lebesgue measure on  $\mathcal{Y}$ , we have with probability 1, that*

$$0 \leq \|\lambda^{\min}(y)\|_0 - \|\lambda^{\text{sparse}}(y)\|_0 \leq 1.$$

Theorem 1 establishes that if  $K$  were the support size of the sparsest distribution consistent with  $y$ , the sparsity of the choice model used by our revenue estimation procedure is either  $K$  or  $K + 1$  for “almost all” data vectors  $y$ . As such, this establishes that the choice model implicitly employed by the robust procedure is essentially also the sparsest model consistent with the observed data.

In addition, the proof of the theorem reveals that the sparsity of the robust choice model consistent with the observed data is either<sup>6</sup>  $m$  or  $m + 1$  for almost all data vectors  $y$  of dimension  $m$ . This yields yet another valuable insight into the choice models implicit in our revenue predictions—the complexity of these models, as measured by their sparsity, grows with the amount of observed data. As such, we see that the complexity of the choice model implicitly employed by the robust procedure scales automatically with the amount of available data, as one would desire from a nonparametric scheme. This provides a potential explanation for the robust procedures’ lack of susceptibility to the overfitting observed for the MMNL model in our empirical study.

<sup>6</sup> Here, we assume that matrix  $A$  has full row rank.



### 6.2. Identifiable Families of Choice Models

We now consider the family of choice models that can be identified via the sparsest fit criterion. For that, we present two abstract conditions that, if satisfied by the choice model generating the data  $y$ , guarantee that the optimal solution to (3) is unique and, in fact, equal to the choice model generating the data.

Before we describe the conditions, we introduce some notation. As before, let  $\lambda$  denote the true underlying distribution, and let  $K$  denote the support size,  $\|\lambda\|_0$ . Let  $\sigma_1, \sigma_2, \dots, \sigma_K$  denote the permutations in the support, i.e.,  $\lambda(\sigma_i) \neq 0$  for  $1 \leq i \leq K$ , and  $\lambda(\sigma) = 0$  for all  $\sigma \neq \sigma_i, 1 \leq i \leq K$ . Recall that  $y$  is of dimension  $m$ , and we index its elements by  $d$ . The two conditions are the following:

*Signature Condition.* For every permutation  $\sigma_i$  in the support, there exists a  $d(i) \in \{1, 2, \dots, m\}$  such that  $A(\sigma_i)_{d(i)} = 1$  and  $A(\sigma_j)_{d(i)} = 0$ , for every  $j \neq i$  and  $1 \leq j \leq K$ . In other words, for each permutation  $\sigma_i$  in the support,  $y_{d(i)}$  serves as its “signature.”

*Linear Independence Condition.* For any  $c_i \in \mathbb{Z}$  (the set of integers) and  $|c_i| \leq C$ , where  $C$  is a sufficiently large number greater than or equal to  $K$ , we have  $\sum_{i=1}^K c_i \lambda(\sigma_i) \neq 0$ . This condition is satisfied with probability 1 if  $[\lambda_1 \lambda_2 \dots \lambda_K]^T$  is drawn uniformly from the  $K$ -dim simplex or, for that matter, any distribution on the  $K$ -dim simplex with a density.

When these two conditions are satisfied by a choice model, this choice model can be recovered from observed data as the solution to problem (3). Specifically, we have the following theorem:

**THEOREM 2.** *Suppose we are given  $y = A\lambda$ , and  $\lambda$  satisfies the signature and linear independence conditions. Then,  $\lambda$  is the unique solution to the program in (3).*

The proof of Theorem 2 is provided in Online Appendix A.2. The proof is constructive in that it describes an efficient scheme to determine the underlying choice model. Thus, the theorem establishes that whenever the underlying choice model satisfies the signature and linear independence conditions, it can be identified using an efficient scheme as the optimal solution to the program in (3). We next characterize a family of choice models that satisfy the signature and linear independence conditions. Specifically, we show that *essentially all* choice models with sparsity  $K(N)$  satisfy these two conditions as long as  $K(N)$  scales as  $\log N$ ,  $\sqrt{N}$ , and  $N$  for comparison data, top-set data, and ranking data, respectively. To capture this notion of “essentially” all choice models, we introduce a natural generative model. We discuss how restrictive the above values of  $K(N)$  are at the end of the section.

*A Generative Model.* Given  $K$  and an interval  $[a, b]$  on the positive real line, we generate a choice model  $\lambda$  as follows: Choose  $K$  permutations,  $\sigma_1, \sigma_2, \dots, \sigma_K$ ,

uniformly at random with replacement,<sup>7</sup> choose  $K$  numbers uniformly at random from the interval  $[a, b]$ , normalize the numbers so that they sum to 1,<sup>8</sup> and assign them to the permutations  $\sigma_i, 1 \leq i \leq K$ . For all other permutations  $\sigma \neq \sigma_i, \lambda(\sigma) = 0$ .

Depending on the observed data, we characterize values of sparsity  $K = K(N)$  up to which distributions generated by the above generative model can be recovered with a high probability. We derive the sparsity bound for three different types of partial information: comparison data, top-set data, and ranking data. We introduced comparison data in §2. We define ranking data and top-set data as follows:

- *Ranking Data.* These data represent the fraction of customers that rank a given product  $i$  as their  $r$ th choice. Here the partial information vector  $y$  is indexed by  $i, r$  with  $0 \leq i, r \leq N$ . For each  $i, r, y_{ri}$  is thus the fraction of customers that rank product  $i$  at position  $r$ . The matrix  $A$  is then in  $\{0, 1\}^{N^2 \times N!}$ . For a column of  $A$  corresponding to the permutation  $\sigma, A(\sigma)$ , we thus have  $A(\sigma)_{ri} = 1$  iff  $\sigma(i) = r$ .

- *Top-Set Data.* These data refer to a concatenation of the comparison data and information on the fraction of customers who have a given product  $i$  as their topmost choice for each  $i$ . Thus,  $A^T = [A_1^T A_2^T]$  where  $A_1$  is simply the  $A$  matrix for comparison data, and  $A_2 \in \{0, 1\}^{N \times N!}$  has  $A_2(\sigma)_i = 1$  if and only if  $\sigma(i) = 1$ . With these definitions, we can now state the following result.

**THEOREM 3.** *Suppose  $\lambda$  is a choice model of support size  $K$  drawn from the generative model. Then,  $\lambda$  satisfies the signature’ and linear independence conditions with probability  $1 - o(1)$  as  $N \rightarrow \infty$  provided  $K = o(\log N)$  for comparison data,  $K = o(\sqrt{N})$  for the top-set data, and  $K = O(N)$  for ranking data.*

Theorem 3 implies that essentially all choice models of sparsity  $\log N$  (and higher) can be recovered from the types of observed data discussed in the theorem. A natural question that arises at this juncture is what a reasonable value of  $K(N)$  might be. To give a sense of this, we provide the following approximation result: A good approximation to *any* choice model for the purposes of revenue estimation is obtained by a sparse choice model with support scaling as  $\log N$ . Specifically, let us restrict ourselves to offer sets that are small, i.e., bounded by a constant  $|\mathcal{M}| \leq C$ ; this is legitimate from an operational perspective and in line with many of the applications we have described. We now show that *any* customer choice model can be well approximated by a choice model with *sparse* support

<sup>7</sup> Though replacement makes repetitions likely, for large  $N$  and  $K \ll \sqrt{N!}$ , they happen with a vanishing probability.

<sup>8</sup> Any distribution with a density on the  $K$ -dim simplex may be picked; we picked uniform for concreteness.

for the purpose of evaluating revenue of any offer set  $\mathcal{M}$  of size up to  $C$ . In particular, we have the following theorem:

**THEOREM 4.** *Let  $\lambda$  be an arbitrary given choice model. Then, there exists a choice model  $\hat{\lambda}$  with support  $O((2C^2 p_{\max}^2 / \varepsilon^2)(\log 2C + C \log N))$  such that*

$$\max_{\mathcal{M}: |\mathcal{M}| \leq C} \left| R(\mathcal{M}) - \sum_{j \in \mathcal{M}} p_j \hat{\lambda}_j(\mathcal{M}) \right| \leq \varepsilon.$$

The proof is provided in Online Appendix A.3. Along with Theorem 3, the above result establishes the potential generality of the signature and linear independence conditions.

In summary, this section visited the issues of explicitly selecting a choice model consistent with the observed data. This is in contrast to our work thus far, which has been simply making revenue predictions. We showed that the robust procedure we used in making revenue predictions may also be seen to yield what is essentially the sparsest choice model consistent with the observed data. Finally, by presenting a family of models for which the sparsest fit to the observed data was unique, and studying the properties of this unique solution, we were able to delineate a data-dependent family of choice models for which the sparsest fit criterion actually yields identification. This formalized the intuitive notion that the complexity of the choice model that can be recovered scales with the amount of data that is available.

## 7. Conclusion and Potential Future Directions

This paper presented a new approach to the problem of using historical sales data to predict expected sales/revenues from offering a particular assortment of products. We depart from traditional parametric approaches to choice modeling in that we assume little more than a weak form of customer rationality; the family of choice models we focus on is essentially the most general family of choice models one may consider. In spite of this generality, we have presented schemes that succeed in producing accurate sales/revenue predictions. We complemented those schemes with extensive empirical studies using both simulated and real-world data, which demonstrated the power of our approach in producing accurate revenue predictions without being prone to overfitting and underfitting. We believe that these schemes are particularly valuable from the standpoint of incorporating models of choice in decision models frequently encountered in operations management. Our schemes are efficient from a computational standpoint and raise the possibility of an entirely data-driven approach to the modeling of choice for use

in those applications. We also discussed some ideas on the problem of identifying sparse or simple models that are consistent with the available marginal information.

With that said, this work cannot be expected to present a panacea for choice modeling problems. In particular, one merit of a structural/parametric modeling approach to modeling choice is the ability to extrapolate. That is to say, a nonparametric approach such as ours can start making useful predictions about the interactions of a particular product with other products only once *some* data related to that product are observed. With a structural model, one can hope to say useful things about products never seen before. The decision of whether a structural modeling approach is relevant to the problem at hand or whether the approach we offer is a viable alternative thus merits a careful consideration of the context. Of course, as we discussed previously, resorting to a parametric approach will typically require expert input on underlying product features that “matter,” and is thus difficult to automate on a large scale. In addition, although our modeling approach is very general, it does not account for competition (see Berry et al. 1995) or intertemporal choice behavior (time dependence in the presence of anticipated discounts or end of the season clearances; see Li et al. 2011).

We believe this paper presents a starting point for a number of research directions. There are numerous directions to pursue from an applications perspective:

1. The focus of this paper has been the estimation of the revenue function  $R(\mathcal{M})$  with the rationale that it forms a core subroutine in essentially any revenue optimization problem seeking to optimize revenues in the face of customer choice. A number of generic algorithms (such as local search) can potentially be used in conjunction with the subroutine we provide to solve such optimization problems. It would be interesting to study such a procedure in the context of problems such as network revenue optimization in the presence of customer choice.

2. Having learned a choice model that consists of a distribution over a small number of rank lists, there are a number of qualitative insights one might hope to draw. For instance, using fairly standard statistical machinery, one might hope to ask for the product features that most influence choice from among thousands of potential features by understanding which of these features best rationalize the rank lists learned. In a different direction, one may use the distribution learned as a “prior,” and given further interactions with a given customer infer a distribution specialized to that customer via Bayes rule. This is effectively a means to accomplishing “collaborative filtering.”

There are also interesting directions to pursue from a theoretical perspective: First, extending our understanding of the limits of identification. In particular, it would be useful to characterize the limits of recoverability for additional families of observable data beyond those discussed in Theorem 3. Second, Theorem 4 points to the existence of sparse approximations to generic choice models. Can we compute such approximations for any choice model but with limited data? Finally, the robust approach in §3 presents us with a family of difficult optimization problems for which the present work has presented a generic optimization scheme that is in the spirit of cutting plane approaches. An alternative to this is the development of strong relaxations that yield uniform approximation guarantees (in the spirit of the approximation algorithms literature).

### Acknowledgments

The authors thank Paat Rusmevichientong for sharing the specifics of the model fit to data from Amazon.com. The authors also thank the anonymous reviewers, whose reviews helped improve the paper.

### References

Anderson SP, De Palma A, Thisse JF (1992) *Discrete Choice Theory of Product Differentiation* (MIT Press, Cambridge, MA).

Belobaba PP, Hopperstad C (1999) Boeing/MIT simulation study: PODS results update. 1999 AGIFORS Reservations and Yield Management Study Group Sympos., April 27–30, London.

Ben-Akiva ME, Lerman SR (1985) *Discrete Choice Analysis: Theory and Application to Travel Demand* (MIT Press, Cambridge, MA).

Berry S, Levinsohn J, Pakes A (1995) Automobile prices in market equilibrium. *Econometrica* 63(4):841–890.

Bierlaire M (2003) BIOGEME: A free package for the estimation of discrete choice models. *Proc. 3rd Swiss Transportation Res. Conf., Ascona, Switzerland*.

Bierlaire M (2008) An introduction to BIOGEME Version 1.7. Tutorial, <http://biogeme.epfl.ch>.

Birkhoff G (1946) Tres observaciones sobre el algebra lineal. *Univ. Nac. Tucuman Rev. Ser. A* 5(1946):147–151.

Calafiore G, Campi MC (2005) Uncertain convex programs: Randomized solutions and confidence levels. *Math. Programming* 102(1):25–46.

Candes EJ, Romberg JK, Tao T (2006) Stable signal recovery from incomplete and inaccurate measurements. *Comm. Pure Appl. Math.* 59(8):1207–1223.

Chandukala SR, Kim J, Otter T, Rossi PE, Allenby GM (2008) Choice models in marketing: Economic assumptions, challenges and trends. *Foundations and Trends® in Marketing* 2(2):97–184.

Cormode G, Muthukrishnan S (2006) Combinatorial algorithms for compressed sensing. Flocchini P, Gąsieniec L, eds. *Structural Information and Communication Complexity* (Springer, New York), 280–294.

Debreu G (1960) Review of R. D. Luce, “Individual choice behavior: A theoretical analysis.” *Amer. Econom. Rev.* 50(1):186–188.

Gallego G, Iyengar G, Phillips R, Dubey A (2004) Managing flexible products on a network. CORC Technical Report TR-2004-01, Cornell University, Ithaca, NY.

Goyal V, Levi R, Segev D (2009) Near-optimal algorithms for the assortment planning problem under dynamic substitution and stochastic demand. Working paper, Columbia University, New York.

Guadagni PM, Little JDC (1983) A logit model of brand choice calibrated on scanner data. *Marketing Sci.* 2(3):203–238.

Kök AG, Fisher ML, Vaidyanathan R (2009) Assortment planning: Review of literature and industry practice. Agrawal N, Smith SA, eds. *Retail Supply Chain Management: Quantitative Models and Empirical Studies* (Springer Science + Business Media, New York), 99–154.

Li J, Granados N, Netessine S (2011) Are consumers strategic? Structural estimation from the air-travel industry. INSEAD Working Paper 2011/104/TOM/ACGRE, INSEAD, Fontainebleau, France.

Luce RD (1959) *Individual Choice Behavior: A Theoretical Analysis* (Wiley, New York).

Mahajan S, van Ryzin GJ (1999) On the relationship between inventory costs and variety benefits in retail assortments. *Management Sci.* 45(11):1496–1509.

Mas-Colell A, Whinston MD, Green JR (1995) *Microeconomic Theory* (Oxford University Press, New York).

McFadden D (1980) Econometric models for probabiistic choice among products. *J. Bus.* 53(3):S13–S29.

Mosteller F, Tukey J (1968) Data analysis, including statistics. Lindzey G, Aronson E, eds. *Handbook of Social Psychology*, Vol. 2 (Addison-Wesley, Reading, MA).

Plackett RL (1975) The analysis of permutations. *Appl. Statist.* 24(2):193–202.

Ratliff RM, Rao V, Narayan CP, Yellepeddi K (2008a) A multi-flight recapture heuristic for estimating unconstrained demand from airline bookings. *J. Revenue Pricing Management* 7(2): 153–171.

Ratliff RM, Rao BV, Narayan CP, Yellepeddi K (2008b) A multi-flight recapture heuristic for estimating unconstrained demand from airline bookings. *J. Revenue Pricing Management* 7(2): 153–171.

Rusmevichientong P, Topaloglu H (2012) Robust assortment optimization in revenue management under the multinomial logit choice model. *Oper. Res.* 60(4):865–882.

Rusmevichientong P, Shen Z-J, Shmoys DB (2010) Dynamic assortment optimization with a multinomial logit choice model and capacity constraint. *Oper. Res.* 58(6):1666–1680.

Rusmevichientong P, Van Roy B, Glynn PW (2006) A nonparametric approach to multiproduct pricing. *Oper. Res.* 54(1):82–98.

Saure D, Zeevi A (2009) Optimal dynamic assortment planning. Working paper, Columbia University, New York.

Talluri K, van Ryzin GJ (2004a) Revenue management under a general discrete choice model of consumer behavior. *Management Sci.* 50(1):15–33.

Talluri KT, van Ryzin GJ (2004b) *The Theory and Practice of Revenue Management* (Springer Science + Business Media, New York).

van Ryzin GJ, Vulcano G (2008) Computing virtual nesting controls for network revenue management under customer choice behavior. *Manufacturing Service Oper. Management* 10(3): 448–467.

Vulcano G, van Ryzin G, Chaar W (2010) OM practice—Choice-based revenue management: An empirical study of estimation and optimization. *Manufacturing Service Oper. Management* 12(3):371–392.

Vulcano G, van Ryzin GJ, Ratliff R (2012) Estimating primary demand for substitutable products from sales transaction data. *Oper Res.* 60(2):313–334.

Wierenga B (2008) *Handbook of Marketing Decision Models* (Springer Science + Business Media, New York).