

# Estimating Preferences and Substitution Patterns from Second-Choice Data Alone\*

Christopher Conlon<sup>†</sup>

Julie Holland Mortimer<sup>‡</sup>

Paul Sarkis<sup>§</sup>

May 29, 2023

## Abstract

We consider identification and estimation of a model of consumer choice where the main source of variation is in the set of products made available to consumers. Instead of relying on variation in the choice environment (prices, product characteristics) we require first-choice probabilities and a subset of (conditional) second-choice probabilities. We develop a semi-parametric low-rank approximation to the matrix of second-choice probabilities that is consistent with mixed logit models of demand but is defined in “product space” and does not require that product characteristics explain substitution patterns. In Monte Carlo experiments we show that our model can replicate a nested logit or random coefficients logit model. We apply our model to a single year of automobile data from Grieco et al. (2021) and show that we can fit substitution patterns with higher accuracy. In a second exercise, we apply our model to the second-choice vending machine experiments in Conlon et al. (2023); Conlon and Mortimer (2021b), and show that we can “complete” the matrix of substitutes.

---

\*[PRELIMINARY and INCOMPLETE] We thank Mark Stein and Mark Vend Company, and seminar participants at NYU Stern, University of Virginia, and IIOC 2021.

<sup>†</sup>New York University, Stern School of Business and NBER: cconlon@stern.nyu.edu

<sup>‡</sup>University of Virginia and NBER: juliemortimer@virginia.edu

<sup>§</sup>Boston College: sarkisp@bc.edu

## 1. Introduction

A large literature in industrial organization (and economics more broadly) is concerned with estimating demand systems for differentiated products. These demand systems often have two key deliverables: the first is the own-price elasticity of demand, and the second is the pattern of substitution across products.

One key challenge in this literature is that the number of cross-product effects grows with the square of the number of products. A demand system with  $J$  products requires estimating  $J^2$  cross-elasticities. This is true for common linear and log-linear demand systems, as well as the second-order flexible Almost Ideal Demand System of Deaton and Muellbauer (1980). This presents practical challenges both in terms of the required amount of variation in the data, and the need for additional instrumental variables.

The literature has primarily addressed this challenge by relying on parametric restrictions or treating products as bundles of characteristics, and formulating demand (and substitution patterns) in terms of these characteristics. For example, the logit model (McFadden, 1974) exhibits the independence of irrelevant alternatives (IIA) property, which restricts substitution to be proportional to observed market shares. The related nested logit (McFadden, 1978) exhibits proportional substitution within pre-determined groups and also across groups themselves. Finite (or continuous) mixtures of IIA logits can generate arbitrarily flexible substitution patterns under additional assumptions (McFadden and Train, 2000). In order to obtain that flexibility, the most common approach is to consider continuous mixtures of consumer preferences for observed characteristics and project substitution patterns onto product characteristics and their interaction with consumer demographics. This is the approach most commonly adopted in industrial organization (Berry et al., 1995, 1999; Nevo, 2000). While these models can be flexible, they have some drawbacks. The first is that their flexibility is generally limited by how well the underlying substitution patterns are explained by the observed product characteristics. The second is that even with modern tools, they can be computationally intensive and challenging to estimate, particularly as they become more flexible (Conlon and Gortmaker, 2020).

We propose a simple estimator, designed for an environment where a researcher observes aggre-

gate data on *market shares* (the probability that a consumer selects a particular product) and the matrix of *second choices* (the conditional probability that a consumer chooses a particular product if their first choice product is not available) from a *single market*. We provide a simple, easy to implement estimator based on a semiparametric mixture of logits to this first- and second-choice data. The key feature of our estimator is that we formulate the problem in *product space* rather than *characteristic space*, but avoid the problem of estimating  $J^2$  elasticities by restricting the rank of the substitution matrix. Instead, we build on a result from Conlon and Mortimer (2021a) and write second-choice probabilities in terms of the first-choice probabilities for a finite number of types  $I$ . This reduces the number of parameters to  $I \times J$ , and restricts the rank of the matrix of substitutes to be no more than  $I$ .

Our approach is likely to be successful if the underlying substitution patterns exhibit a *low rank structure* where  $I \ll J$ . When this is the case, we can extend our estimator to the case where second-choice probabilities are only partially observed and redefine our problem as one of *matrix completion*. Our estimator also tends to produce choice probabilities that are *sparse*, so that not every individual type chooses each option with positive probability. This can be viewed as either a feature: it produces choice probabilities and substitution patterns that are more extreme than most mixed logits and is robust to not all individuals “considering” all products; or a bug: under full consideration it is no longer consistent with full support IID error terms.

Though our estimator recovers individual specific first-choice probabilities, we show it is straightforward to recover indirect utilities and tastes for product characteristics. This builds on a large literature of second-stage approaches in these kinds of models (Nevo, 2001; Bayer et al., 2007; Bayer and Timmins, 2007; Grieco et al., 2023). The most important parameters to recover in a second-stage are likely those governing (endogenous) prices. We show those parameters can be estimated using: (a) standard restrictions arising from instrumental variables; (b) auxilliary information on price cost margins; or (c) auxillary estimates of own-price elasticities. Together with the substitution patterns, these provide semiparametric estimates for most of the objects necessary to make counterfactual equilibrium predictions (mergers, taxes and subsidies, etc.). While academic researchers rarely have access to high quality data on price-cost margins, such information

is routinely provided to antitrust agencies as part of merger investigations.

What sets our estimator apart from the prior literature is that the data requirements are quite different. Rather than utilizing cross-market variation in the characteristics and assortment of products (see Berry and Haile (2014) on “aggregate data” for formal identification results), we rely on observing second-choice data for a single market (more in line with Berry and Haile (2022) on “micro data”). A valid concern is whether such second-choice data is readily available and what the quality of the data is likely to be. An obvious source of such data would be stated-preference second-choice data arising from consumer surveys. For example, the UK Competition and Markets Authority (CMA) frequently conducts surveys asking customers questions such as “if this supermarket were to close, where would you shop?” (Reynolds and Walters, 2008). Likewise, Berry et al. (2004a); Grieco et al. (2021) observe second-choice probabilities for a subset of automobile consumers and use these to inform mixed logit estimates of demand parameters. Another (inexpensive) alternative would be to design purpose-driven surveys such as Conlon and Gortmaker (2023) who survey consumers on second-choice soft-drink choices, or Magnolfi et al. (2022) on ready-to-eat breakfast cereal.<sup>1</sup> A long literature in marketing studies the design and implementation of *conjoint*s for similar purposes Allenby et al. (2019). A (possibly more credible) alternative would be to rely on revealed preference approaches to estimate second choices. For example, Conlon and Mortimer (2013) show how to exploit exogenous variation in the timing of stock-out events, while Conlon et al. (2023) experimentally manipulate product assortment to estimate second-choice probabilities. An online alternative might be to experimentally manipulate the search results facing consumers in order to recover revealed-preference second-choice probabilities.

There are two interpretations of our approach: the first is as a method to use data on second-choices to estimate the parameters of a demand system. The second is that it provides a way to rationalize surveys and other second-choice data in a way that is consistent with a discrete choice demand framework. While this aspect is not always an issue for academic researchers, it is frequently relevant for antitrust agencies. As an example, the CMA might be inclined to treat

---

<sup>1</sup>The latter asks consumers to rank which products are more similar to one other rather than which product they prefer, though the design and implementation of the online survey could be easily adapted. The former survey was completed for less than \$200 over the course of a week.

second-choice surveys as if they were the diversion ratios that enter the first-order conditions of merging parties and calculate UPP. Or, the DOJ or FCC may be inclined to use observational data on customer flows (sometimes called “win/loss data”) in place of diversion ratios.<sup>2</sup> Likewise Farrell and Shapiro (2010) supposes that diversion ratios might be observed in the “normal course of business,” implying that diversion ratios are data that firms might track internally, and could be requested by antitrust authorities as part of an investigation. Conlon and Mortimer (2021a) point out that diversion ratios measured from small price changes, quality changes, and second-choice data are related but not identical. Our framework provides a way to take second choice data and recover the underlying primitives of the demand system, which can then be easily translated into the object(s) of interest, such as the appropriate UPP measure, or merger simulation.

In order to compare our semi-parametric method against other commonly used parametric approaches, we first estimate our semi-parametric model and mis-specified parametric models on data generated from a given model and compare the models’ fit on out-of-sample second-choice probabilities. Using data from Conlon and Mortimer (2021b), we estimate a nested-logit model of demand with nests corresponding to product categories and an outside good in a separate nest, and a random-coefficients model of demand with independent normally distributed tastes  $\beta_i$  on the constant term and three observable product characteristics. We use these estimated models to generate two fake datasets using analytic formulas. We then compare our model’s out-of-sample fit on predicted substitution patterns for rows of the second-choice matrix not used in estimation with three common parametric specifications: a simple logit, a random coefficients logit model with iid normal tastes over nests (RCN), and a random coefficients model with iid normal tastes over characteristics (RCC), estimated via Maximum Likelihood with and without additional information on second-choice probabilities.

When fitting the fake data generated by the nested-logit model, the parametric random-coefficients nested logit (RCN) model fits the data well. This model is almost correctly specified and differs only in whether the functional form follows a normal distribution or an extreme value distribution. The addition of second-choice moments improves the model’s prediction error, resulting in nearly

---

<sup>2</sup>See Qiu et al. (2021) for an in depth analysis of win/loss data.

perfect fit. On the same data, the random-coefficients model based on characteristics (RCC) appears to be a poor fit while our semi-parametric model performs as well as, or better than, the RCN model estimated from observational data alone (and much better than the RCC model with or without additional moments). Similarly, when fitting fake data generated by the model based on characteristics, the RCC model fits well, the RCN model performs poorly, and our semi-parametric approach performs very well.

As a second exercise, we apply our estimator to a subset of data from Grieco et al. (2021). The 2015 MaritzCX survey provides a matrix of second choice probabilities. This survey includes recent car purchasers and is based on new vehicle registrations. Although this matrix includes survey data on all products, meaning that all rows are observed, we view the exercise as a demonstration of our estimator's effectiveness when dealing with a large number of products ( $J = 318$ ) in a context for which substitution patterns may be influenced by endogenous characteristics, such as prices. To evaluate our estimator, we adopt the same cross-validation approach used in the previous exercise and compare our preferred specification to the parametric model estimated in the Grieco et al. (2021) study.

In this context, our semi-parametric estimator outperforms state-of-the-art parametric models that use significantly more data. With our preferred specification (selected by cross-validation), our model can rationalize both extremely high diversion ratios and more spread-out substitution patterns. For instance, our model can match the substitution between two cargo vans at around 67% and consistently keep the substitution between similar sedans below 10%. In contrast, logit-type parametric models are unlikely to achieve this, as shown by our comparison with the estimates from Grieco et al. (2021).

As a final exercise, we apply our semi-parametric estimator to experimental second-choice data from Conlon and Mortimer (2021b). The experiment involved multiple treatment arms for which one or two products were exogenously removed from vending machines in downtown Chicago office buildings. This allows us to measure substitution to remaining products without any parametric restrictions on demand. Thus, the resulting data consists of second-choice probabilities from only a few products relative to the whole set of products available, so that a significant portion of the

matrix of second choices for all products is missing. This provides an illustration of the *low rank matrix completion* aspect of our approach.

Our semi-parametric model outperforms a plain or mixed logit fit to the same dataset by identifying pairs of products with high diversion ratios and including diversion to the outside good. It infers relatively large amounts of substitution between similar products such as Twix and Raisinets, even though neither product is exogenously removed, and we don't use any information about the characteristics of either product. In addition, we estimate little to no substitution from salty snacks (e.g., Lay's potato chips) to popular products like Snickers and Twix. While this substitution pattern would be possible in a nested logit, we haven't provided the estimator with any information about product categories either. Furthermore, our estimator produces sparse consumer choices, indicating that some consumers have very low or zero probability of selecting certain options. We do not differentiate whether this sparsity is due to consumers not considering certain products or consumer preferences (where consumers always prefer the outside option to certain products).

Our results demonstrate that the matrix of second-choice probabilities has a regular pattern that can be approximated with low-rank methods under weak parametric restrictions. Therefore, if some substitution patterns are observed, estimating substitution between remaining products becomes a matrix completion exercise. This indicates that our estimator can be highly adaptable and easily estimated by antitrust authorities under time constraints and with limited data.

The remainder of the paper proceeds as follows: Section 2 reviews the results in Conlon and Mortimer (2021a) and describes the estimator. Section 3 provides the results of the Monte Carlo simulations. Section 4 provides the results of the estimation exercise on the survey second-choice data from the U.S. auto industry. Section 5 and Section 5.1 describe the field experiments with Mark Vend Company and provide results of the estimator in that setting. Section 6 provides extensions and Section 7 concludes.

## 2. Model

### 2.1. Our Estimator

Throughout we use *lower case bold* to denote the  $\dim(\mathcal{J})$  vector of all choices (e.g.  $\mathbf{s}_i, \mathbf{p}$ ) and ***UPPER CASE BOLD*** to denote the  $\dim(\mathcal{J}) \times \dim(\mathcal{J})$  matrix of all choices (e.g.  $\mathbf{D}$ ), and *caligraphic font*  $\mathcal{S}_j, \mathcal{D}_{j \rightarrow k}$  to denote observed data. We will at times abuse notation and allow  $\mathcal{J}$  to denote both a set and its cardinality  $\dim(\mathcal{J}) = \mathcal{J}$ .

We begin with a researcher who has data on consumers who make a discrete choice among a set of products in  $\mathcal{J}$  including an “outside” or “no purchase” option which we denote with  $j = 0$ . Furthermore, we assume that the researcher has access only to aggregate data on the first choices  $\mathcal{S}_j$  and second choices  $\mathcal{D}_{j \rightarrow k}$  for consumers facing a single choice environment (market). We denote the first and second-choice probabilities below:

$$\mathcal{S}_j = \mathbb{P}(\text{chooses } j \in \mathcal{J}) \tag{1}$$

$$\mathcal{D}_{j \rightarrow k} = \mathbb{P}(\text{chooses } k \in \mathcal{J} \setminus \{j\} \mid \text{chooses } j \in \mathcal{J}) \tag{2}$$

Later, we allow for a set of  $(j, k)$  tuples for which  $\mathcal{D}_{j \rightarrow k}$  is observed, denoted  $P_\Omega$ , and for which it is unobserved denoted  $P_{\bar{\Omega}}$ .

A common approach is to assume that these choice probabilities are generated by consumers  $i$  making discrete choices to maximize (random) indirect utilities:  $u_{ij} = V_{ij} + \varepsilon_{ij}$ . The most common specification assumes that  $\varepsilon_{ij}$  is IID and type I extreme value (Gumbel) distributed, so that the probability that  $i$  chooses  $j$  (conditional on the vector of  $V_{ij}$ 's denoted  $\mathbf{V}_i$ ) is given by:<sup>3</sup>

$$\mathbb{P}(u_{ij} > u_{ij'}; \text{ for all } j \neq j' \mid \mathbf{V}_i) = \frac{e^{V_{ij}}}{\sum_{j' \in \mathcal{J}} e^{V_{ij'}}} \equiv s_{ij}(\mathbf{V}_i). \tag{3}$$

The unconditional probabilities require integrating out over the distribution of  $V_{ij}$ , which is typically

---

<sup>3</sup>An alternative to the multinomial logit model is the multinomial probit model. This lets  $\varepsilon_i \sim N(0, \Sigma)$  where  $\Sigma$  is a  $J \times J$  matrix. Like the log-linear model, this requires estimating on the order of  $\frac{J \cdot (J-1)}{2}$  parameters, which makes it impractical for large  $J$ .

done by discretizing the distribution of heterogeneity with some probability weights corresponding to different vectors of  $\mathbb{P}(\mathbf{V}_i = \mathbf{v}_i) = \pi_i$  so that  $\pi_i \geq 0$  and  $\sum_{i=1}^I \pi_i = 1$  (so that  $\pi$  constitutes a valid probability measure):

$$\mathbb{P}(u_{ij} > u_{ij'}; \text{ for all } j \neq j') = \int s_{ij}(\mathbf{V}_i) f(\mathbf{V}_i) \partial \mathbf{V}_i \approx \sum_{i=1}^I \pi_i s_{ij}(\mathbf{V}_i) \equiv s_j. \quad (4)$$

It is helpful to define some additional notation: let  $\mathbf{s}_i$  be the  $\dim(\mathcal{J})$  vector of type-specific (conditional) choice probabilities with entries  $s_{ij}$  defined in (3), let  $\mathbf{s}$  be the vector of unconditional choice probabilities  $s_j$  defined in (4), and let  $\mathbf{S}$  be a  $\dim(\mathcal{J}) \times I$  matrix with column vectors  $\mathbf{s}_i$ . If  $\pi$  denotes an  $I$  vector with entries  $\pi_i$  then we can write  $\mathbf{s} = \mathbf{S}\pi$ .

In our previous work, Conlon and Mortimer (2021a), we show that for any mixed logit, the second-choice probabilities from (2) can be written in terms of the weights and the conditional and unconditional probabilities  $(\pi_i, s_{ij}, s_j)$  from (3) and (4):<sup>4</sup>

$$D_{j \rightarrow k} = \sum_{i=1}^I \pi_i \cdot \frac{s_{ik}}{1 - s_{ij}} \cdot \frac{s_{ij}}{s_j}. \quad (5)$$

It is convenient to interpret (5) as the  $(j, k)$ th entry in the second-choice matrix  $\mathbf{D}$ .

Our estimator simply matches the observed first and second choice data from (1) and (2) with the predicted versions from (4) and (5). We can accomplish this by minimizing the (potentially weighted) least squares error ( $\ell_2$ /Frobenius norm) so that:

$$\min_{(\mathbf{S}, \pi) \geq 0} \|\mathcal{P}_\Omega(\mathcal{D} - \mathbf{D})\|_{\ell_2} + \lambda \|\mathcal{S} - \mathbf{S}\pi\|_{\ell_2} \text{ with } \pi \cdot \mathbf{1}_I = 1, \quad \mathbf{1}'_{\mathcal{J}} \mathbf{S} = \mathbf{1}_I. \quad (6)$$

The goal is to match the observed second choice probabilities in  $\mathcal{D}$  and also the first-choice probabilities (market shares)  $\mathcal{S}$  subject to some tuning parameter (Lagrange multiplier)  $\lambda$ . A typical challenge in computer science for problems like (6) is to avoid overfitting by restricting either the

---

<sup>4</sup>For the case of “second choice data”, Conlon and Mortimer (2021a) show it doesn’t matter whether second choices are obtained by raising the price to the choke price, decreasing the quality such that no individuals purchase, or removing the product from the choice set. In all cases, we average over all individual diversion ratios  $D_{jk,i} = \frac{s_{ik}}{1 - s_{ij}}$  and weight them in accordance with the fraction of  $j$ ’s sales they represent  $w_i = \frac{s_{ij}}{s_j}$

rank of  $\mathbf{D}$  or its nuclear norm (sum of singular values). Below, we show that the rank of  $\mathbf{D}$  is bounded by the number of types  $I$ , and we can restrict the rank of the matrix by directly limiting the number of types. More generally, we propose that the tuning parameters  $(\lambda, I)$  be chosen by cross-validation.

To see the rank restriction imposed by  $I$  we can simply re-write (5) in matrix form as:<sup>5</sup>

$$\begin{aligned}\mathbf{D} &= \left( \sum_{i=1}^I \pi_i \cdot \mathbf{s}_i \cdot \left[ \frac{1}{(1 - \mathbf{s}_i)} \right]^T \cdot \text{diag}(\mathbf{s}_i / \mathbf{s})^{-1} \right)^T \\ &= \text{diag}(\mathbf{s})^{-1} \cdot \left( \sum_{i=1}^I \pi_i \cdot \left[ \frac{\mathbf{s}_i}{(1 - \mathbf{s}_i)} \right] \cdot \mathbf{s}_i^T \right)\end{aligned}\quad (7)$$

This shows that we can write  $\mathbf{D}$  as the sum of  $I$  rank-one matrices (outer product of vectors). The immediate result from (7) is that it shows us how to construct a low-rank  $I \ll \mathcal{J}$  representation of a potentially large  $\mathcal{J} \times \mathcal{J}$  matrix of substitution patterns. The plain IIA logit model corresponds to  $I = 1$ , in Appendix A.2 we show that the nested logit model limits the rank to be less than or equal to the number of nests.

We can also re-write the constraints in (6) as  $\ell_1$  constraints so that:

$$\min_{(\mathbf{S}, \boldsymbol{\pi}) \geq 0} \|\mathcal{P}_{\Omega}(\mathcal{D} - \mathbf{D})\|_{\ell_2} + \lambda \|\mathcal{S} - \mathbf{S} \boldsymbol{\pi}\|_{\ell_2} \quad \text{with} \quad \|\boldsymbol{\pi}\|_{\ell_1} \leq 1, \quad \|\mathbf{s}_i\|_{\ell_1} \leq 1. \quad (8)$$

It should be clear that this represents a *non negative LASSO* problem (Wu et al., 2014). This means we are likely to get *sparse solutions* to the program above so that  $s_{ij} = 0$  for many  $(i, j)$ . An economic interpretation might be that this is a product that type  $i$  really despises  $V_{ij} \rightarrow -\infty$ , or that type  $i$  is unaware of or does not consider  $j$ .<sup>6</sup> In a sense, we are agnostic about (and robust to)

---

<sup>5</sup>Here  $\text{diag}(\mathbf{s}_i)$  is a diagonal matrix with entries  $s_{ij}$  and  $\text{diag}(\mathbf{s})^{-1}$  is a diagonal matrix with entries  $\frac{1}{s_j}$ . The  $\text{diag}(\mathbf{s})^{-1}$  term is serving to row-normalize the matrix so that  $\sum_{j \neq k} D_{jk} = 1$  for each row. The diagonal entries  $D_{jj,i} = \frac{s_{ij}}{1-s_{ij}}$ , while not interpretable as a diversion ratios, are related to the willingness to pay (WTP) for product  $j$  (at least when  $s_{ij}$  isn't too large). See Conlon and Mortimer (2021a).

<sup>6</sup>A large literature examines consideration sets in discrete choice models which may result from rational inattention (Matějka and McKay, 2015; Manzini and Mariotti, 2014), cognitive capacity, or random attention (Dardanoni et al., 2020; Masatlioglu et al., 2012). A significant challenge in this literature is to separate the impact of consumer preferences from heterogeneity in consideration sets. Typically, this requires either the presence of an exclusion restriction that shifts consideration independently of preference (Goeree, 2008), or exploiting differences in functional forms (Abaluck and Adams-Prassl, 2021). Other recent approaches can account for unobserved and limited consideration, but at the cost of potentially losing point identification (Barseghyan et al., 2021a,b).

the particular reason for which  $s_{ij} = 0$  (though it will largely be a result of the  $\ell_1$  penalty). Because second-choices depend on  $\frac{s_{ik}}{1-s_{ij}}$ , it is also worth noting that sparsity in  $\mathbf{s}_i$  will tend to create sparsity in  $\mathbf{D}$ , particularly when the observed data  $\mathcal{D}$  is sparse (or nearly sparse).<sup>7</sup> A common critique of logit and mixed logit demand systems is that all products are necessarily substitutable (at least a little) with one another.<sup>8</sup>

We should also note that the estimator in (6) or (8) is a *minimum distance* estimator, and the underlying asymptotic thought experiment would require either: (a) observing the complete matrix  $\mathcal{D}$ ; or (b) taking the number of choices  $\mathcal{J} \rightarrow \infty$  as in Berry et al. (2004b).

[Should we show that standard MD inference applies here?]

## 2.2. Second Stage

Our estimator recovers estimates of  $\widehat{\mathbf{S}}$  or  $\widehat{s}_{ij}$  and  $\widehat{\pi}$ , but these alone are not sufficient to calculate price-elasticities and consumer welfare. Following a long literature that separates the estimation of heterogeneous preferences from addressing the endogeneity of prices (Goolsbee and Petrin, 2004; Bayer et al., 2007; Bayer and Timmins, 2007; Grieco et al., 2023), we consider a second-stage in order to recover coefficients on prices.

In our first-step, we required the mixed logit only in the definition of second-choice probabilities in (5). In order to recover price sensitivities, we must lean harder on the assumption that we can write  $u_{ij} = V_{ij} + \varepsilon_{ij}$  where  $\varepsilon_{ij}$  is an IID Type I extreme value error term. Following the logic in Berry (1994) we can write:

$$V_{ij} - V_{i0} = \begin{cases} \ln \widehat{s}_{ij} - \ln \widehat{s}_{i0} & \text{if } \widehat{s}_{ij} > 0, \\ \text{n.a.} & \text{if } \widehat{s}_{ij} = 0. \end{cases} \quad (9)$$

This differs from the usual setup in two ways: (1) we recover  $i$  specific utility parameters  $V_{ij}$ ; (2) in the case where we have sparse choice probabilities  $s_{ij} = 0$  we need to be careful about interpretation. If the reason that  $s_{ij} = 0$  is that the consumer  $i$  is unaware of  $j$ , then we can simply

---

<sup>7</sup>The expression in (5) implies  $D_{jk,i} = 0$  IFF  $s_{ij} = 0$  or  $s_{ik} = 0$ , which creates neither a theoretical nor practical problem.

<sup>8</sup>One way this has been addressed previously is the *pure characteristics model* of Berry and Pakes (2007).

ignore the fact that  $s_{ij} = 0$ . If instead,  $s_{ij} = 0$  because a consumer is highly price sensitive and never chooses a luxury car, then we have to model the selection rule more carefully.

### Instrumental Variables

Most of the literature assumes that  $V_{i0} = 0$  for each type  $i$ , though it isn't clear that is necessarily required here.<sup>9</sup> Our goal is to recover  $\beta_i^p = \frac{\partial V_{ij}}{\partial p_j}$ , and we present several approaches under different data environments. The most familiar approach would be to construct a second minimum distance estimator, where we stack up all  $(i, j)$ :

$$\min_{\beta < 0, \xi} \|z'_j(\ln \hat{s}_{ij} - \ln \hat{s}_{i0} - x_j \beta_i - p_j \beta_i^p + \xi_j)\| \quad (10)$$

This approach would require instruments for prices  $z_j$  as in Berry et al. (1995) or Berry and Haile (2014). It is unlikely is that our instruments will vary across  $i$  as well as across  $j$ . Though if we were willing to set  $\xi_j = 0$ , we could run  $I$  separate cross-product regressions (which is the primary source of variation here). Likewise if  $\mathcal{J}$  is large, and we are not interested in the interpretation of non-price  $\beta_i$  coefficients then we can potentially estimate a partially linear model separately for each  $i$ :

$$\min_{\beta_i^p, f_i(\cdot)} \|z'_j(\ln \hat{s}_{ij} - \ln \hat{s}_{i0} - f_i(x_j) - p_j \beta_i^p)\|_{\ell_2}.$$

### Observed Elasticities

An alternative approach would be to calibrate  $\beta_i^p = \frac{\partial V_{ij}}{\partial p_j}$  using observed own-price elasticities. These might be estimated in any number of ways: from a (quasi)-experiment; from another study; or from a simpler (ie: plain logit) demand system on observational data.

For any mixed logit, the own-price elasticities follow a well-known format, which depends on objects we can estimate in the first-stage and  $\beta_i^p$ , and we can construct another minimum distance estimator:

$$\min_{\beta_i^p < 0} \left\| \mathcal{E}_{jj} - \frac{p_j}{s_j} \sum_{i=1}^I \beta_i^p \cdot \hat{\pi}_i \cdot \hat{s}_{ij} \cdot (1 - \hat{s}_{ij}) \right\|_{\ell_2}. \quad (11)$$

---

<sup>9</sup>This would be sumsumed into the type-specific coefficient on the constant  $\beta_{i0}$  below.

We need to observe at least as many elasticities as types  $I$ . Things simplify further if  $\beta_i^p = \beta^p$ , in which case a single elasticity (or the average elasticity) is sufficient to identify  $\beta^p$ . In our empirical example, we estimate  $\beta_i^p$  using the own-elasticities estimated in .

### Observed Price Cost Margins

Another possibility is that average price cost margins ( $p_j - c_j$ ) are often provided to antitrust authorities as part of merger review. Ideally these would be observed at the product level, but also possibly at the firm level. We need to construct  $\Delta$  the  $\mathcal{J} \times \mathcal{J}$  matrix of demand derivatives with entries  $\frac{\partial q_j}{\partial p_k}$ . We also need one additional piece of information, the ownership matrix  $\mathcal{H}$  which typically has entries equal to 1 if products  $(j, k)$  have the same owner and zero otherwise. This let's us write:

$$\mathbf{c} = \mathbf{p} - \left( \mathcal{H} \odot \left( \sum_{i=1}^I \pi_i \cdot \Delta_i(\beta_i^p) \right) \right)^{-1} \mathbf{s}, \quad (12)$$

$$\Delta_i(\beta_i^p) = \beta_i^p \left( -\mathbf{s}_i \mathbf{s}_i^T + \text{diag}[\mathbf{s}_i] \right).$$

This would enable us to match observed and predicted price-cost margins, or observed marginal costs to predicted marginal costs.

$$\min_{\beta_i^p < 0} \left\| \mathcal{C}_{jj} - c_j(\beta_i^p, \mathcal{H}, \hat{S}, \hat{\pi}) \right\|_{\ell_2}. \quad (13)$$

As with the elasticity we could match price-cost margins for a subset of products (or for the average by firm, etc.).

### 2.3. Optional $\ell_2$ Penalization

Our estimator is already using a non-negative LASSO  $\ell_1$  penalty term on  $s_{ij}$  and  $\pi$  as part of the “adding up” constraint that guarantees these are valid probability measures. There is a significant literature on joint  $\ell_1$  and  $\ell_2$  penalization (called *elastic net*) and see Heiss et al. (2022)

For all three cases in Section 2.2, we may want to add a penalty term to our minimum distance estimators of the form  $\lambda_p \cdot \|\beta_i^p - \sum_i \pi_i \beta_i^p\|_{\ell_2}$  or  $\lambda_p \cdot \|\beta_i^p\|_{\ell_2}$  to shrink the recovered  $\beta_i^p$  parameters towards the overall mean, or towards zero respectively. This would have the effect of penalizing

“outliers” in  $\beta_i^p$ , which may be useful if we rely on cross-sectional variation with a limited number of observations, or to prevent extreme and imprecise  $\beta_i^p$  values for types where  $\pi_i$  is very small.

Likewise, even though we are already placing an  $\ell_1$  penalty term on  $s_{ij}$  such that  $\sum_j s_{ij} = 1$ , we may also want to consider an  $\ell_2$  penalty term of the form:  $\|s_{ij} - \mathbf{s}\|_{\ell_2}$  or  $\|s_{ij} - \mathcal{S}\|_{\ell_2}$ . This would have the effect of shrinking the estimated  $s_{ij}$  back towards the plain logit “prior”, and preventing the individual types from getting “too extreme”. This might be useful as  $I$  becomes large and the variance of estimates increases, though it may be easier to estimate a model with fewer types. This may not be good idea if the true distribution of  $\mathbf{s}_i$  includes extreme types.

Finally, we could consider an  $\ell_2$  penalty on the type-weights  $\pi_i$  so that  $\sum_{i=1}^I \pi_i^2 \leq c_\pi$ . This would effectively penalize the concentration/HHI of the types and push us towards estimating types with weights  $\frac{1}{I}$  and away from types with very large or very small weights. This might be of practical use if models with large numbers of types “collapse”, but it is otherwise not obviously useful.

## 2.4. Comparisons

Our semiparametric model is different from the fixed grid estimator of Fox et al. (2011); Heiss et al. (2022) who formulate the problem in characteristic space and construct a “prior” on mixture distribution for the coefficients  $\beta_i^* \sim g(\beta)$ . They draw a large number of points  $\beta_i$  (over 1000) from the prior, and then compute  $s_{ij}^*(\beta_i)$  ahead of estimation on a fixed grid. They use a non-negative Lasso/ $L_1$  penalty to impose sparsity such that a small number of types (20-50) receive positive weight. They search over  $\pi_i$  in order approximate the true  $f(\beta_i)$  via a finite mixture.

$$\min_{\pi} \sum_{j,t} \left( \mathcal{S}_{jt}(\mathbf{x}_t) - \sum_i \pi_i \cdot s_{ijt}^*(\beta_i^*, \mathbf{x}_t) \right)^2 \quad \text{subject to} \quad s_{ijt}^*(\beta_i^*, \mathbf{x}_t) = \frac{e^{\beta_i^* x_{jt}}}{1 + \sum_{j'} e^{\beta_i^* x_{j't}}} \\ 0 \leq \pi_i \leq 1, \quad \sum_i \pi_i = 1 \tag{FKRB}$$

Both models estimate a constrained least squares problem for the aggregate shares, though theirs requires variation in  $\mathbf{x}_t$  across markets, and ours requires second choices within a single-market. Another major difference is that we search over both the probabilities of types and the preferences

of the types  $(s_{ij}, \pi_i)$  in *product space*, rather than considering a pre-specified grid of types in *characteristic space* and searching over  $\pi_i$  only.<sup>10</sup> The most important difference is that in their model generates a large number of types  $I$  for the fixed grid, while our approach is meant to keep  $I$  (and the rank of the diversion matrix) as small as possible.

Another semiparametric model is Raval et al. (2017). Their semiparametric model for hospital demand groups consumers into  $g \in G$  bins based on observable characteristics such as income, zip code, severity of diagnosis, and age. They have a tuning parameter on the minimum number of consumers per bin (which governs the number of bins). Within a group they assume all individuals have the same  $\beta_g$  but do not require anything about  $\beta_g$  and  $\beta_{g'}$  other than that  $\xi_j$  is common across groups. They assume that preference follow a plain logit within each bin:

$$s_{g(i),j} = \frac{e^{\beta_g x_j + \xi_j}}{1 + \sum_{j'} e^{\beta_g x_{j'} + \xi_{j'}}}, \quad D_{kj,i} = \frac{s_{g(i),j}}{1 - s_{g(i),k}} \quad (\text{RRT})$$

Raval et al. (2022) estimate diversion for hospitals using the diversion above and use observed second-choice diversion ratios (from natural disaster induced closures) to validate models of hospital demand, but do not use this variation to estimate the parameters of the model. The main difference between our models is that they are able to observe consumers and group them into demographic bins before estimation, whereas we consider aggregate data and must infer the mixing distributions.

In short, Fox et al. (2011) fix a grid of  $\beta_i^*$  (and hence  $s_{ij}^*$ ) and estimate  $\pi_i$  using non-negative Lasso, while Raval et al. (2017) know the assignments of individuals to types  $g(i)$ , and type specific shares  $s_{g(i),j}$ , and estimate  $\beta_g$  for each group. We estimate both:  $(\pi_i, s_{ij})$ . However, our method requires observing at least some second choices, and our goal is to explain substitution with the smallest number of types possible.

Our approach in (6) is also related to a large class of problems in computer science known as “non-negative matrix factorization” which solve:

$$\min_{W \in \mathbb{R}_+^{(J+1) \times I}, H \in \mathbb{R}_+^{I \times (J+1)}} (\mathcal{D}_{jk} - (W H)_{jk})^2. \quad (\text{NMF})$$

---

<sup>10</sup>Arguably it would be straightforward to include product dummies in  $x_j$  and reformulate their approach in *product space* as well.

This looks for a rank  $I$  approximation to  $\mathcal{D}$  in terms of two component matrices  $X$  and  $W$  that each have non-negative elements. Absent the non-negativity constraints it is well known (Eckart–Young–Mirsky theorem) that the singular value decomposition gives the best rank  $I$  approximation to  $\mathcal{D}$ .

On one hand, our discrete choice setup places even more assumptions on (NMF) than simply non-negativity. Rather than factorize into two rank  $I$  matrices  $(W, H)$  we search for a single rank  $I$  matrix  $\mathbf{S}$  and impose a particular relationship between the columns of  $\mathbf{S}$  and the matrix  $\mathbf{D}$ . Our problem may be easier because  $\mathcal{D}$  has a predictable structure and is amenable to a low-rank approximation, and discrete choice theory provides some structure on the component vectors  $\mathbf{s}_i$ , including that  $\mathbf{s}_j = \sum_i \pi_i \cdot \mathbf{s}_i$ .<sup>11</sup>

[Do we need a formal identification result?? – seems likely for  $I \ll \mathcal{J}$  case.]

## 2.5. Computational Details

The minimum distance problem in (6) is non-convex, but easy to estimate relative to typical GMM or maximum likelihood estimators, and can generally be solved within a few seconds with a standard constrained L-BFGS optimization routine, or with an unconstrained Stochastic Gradient method (such as Adam (Kingma and Ba, 2017)). The objective is quadratic in parameters, and most of the constraints are quadratic or linear with the exception of the  $\mathbf{D}$  matrix in (7) which are nonlinear and non-convex (though the derivatives are simple).

Typically, latent class logit models parameterize  $V_{ij}(x_j) = \beta_i x_j + \xi_j$ , and are estimated via full information maximum likelihood or via the EM algorithm as described in Greene and Hensher (2003). Estimation is often challenging for finite mixture models which estimate both  $\pi_i$  and  $\beta_i$ .<sup>12</sup> Our approach is both much easier to estimate, and is not necessarily restricted to interactions with observed covariates  $x_j$ . The main difference between our approach and the typical latent class logit

---

<sup>11</sup>One challenge of (NMF) is that it may not be identified (Ke et al., 2022) for the classic Latent Dirichlet Allocation (LDA) problem of Blei et al. (2003). It isn't clear whether or not the additional structure we impose is sufficient to restore identification.

<sup>12</sup>Maximum likelihood estimation of the latent class logit is notoriously difficult. The `gmm1` R package only offers FIML estimation for small problems and not EM estimation for large problems. Greene and Hensher (2003) propose an EM algorithm which appears to only be available in NLOGIT/LIMDEP.

is that we are estimating the model in *product space*. The summing up constraints are what enforce that the resulting model is consistent with mutually exclusive and exhaustive discrete choice.

In general, the problem scales with the number of products  $J$  and number of types  $I$  but not the number of “markets” or “observations”. With  $I$  types there are  $(J + 1) \times I$  parameters.<sup>13</sup> If  $I = 1$  then we are just fitting an IIA logit by least squares (with auxiliary moments for  $\mathcal{D}$ ). If  $I = 2$  then we have two latent classes like Berry and Jia (2010) use for business and leisure travelers. As we increase  $I$ , we increase the complexity of the model. In the limit we should be able to approximate any  $\mathcal{D}$  as  $I \rightarrow J$ .<sup>14</sup> Our hope is that like many phenomenon,  $\mathcal{D}$  may have a low-rank structure, or at least be amenable to low-rank approximation (Chen et al., 2021; Udell and Townsend, 2019).

### 3. Monte Carlo Simulations

In order to compare our semi-parametric method against other common parametric methods, we generate data from a given model and estimate both our semi-parametric model and mis-specified parametric models, then compare the models’ fit on out-of-sample (on  $\mathcal{P}_{\overline{\Omega}}$ ) second-choice probabilities  $\|\mathcal{P}_{\overline{\Omega}}(\mathcal{D} - \mathbf{D}(\mathbf{S}, \pi))\|$  in both  $\ell_2$  (MSE) and  $\ell_1$  (MAD).

This exercise has two goals: verify whether our estimator can approximate common models, and identify how much data and model complexity is sufficient to achieve reasonable performance. In addition to our semi-parametric model, we estimate two mixed logit models with  $V_{ij}(x_j) = \beta_i x_j + d_j$  where  $\beta_i \sim N(\mu, \Sigma)$  and  $\Sigma$  is diagonal. In the first,  $x_j$  is given by a set of dummies on product categories, which we label *random coefficients on nests* or (RCN). In the second  $x_j$  is given by the observed characteristics in the vending data (salt, sugar, and nut content), which we label *random coefficients on characteristics* (RCC).

#### 3.1. Data Generating Process

We start with the data from Conlon and Mortimer (2021b) which included observational data on 66 vending machines in office buildings in downtown Chicago. For several weeks, the authors ran six experiments where the top-selling products in each category were removed.

---

<sup>13</sup>This ignores the parameters that are pinned down by the constraints in (6) such as  $s_j$ .

<sup>14</sup>We don’t have a formal result here, though it seems like it would merely be a restatement of McFadden and Train (2000).

We begin by estimating a nested logit model  $V_{ijt} = \delta_j + \xi_t + \varepsilon_{ijt}(\rho)$  using the observational data for  $\bar{J} = 45$  and  $G = 6$  nests corresponding to each product category (Salty Snacks, Chocolate Candy, Non-Chocolate Candy, Cookies, Pastry, and Other) as well as an outside good in a separate nest. We calibrate  $\rho = 0.25$  so that diversion to the outside good is approximately 30%.

We repeat this exercise on the same data where instead we calibrate a random coefficients model  $V_{ijt} = \beta_i x_j + \delta_j + \xi_t$  with independent normally distributed tastes  $\beta_i$  on (constant, salt, sugar, and nut content).<sup>15</sup>

We then treat these parameter estimates as the “ground truth”  $\theta_0$ , and use the estimated nested or random coefficients logit model to generate fake data. We construct two datasets: one single market where all products are available and  $T = 250$  markets of  $M = 1,000$  consumers where only  $J = 30$  products are available in any given market. For each dataset, we generate shares and diversion following the analytic formulas given by the nested logit described in Appendix B.

1. Set the parameters of the model, product mean utilities and nonlinear parameters  $\theta_0 = [\delta_j, \rho]$  or  $\theta_0 = [\delta_j, \Sigma]$ .
2. Generate the “true” shares  $S(\mathbf{x}, \theta_0)$  and second-choice diversion ratios  $D(\mathbf{x}, \theta_0)$  assuming full availability  $\bar{J} = 45$  using the nested logit formulas from Appendix B.
3. For each market  $t = 1, \dots, T$ , draw  $J = 30$  products to be available in that market  $\mathbf{x}_t$ , making sure each nest  $g$  contains at least one product.
  - (a) Compute shares and multiply by market size  $M \cdot s_j(\mathbf{x}_t)$  to obtain quantities  $q_j(\mathbf{x}_t)$ .
  - (b) Generate the  $J \times J$  market specific diversion matrix  $D(\mathbf{x}_t)$ .

### 3.2. Comparing semi-parametric and parametric models

To compare our model with common parametric specifications, we estimate two random coefficients models, RCC and RCN on the  $T = 250$  markets from our fake data. We estimate the parametric models via Maximum Likelihood. In alternative specifications, we augment these data with addi-

---

<sup>15</sup>We calibrate the diagonal elements to be somewhat more heterogeneous than those estimated in Conlon and Mortimer (2021b)  $\sigma_0 = 3.52, \sigma_{salt} = 1.00, \sigma_{sugar} = 1.8, \sigma_{nuts} = 0.38$

tional moments matching selected rows from  $\mathcal{D}_{j,\cdot}(\mathbf{x})$  to the second-choice diversion ratios predicted by the model for the full sample of  $J = 45$  products  $D_{j,\cdot}(\mathbf{x}, \theta)$ .

Our semiparametric estimator from (6) uses less data than the parametric models. We don't use any of the observed variation in choice sets across markets  $t$ , nor any product characteristics. Instead, we use only the aggregate shares when all products are available  $S(\mathbf{x}, \theta_0)$  and a subset of  $L \in \{6, 8, 10\}$  rows from the matrix of second-choice diversion ratios  $\mathcal{D}_{j,\cdot}(\mathbf{x}, \theta_0)$ .<sup>16</sup>

Because the models are parametrized differently, rather than compare  $\hat{\theta}$ , we instead compare models based on out-of-sample predicted substitution patterns  $\|\mathcal{D}(\mathbf{x}, \theta_0) - D(\mathbf{x}, \hat{\theta})\|$ . We compare models based only on rows of the diversion matrix not used in estimation. We report the out-of-sample RMSE and median absolute deviation (MAD) in Figure 1. What is immediately obvious is that the RCN model fits the data quite well in terms of MAD and RMSE. Because the data generating process is a nested logit, this model is nearly correctly specified and differs only as to whether the functional form follows a normal distribution or a nested logit distribution. The additional second-choice moments improve the prediction error of the RCN model (in terms of RMSE) so that it fits nearly perfectly. The RCC model appears badly misspecified; the additional moments improve the RMSE significantly but the MAD very little.

For  $I \geq 4$  consumers, the semi-parametric model performs as well or better than the RCN model estimated from observational data alone (and much better than the RCC model with or without the additional moments) in terms of both MAD and RMSE. The number of "observed" rows  $L \in \{6, 8, 10\}$  does not have a significant impact on the prediction error of the model, suggesting that  $L = 6$  rows is "enough" to explain the substitution patterns in the remaining  $J = 45 - L$  rows of the matrix.<sup>17</sup>

We get a similar result in Figure 2 when we use the RCC model as the data generating process. In this case the semi-parametric model substantially outperforms the RCN model. The addition of the second-choice moments does little to improve the predictive power of the RCN model.

We calculate the "true" rank of the  $\mathcal{D}(\mathbf{x}, \theta_0)$  from the nested logit DGP to be  $G = 6$  (the

---

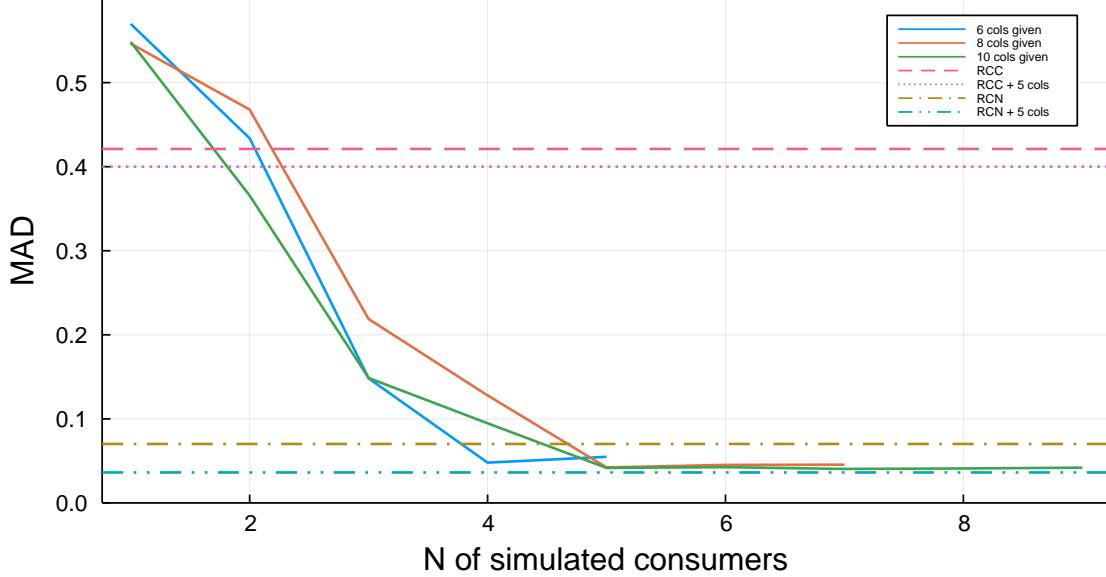
<sup>16</sup>When including  $L = \{5, 6, 8, 10\}$  rows, we use the same ordering of products and try to select the top product from each category.

<sup>17</sup>Figure 1 does not report (or estimate) cases where the rank of the matrix  $I$  exceeds the number of observed rows  $L$ .

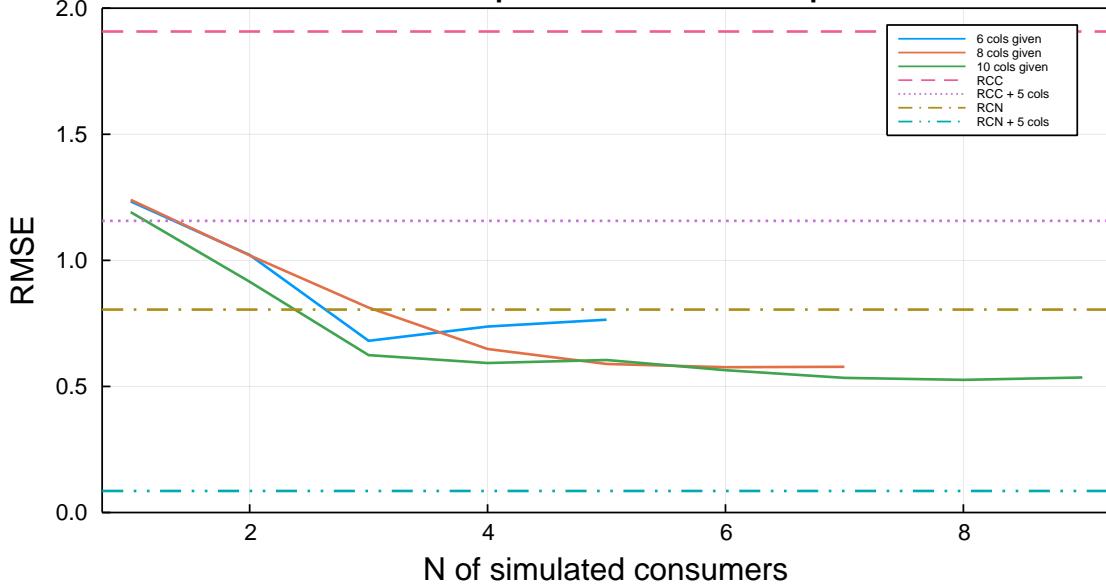
number of nests) and the approximate rank (nuclear norm) to be 2.33. For the RCC DGP the approximate rank 4.5. This helps explain why our semi-parametric approximation of rank  $I = 4$  or  $I = 5$  performs so well—the underlying matrix has a low-rank structure which makes it amenable to our approximation.

Figure 1: Out-of-sample fit for nested logit: Parametric and Semi-parametric Models

### Out- of- sample MAD Comparison

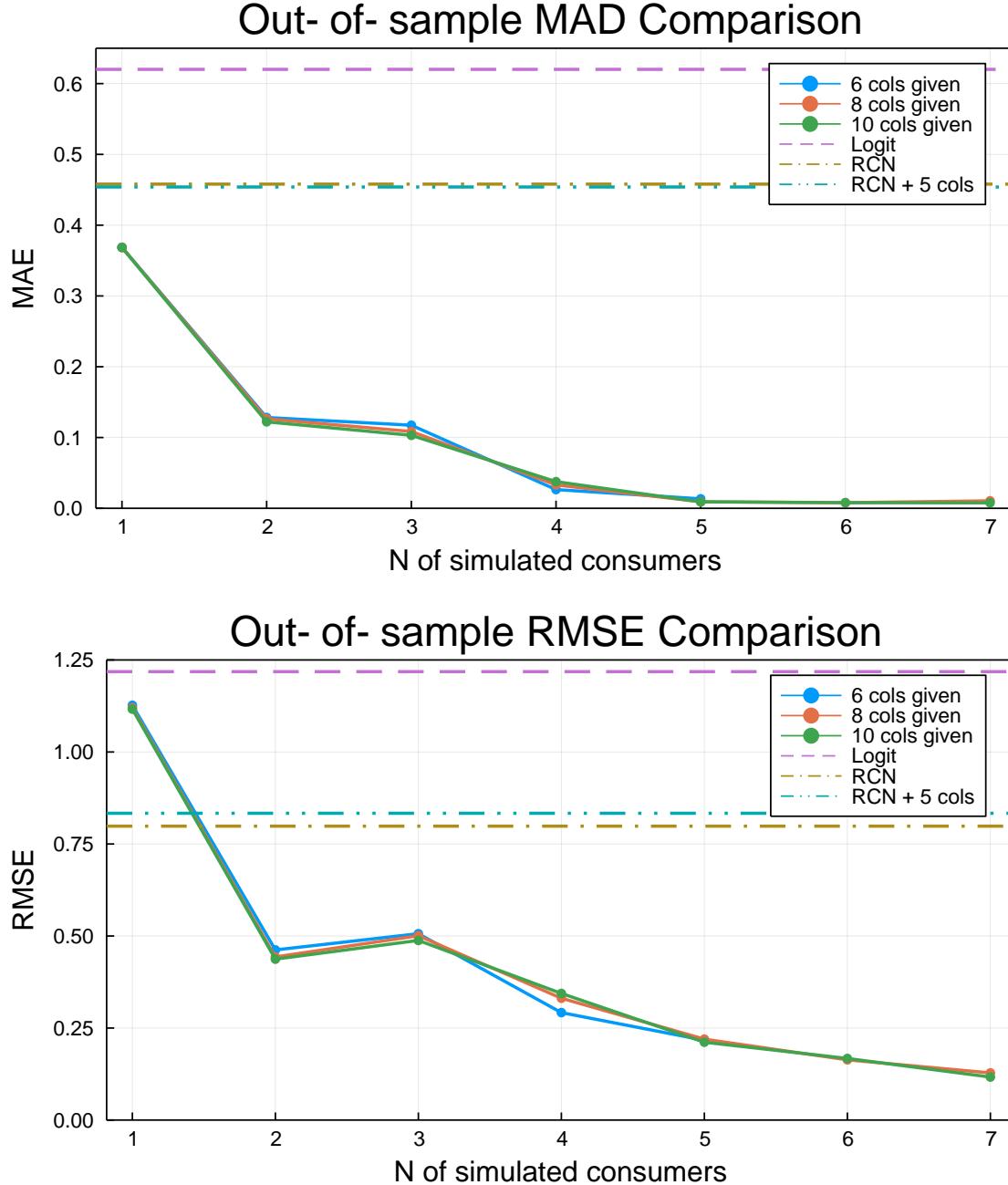


### Out- of- sample RMSE Comparison



Notes: The data generating process is a nested logit model with six categories (rank  $G = 6$ , nuclear norm 2.33). Parametric models (Logit, RCC, and RCN) are fit on full sample  $T = 250$  with or without additional moments based on (five) observed rows of diversion matrix. Semi-parametric model (CMS) fit only on shares  $\mathcal{S}(\mathbf{x}, \theta_0)$  and  $L$  rows of diversion matrix  $\mathcal{D}_{l,.}(\mathbf{x}, \theta_0)$ . Out of sample fit is now on remaining  $J - L$  rows where  $L \in \{6, 8, 10\}$ .

Figure 2: Out-of-sample fit for RCC: Parametric and Semi-parametric Models



Notes: The DGP is a random coefficients logit ( $\sigma_0 = 3.52$ ,  $\sigma_{salt} = 1.00$ ,  $\sigma_{sugar} = 1.8$ ,  $\sigma_{nuts} = 0.38$ ) with rank  $X$ , nuclear norm 4.5.

Parametric models (Logit, RCN) are fit on full sample  $T = 250$  with or without additional moments based on (five) observed rows of diversion matrix.

Semi-parametric model (CMS) fit only on shares  $\mathcal{S}(\mathbf{x}, \theta_0)$  and  $L$  rows of diversion matrix  $\mathcal{D}_{l,\cdot}(\mathbf{x}, \theta_0)$ .

Out of sample fit is on remaining  $J - L$  rows where  $L \in \{6, 8, 10\}$ .

#### 4. U.S. Auto Industry Application

In this section, we estimate the full matrix of second-choice diversion ratios for the U.S. auto industry, using consumers' second-choice data from the 2015 MaritzCX survey. MaritzCX is an automobile industry research and marketing firm that surveys recent car purchasers based on new vehicle registrations. The survey includes a question about the cars that the respondents considered but did not purchase. We consider the first listed car as the purchaser's second choice. These data have been previously used by Grieco et al. (2021) (henceforth GMY) to supplement their parametric model with additional moment conditions. We use the 2015 survey data, which encompasses 53,328 purchases.

For market shares,  $\mathcal{S}$ , we use the 2015 sales of all U.S. vehicles matched with models observed in the survey data and a market size of 20,765,000 as per GMY's methodology. Importantly, we do not make any assumptions to transform these second choices into diversion ratios; we treat them as data.

The GMY model incorporates elements of Berry et al. (1995), Petrin (2002), and Berry et al. (2004a), and represents the state-of-the-art in terms of BLP-style demand estimation. It includes a discrete-choice problem, a firm profit-maximization problem, and micro moments based on demographic information and the second-choice data presented above. Consumer utility is determined by a linear index of car characteristics (e.g. price, footprint, segment).<sup>18</sup> It also accounts for year-to-year variation in the utility of the outside good and average unobserved quality of new cars. Consumer heterogeneity is generated by interacting household characteristics and unobserved preferences with car attributes, allowing for different substitution patterns by demographics. GMY compute both first-choice probabilities (market shares) and second-choice shares conditional on the first choice (diversion matrix) using this consumer-choice model. The supply model assumes simultaneous multi-product pricing given a constant marginal cost function using a Bertrand-Nash assumption.

---

<sup>18</sup>A complete list of characteristics entering the utility function includes: price, footprint, horsepower, mpg, height, curbweight, number of trims, years since last redesign and dummies for van, SUV, truck, luxury, sport, electric, european brand, US brand, new release. Some of these characteristics are interacted with demographic variables such as income, age, rural status and family size.

To properly compare our model with GMY, we must first understand how the latter model constructs micro-moments related to second-choice data. GMY measures correlations in car characteristics between observed first and second choices from survey data and matches them using the same correlations implied by their model. This differs from our estimator, as GMY matches summary statistics instead of the matrix of second-choice probabilities. Our model may have an unfair advantage in that we are assessing performance on its ability to fit second-choice probabilities, while the GMY model is estimating parameters on a much larger set of data while trying to fit additional moments. Still, this serves as a useful benchmark, as we should think of GMY as the state-of-the-art and the best one can hope to do with a BLP-type model.

#### 4.1. Cross-validation exercise

We utilize a cross-validation method to determine the rank of our matrix completion estimator. The 318 products (car models) are randomly split into 20 folds, and we estimate our model 20 times, leaving out one fold each time. To assess the out-of-sample fit on the omitted diversion ratios, we use Mean Absolute Deviation (MAD) and Root Mean Squared Error (RMSE). In Figure 3, we present the performance of our model for various rank choices and compare it to the diversion estimates obtained with the state-of-the-art parametric model used in GMY.

### Cross-validation Results

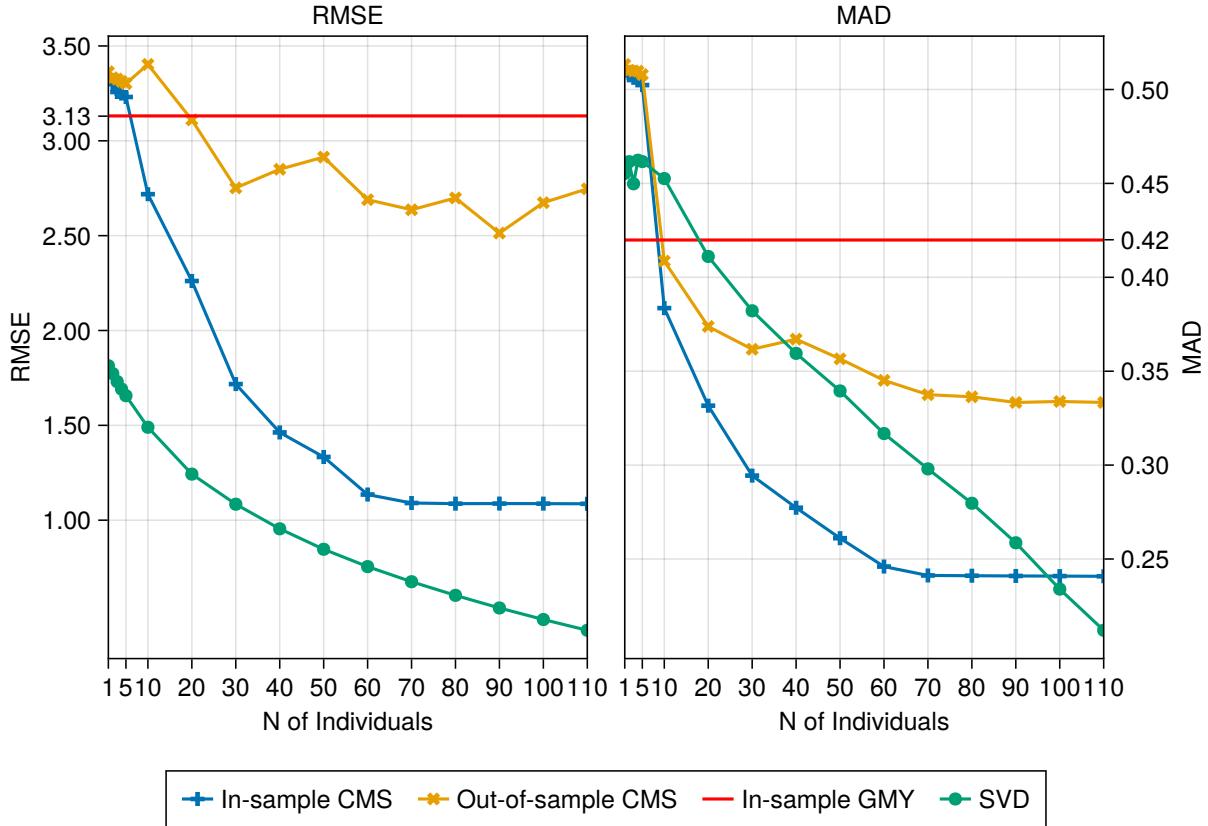


Figure 3: Out-of-sample fit on  $\mathcal{D}$  for different rank  $I$

Cross-validation exercise selects the model with lowest out-of-sample prediction error:  $I = 90$

Our semiparametric model shows superior in-sample performance compared to its parametric counterpart, even with a relatively low rank. Out-of-sample, our model outperforms the parametric model starting from a rank of  $I = 20$  for RMSE fit, and  $I = 10$  for MAD fit. It is important to remember here that the parametric model used in GMY employs a vast amount of data while our estimator uses only second choices and aggregate sales for a single year. According to our cross-validation exercise, our preferred specification is a rank of  $I = 90$ . In the following subsections, we will present results based on estimating this specification, as well as  $I = 1$  (logit) and  $I = 30$  as it could be a local minimum if the researcher were to stop the search at a lower rank than we did.

Figure 4 shows a comparison between our preferred specification and the data, with an additional matrix (far right) displaying the error that remains from prediction. In the following plots, we adopt

the visual convention that the product  $j$  removed corresponds to a column, and each cell represents diversion from  $j$  to  $k$ . The blocks that show up on the diagonal of the diversion matrix are clusters of high substitution among products. They are organized in segments with respect to their size, starting from convertibles (top left) to cargo vans (bottom right). This organization of the data is purely visual and does not impair the performance of our estimator.

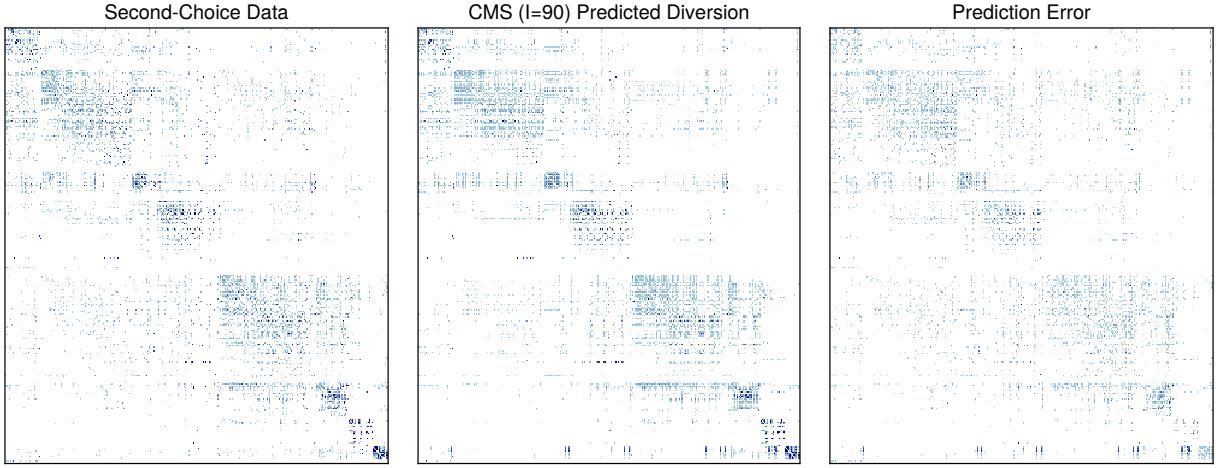


Figure 4: Side-by-side comparison between data, predicted diversion and absolute error

As is generally the case with image processing (another field that uses matrix completion methods), prediction errors tend to occur along “edges” of the image: in our case, products for which diversion is high relative to surrounding products.

#### 4.2. Comparison of implied diversion

Figure 5 presents a comparison between the second-choice survey data and the implied diversion from three models: GMY, our preferred specification, and a logit model. Visually, our semiparametric model with a rank of  $I = 90$  is the closest to the data used for estimation. GMY does a good job of identifying “groups” of vehicles with strong substitution (visible as blocks on the diagonal—mostly because it includes normally distributed random coefficients on product categories such as: sedans, SUV’s, Pickup Trucks, CUV’s, etc.). One challenge for GMY (and most parametric logit models) is that it produces logit-like substitution within each category where the most popular SUV is the best substitute for other SUV’s (visually this appears as a gradient within each block). At the same time, because of the logit error, it also predicts too much substitution to the most popular overall products such as the Camry, Accord, and F-150 even from different groups. This

phenomenon is also observed in other contexts such as the vending experimental data (presented in the next section), although in this case, GMY does a significantly better job, in part because the model has quite a few parameters, and categories do a good job at explaining substitution. The main challenge of these parametric models is that while they over-predict substitution to the field, they under-predict substitution to the closest substitutes.

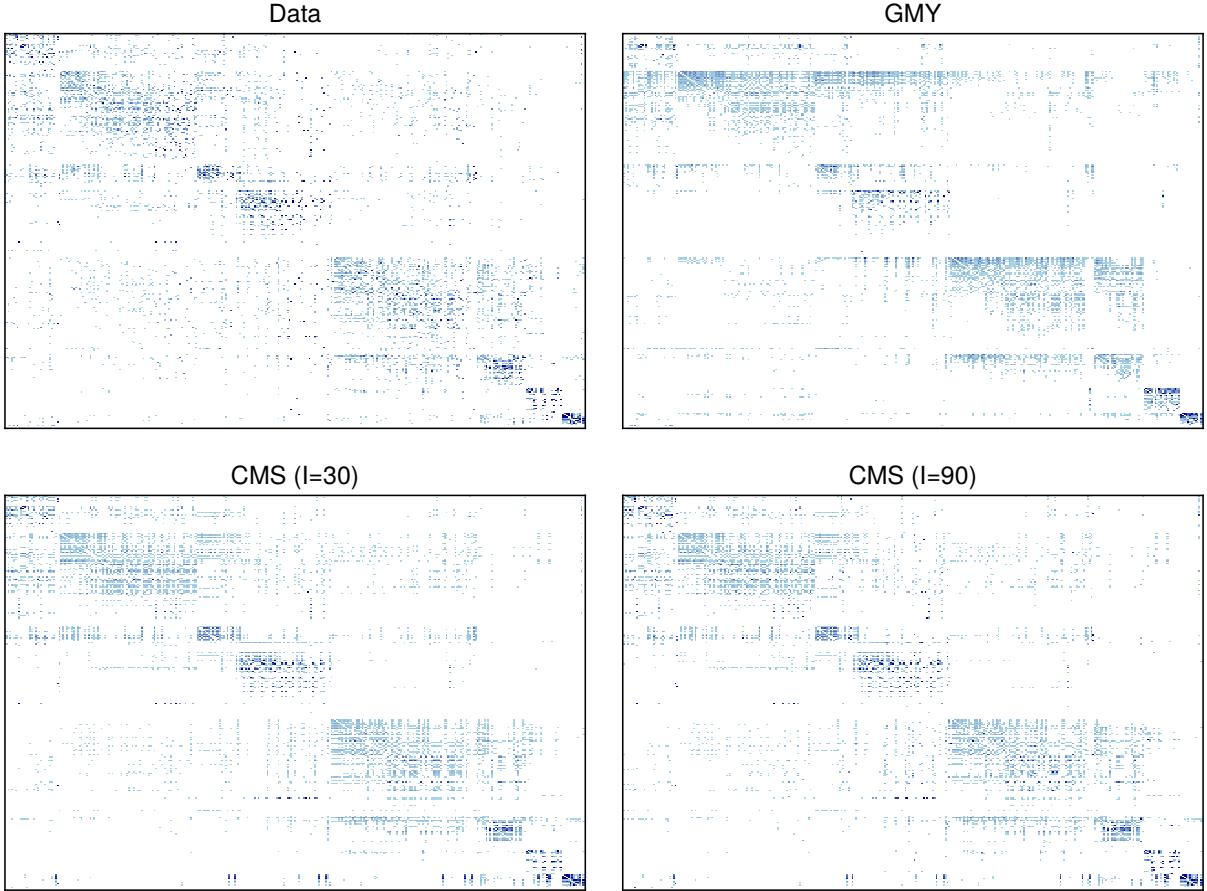


Figure 5: Estimated Diversion Matrices

Visually, our preferred specification matches the sparsity patterns of the data quite closely, while lower rank models and GMY tend to show “fat tails” in that second-choice diversion off the diagonal is higher than in the data.

Our model is capable of rationalizing high second-choice diversion ratios within groups of highly substitutable products and zero diversion with products outside the groups. To illustrate this fact, we compare implied diversion from multiple estimators for three vehicles: the Honda Accord, the Ford F-Series (the best selling vehicle in the U.S.), and the Mercedes-Benz Sprinter Van. The

first vehicle displays relatively spread out second-choice diversion to other models in its category, while the second displays diversion concentrated between only a few vehicles. Finally, the third vehicle has one substitute with two thirds of recorded substitution. Tables 1, 2 and 3 show that our preferred specification is able to rationalize these three different substitution patterns with a low error rate compared to the parametric model in GMY or a simple logit model.

Model	Raw	Logit	CMS I=30	CMS I=90	GMY
Subaru Legacy	10.27	1.01	8.05	8.21	1.3
Toyota Camry	9.1	0.84	6.7	6.85	9.48
Acura Tlx	6.07	0.71	1.83	2.07	0.46
Honda Civic	5.97	0.91	2.9	2.75	3.89
Mazda Mazda6	5.68	0.52	4.77	4.87	1.32
Volkswagen Passat	4.01	0.74	4.34	4.41	1.22
Nissan Altima	3.52	0.6	3.87	3.89	7.22
Hyundai Sonata	3.52	0.68	5.61	5.68	5.09
Volkswagen Jetta	3.33	0.97	4.4	4.23	1.48
Mazda Mazda3	2.15	1.08	2.08	1.61	1.49
Toyota Corolla	1.96	0.71	2.46	2.32	4.66

Table 1: Top Substitutes: Honda Accord

Model	Raw	Logit	CMS I=30	CMS I=90	GMY
Ram Pickup	24.59	1.36	23.38	23.37	19.4
Gmc Sierra	20.29	1.28	21.0	21.02	17.27
Chevrolet Silverado	15.62	1.21	16.73	16.75	33.62
Toyota Tundra	12.98	0.76	12.69	12.69	2.29
Toyota Tacoma	6.31	1.13	3.6	3.62	2.83
Chevrolet Colorado	4.64	1.08	3.37	3.38	2.87
Gmc Canyon	2.3	0.62	1.71	1.73	1.02
Nissan Frontier	1.63	0.67	0.83	0.84	0.61
Jeep Wrangler	1.59	0.62	0.96	0.81	0.06
Nissan Titan	0.7	0.07	0.79	0.81	0.18
Ford Explorer	0.63	0.4	0.05	0.03	0.71

Table 2: Top Substitutes: Ford F-Series

Model	Raw	Logit	CMS I=30	CMS I=90	GMY
Ford Transit Wagon	66.67	0.19	47.63	66.19	0.04
Ram Promaster	16.67	0.02	0.0	16.19	1.76
Ford Transit Connect	8.33	0.18	7.33	7.85	0.01
Nissan Nv	8.33	0.17	29.96	7.58	6.52
Mini Cooper	0.0	0.45	0.0	0.0	0.0
Volkswagen Beetle Ii Cabrio	0.0	0.25	0.0	0.0	0.0
Audi A5	0.0	0.28	0.0	0.0	0.02
Mazda Mx-5 Miata	0.0	0.19	0.0	0.0	0.0
Audi S5	0.0	0.15	0.0	0.0	0.0
Porsche Boxster	0.0	0.07	0.0	0.0	0.0
Volkswagen Eos	0.0	0.07	0.0	0.0	0.0

Table 3: Top Substitutes: Mercedes-Benz Sprinter Van

Figure 6 presents a comparison of in-sample performance for various ranks using different fit measures: RMSE, MAE, the fraction of correctly predicted top 10 substitutes, and the fraction of correctly predicted pairwise comparisons. We selected the latter two measures because they represent important features of the diversion matrix, although they are not directly targeted by either class of models in estimation. The % Correct Top 10 fit measures how many top 10 substitutes (unranked) were correctly predicted (0/1) on average across products. The pairwise comparisons ask for two substitutes  $k$  and  $k'$  whether  $\mathbb{I}[D_{jk} > D_{jk'}]$  is predicted correctly (0/1) and averages across all  $(j, k, k')$ . As an example, consider second-choice diversion from the Honda Accord in Table 1: the Subaru Legacy dominates the Toyota Camry in the data, which is correctly predicted by both CMS  $I = 30$  and  $I = 90$ , but not by GMY nor the Logit model. The % Correct Pairwise Fit is the average proportion of such correctly predicted pairwise comparisons across all products.

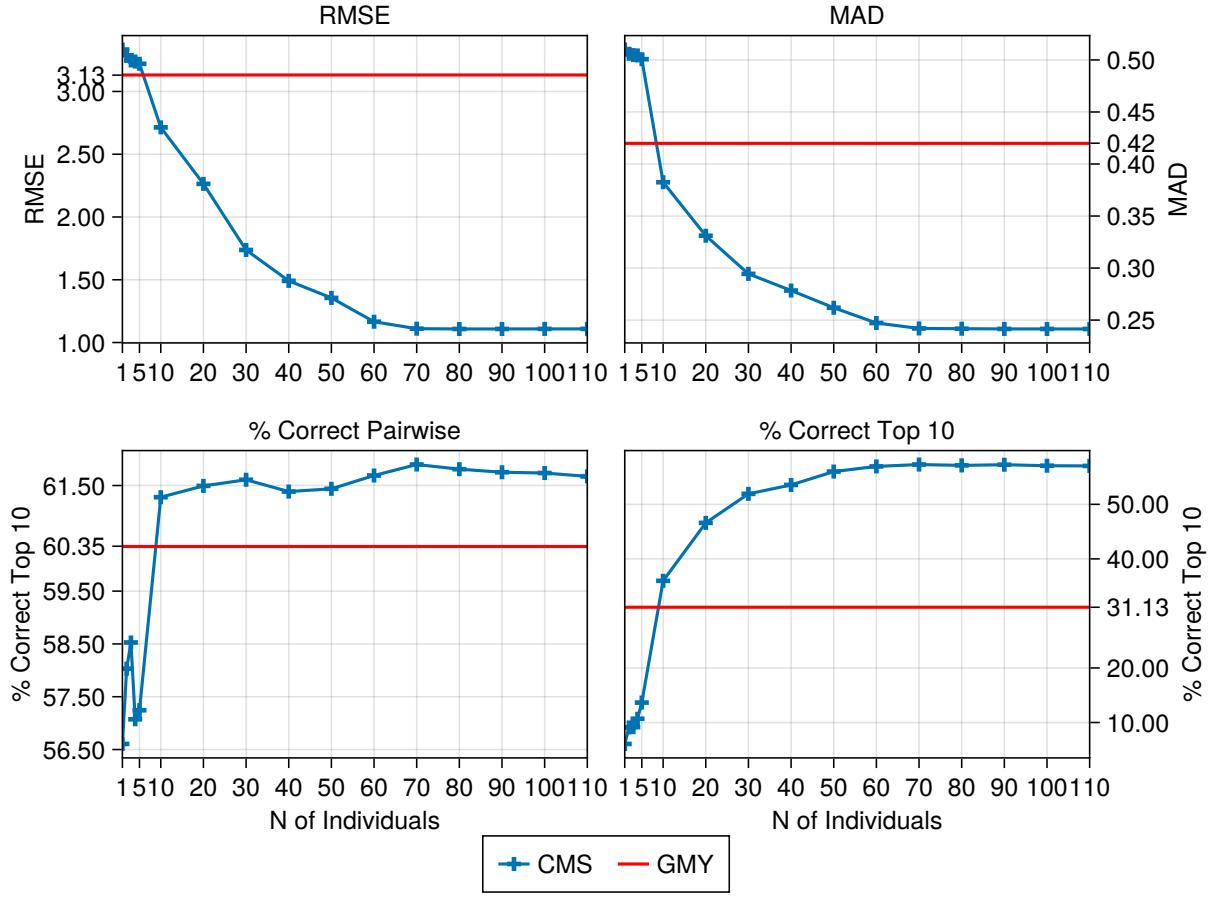


Figure 6: In-sample fit comparison across ranks

Our semiparametric model outperforms its parametric counterpart in all measures starting from relatively low ranks (between 5 and 10). Note: the two bottom fit measures are explained in the previous paragraph.

Our model is able to produce these extreme patterns for a few reasons. One is that it tends to produce highly sparse  $\mathbf{s}_i$  vectors with no substitution to a large number of products at the individual level. The second is that individual diversion ratios are  $D_{j,k,i} = \frac{s_{ik}}{1-s_{ij}}$ , so we can get extremely high rates of substitution when for some individual  $i$ , there is a high value of both  $s_{ij}$  and  $s_{ik}$  (such as in the Mercedes-Benz Sprinter Van example Table 3). This is particularly true if  $s_{ij} = 0$  for most of the other types, such that only a small number of types have non-trivial choice probabilities for the Sprinter Van. However, our model still places some significant restrictions on substitution patterns. In Table 1, we see that the Honda Accord is highly substitutable both the Toyota Camry ( $D_{j \rightarrow k}$ ) and the Subaru Legacy ( $D_{j \rightarrow k'}$ ). For small values of  $I$ , this will imply and without strong

substitution between  $D_{k \rightarrow j}$  (ie: from Camry and Legacy to Accord) or  $D_{k \rightarrow k'}$  (between Legacy and Camry). If this pattern doesn't arise in the data, we can rationalize it by having one type with high  $s_{ij}$  for Camry and Accord (but not Legacy) and another with high  $s_{ij}$  for Legacy and Accord (but not Camry). As the rank  $I$  gets larger, it becomes easier to accomodate this kind of behavior, however this may also be the source of “overfitting” the observed substitution patterns, one advantage of the low-rank structure might be that the model is able to “learn” these kinds of patterns.

## 5. Vending Application

In this section, we apply our estimator for estimating the matrix of diversion ratios from experimental second-choice data using the data from Conlon and Mortimer (2021b). In that data, we ran a field experiment with multiple treatment arms in which we exogenously remove one or two products at a time from vending machines located in office buildings in downtown Chicago. The product removals allow us to measure subsequent diversion to the remaining products without any parametric restrictions on demand.

We provide a complete description of the industry and data in Appendix C.1. We discuss our experimental design in Appendix C.2, and describe our experimentally generated data in Appendix C.3. Our data span the period from June 2007 to September 2008, and the product removal experiments took place between the months of May and October in each year. The data record product-level sales for each of 66 machines for each service visit. We convert the visit level data into a standardized weekly measure by assuming sales are uniform across business days within the visit.

Unlike our previous work in Conlon and Mortimer (2013), which looked at vending machines at Arizona State University, products are rarely out of stock in the MarkVend dataset. Moreover, we observe no price variation within a product across machines or over time. This means that the primary source of variation in  $\mathbf{x}$  is variation in assortment (across time and machines) for less popular products. To augment this variation, we ran a series of experiments where we removed one of the best-selling products in each category for 2.5-3 weeks (Snickers, M&M Peanut, Doritos

	Control	Snickers	M+M Peanut	Famous Amos	Animal Crackers
# Machines	66	56	62	62	62
# Weeks	160	6	6	5	4
# Machine-Weeks	8,525	223	190	161	167
# Products	76	66	67	65	67
Total Sales	700,404.0	19,005.2	16,232.5	14,394.0	13,910.5
—Per Week	4,377.5	3,167.5	2,705.4	2,878.8	3,477.6
—Per Mach-Week	82.2	85.2	85.4	89.4	83.3
Total Focal Sales*		44,026.3	42,047.8	26,113.2	21,578.4
—Per Week		273.5	262.8	163.2	134.0
—Per Mach-Week		5.2	4.9	3.1	2.5

<sup>†</sup> Numbers for Snickers removal. Summary statistics for other removals differ minimally because of different definition of the starting day of the week.

\* Focal sales during the control period. Focal sales during the treatment are close to zero. Any deviation from zero occurs because of the apportionment of service visit level sales to weekly sales.

Table 4: Summary Statistics

Nacho, Cheetos, Zoo Animal Crackers and Famous Amos Cookies). When we remove products, we remove them from all machines at a particular client location (all floors within the same office building).

We summarize the data in Table 4. For each experiment we treat around  $900 - 1,000$  would-be consumers of the focal product (the product we remove in our experiment). The idea is measure these consumers  $\Delta q_j(\mathbf{x})$  and to measure the products to which they substitute  $\Delta q_k(\mathbf{x})$  in order to compute an estimate of the diversion ratio  $\mathcal{D}_{jk} = \frac{\Delta q_k}{\Delta q_j}$ . The main challenge to directly constructing these measures is that both overall sales and product level sales can be quite volatile from week to week. That is, a product may appear to be a substitute simply because its sales increased from one week to the next, or because the machine became busier overall.

To address these concerns, we develop a simple matching estimator based on the assumption that consumers make discrete choices under the follow assumptions: (1) We only use untreated weeks at the same machine and where the substitute product  $k$  is available as controls; (2) Removing a product cannot increase the total sales at the machine, and cannot decrease the total sales by more than the expected sales of the product that is removed  $\Delta Q \in [-\mathbb{E}[q_j], 0]$ . These assumptions enable us to construct estimates of  $\Delta q_j$  and  $\Delta q_k$  for each substitute and each machine-week in our treatment group. We construct aggregate measures of  $(\Delta q_j, \Delta q_k)$  simply by aggregating across

treated machine-weeks.<sup>19</sup>

We make the additional assumption that all products are substitutes and apply a simple non-parametric Bayesian shrinkage estimator, which we estimate via MCMC in `Turing.jl`.<sup>20</sup>

$$\begin{aligned} \Delta q_k &\sim \text{Binomial}(\Delta q_j; \mathcal{D}_{j,k}) \text{ for } k = 0, 1, \dots, K \\ (\mathcal{D}_{j,1}, \dots, \mathcal{D}_{j,K}, \mathcal{D}_{j,0}) &\sim \text{Dirichlet}(\mathcal{S}_1, \dots, \mathcal{S}_K, \mathcal{S}_0) \end{aligned} \quad (14)$$

The use of the Dirichlet prior ensures that all of the estimated Diversion ratios are bounded between  $[0, 1]$  and that  $\sum_k \mathcal{D}_{jk} = 1$ . We center the prior on the observed aggregate shares, but our choice of prior has little impact (we get nearly identical results with a uniform prior  $\frac{1}{J}$ ). We also use a very weak prior (around 3-4 pseudo-observations from the prior and more than 800 from the data). In general the confidence intervals of the estimated diversion ratios are precise (see Appendix C.5).

We report the point estimates for  $\mathcal{D}_{jk}$  for our six product removals in Table 5. We repeated two of our product removals (M&M Peanut and Cheetos) and also ran two experiments where we removed two products at once from the same category (Salty Snacks: Cheetos and Doritos; Chocolate: M&M Peanuts and Snickers). The first column of Table 5 reports the aggregate market shares  $\mathcal{S}_j$ , together with the rows of  $\mathcal{D}_{jk}$ , these for the inputs into our estimator in (6).

## 5.1. Results

In this section, we provide results from our semi-parametric model described in (6), as well as comparisons with popular alternative approaches such as mixed logit models. Unless specified otherwise, the data used to estimate our model consists of only the information in Table 5: a single vector of aggregate shares  $\mathcal{S}$  and a varying number of rows from the diversion matrix  $\mathcal{D}_{j,..}$ . For comparison, we estimate several parametric models RCC (with random coefficients on salt, sugar, nut content, and the constant) and RCN (with nests for each product category) using the full set of observational sales data for several years  $\mathbf{q}(\mathbf{x}_t)$ .

---

<sup>19</sup>Since we've matched on total sales  $Q_t = \sum_{j=1}^J q_{jt}$ , we can identify  $\Delta q_{0t}$  even though  $q_{0t}$  itself is unknown.

<sup>20</sup>We work with the marginal distribution for each substitute  $k$  rather than the joint (Multinomial) distribution because it may be the case that certain substitutes were not available for the full set of "treated" customers and  $\Delta q_j$  is  $k$  specific.

Product	Shares	Snickers	ZAnimal	Doritos	M&M	Cheetos	CC FamA	Cheetos II	M&M II	Salty	Choc
Snyders	4.15	0.02	0.03	0.04	0.04	0.03	0.03	4.01	0.02	0.03	0.01
Cheetos	3.22	0.01	0.02	1.08	0.02		0.03		0.02		0.01
Ruffles	2.8	0.01	0.02	0.95	0.65	6.98	3.85	0.03	0.02	6.77	0.01
Dorito Nacho	2.7	0.02	0.02		0.02	5.9	1.3	0.48	0.02		0.01
Rold Gold	2.56	5.93	5.15	0.68	0.32	0.02	0.02	2.36	1.14	2.66	0.01
Baked	2.39	0.47	3.04	0.02	0.01	0.02	0.03	2.69	5.04	8.85	1.52
Salty Other	2.26	0.01	0.03	0.02	0.01	0.01	0.03	0.02	0.04	0.03	0.02
Sun Chip	2.12	2.23	0.02	5.13	1.98	13.16	26.81	3.62	0.02	0.03	0.01
Cheez-It	2.0	0.01	0.02	3.26	0.02	0.01	1.16	7.25	0.02	2.7	0.61
Jays	1.92	0.01	0.03	0.02	0.01	0.02	0.02	0.01	0.02	0.02	0.01
Frito	1.86	0.01	0.03	0.02	0.03	4.17	0.03	2.05	0.03	6.74	0.99
FritoLay	1.48	1.57	1.21	0.03	1.18	0.72	0.81	4.8	0.02	7.74	0.01
Smartfood	1.32	0.01	1.24	1.88	0.04	0.03	0.02	0.03	0.02	1.51	0.02
Lays	1.01	0.01	0.03	0.02	0.03	0.01	0.03	2.67	0.02	1.09	0.59
Cheetos Flamin	0.62	0.01	0.03	0.04	0.02	2.55	0.17	0.03	0.03	2.36	0.01
Dorito Blazin	0.61	0.02	0.03	2.27	0.03	1.2	8.01	0.55	0.02	3.07	0.17
Popcorn	0.42	0.02	2.81	1.34	0.04	0.32	0.17	0.04	0.02	0.03	0.02
Ritz Bits	0.38	0.02	0.77	0.04	0.03	0.02	0.02	0.03	0.02	0.13	0.01
M&M Peanut	4.14	16.44	9.72	2.69		8.27	0.03	0.04		0.01	
Snickers	3.96		6.93	1.6	17.58	4.58	0.03	4.9	17.74	7.52	
Twix Caramel	2.42	20.25	7.81	2.99	4.96	0.02	0.04	0.05	7.23	0.02	11.64
Raisinets	1.6	0.59	2.12	1.81	5.21	0.01	0.02	0.03	1.61	0.02	3.28
M&M Milk Choc	1.16	5.3	2.6	5.67	10.69	2.57	0.03	0.04	3.05	0.03	8.33
Choc Mars	1.11	1.72	0.02	0.03	0.02	0.01	0.02	0.68	0.03	0.05	6.79
Reeses PB Cups	0.59						0.03	0.66	4.03	2.45	6.57
Butterfinger	0.5	4.13	1.88	6.28	5.41	1.39	1.33	3.67	0.02	0.02	5.22
Choc Herhsey	0.22	2.57	3.68	3.98	1.45	6.32	7.87	3.11	0.31	0.04	2.16
Skittles Original	1.13	0.01	0.02	0.02	0.02	0.03	0.02	0.02	0.03	0.02	0.01
Nonchoc Other	0.78	0.01	0.03	0.02	0.02	0.04	0.04	0.03	0.04	0.21	2.63
Twizzlers	0.33	0.01	0.04	0.06	1.37	2.14	1.26	1.74	4.37	0.02	0.32
ZAnimal Cracker	2.47	1.57		0.03	0.03	0.02	0.02	0.02	0.02	0.03	0.01
CC Fam Amos	2.05	2.87	5.99	4.33	0.57	0.02		0.02	0.01	0.01	2.48
Ruger Wafer	1.72	1.42	0.03	0.01	0.03	0.03	0.03	1.22	0.01	0.92	0.77
Grandmas CC	1.09	0.01	0.02	0.03	0.03	0.02	3.37	0.04	1.78	2.01	0.88
Rasbry Knotts	0.7	0.01	2.38	0.03	0.02	0.03	4.07	3.54	1.65	0.02	0.08
Choc Fam Amos	0.4	1.75	2.45	1.98	1.16	0.03	3.54	1.12	1.86	0.03	0.38
Nabisco	0.39	0.02	2.0	5.71	1.63	1.67	2.52	0.02	0.57	0.02	1.18
Pop-Tarts	1.74	0.01	0.01	0.04	0.02	0.03	0.02	0.02	0.03	0.01	1.27
Rice K Treats	0.27	1.18	3.63	3.3		1.08	0.02	5.85	2.62	2.59	1.94
Nature Valley	2.0	0.02	0.03	0.07	0.01	0.02	0.03	7.11	4.78	1.98	0.01
Planters	1.92	3.96	6.15	9.33	8.99	8.3	11.57	0.03	2.36	0.01	2.02
KarNuts	1.7	0.01	2.55	0.03	0.02	0.02	0.04	0.11	5.91	2.49	1.23
Farleys Fruit Snax	1.04	0.01	0.03	0.02	0.75	0.16	0.03	3.98	0.02	0.08	0.41
Cherry Fruit Snax	0.37	0.01	0.02	0.02	0.02	0.01	0.04	0.03	0.11	0.04	0.01
Cliff	0.28	1.34	1.48	5.94	1.68	0.04	0.03	0.05	0.03	0.03	0.01
Outside Good	30.12	24.42	23.86	27.13	33.83	27.96	21.41	31.21	33.21	35.57	36.33

Table 5: Aggregate Shares and Estimated Second Choice Diversion (Ground Truth)

Each column corresponds to a different product removal experiment.

Salty: removes both Cheeto and Doritos.

Choc: removes both Snickers and M&M Peanut.

Our notion of fit is the same as in Section 3, we withhold a row of the diversion matrix  $\mathcal{D}_{j,\cdot}$  from estimation and compare the out-of-sample predictions:  $\|\mathcal{D}_{j,\cdot}(\mathbf{x}) - D_{j,\cdot}(\mathbf{x}, \hat{\theta})\|$ .

## 5.2. Cross-validation and out-of-sample fit

To avoid potential over-fitting, we choose the complexity of our model  $I$  in (6) by cross-validation. For each fold, we estimate the model while omitting a single experiment (such as Snickers or Animal Cracker removal) and measure the out-of-sample goodness-of-fit on the omitted diversion ratios. Figure 7 compares two out-of-sample measures of fit (MAD and RMSE) between our model and parametric models. For our model, the whiskers represent the 25th and 75th percentile of fit across all folds. In general, cross-validation tends to select relatively small models, here a model of rank  $I = 2$  or  $I = 3$ .

Our model in (6) performs substantially better than plain or mixed logit fit to the same dataset. In particular, we compare our method to a simple logit model, as well as a random coefficients model on product characteristics (RCC) and on nests dummies (RCN). Note that while our estimator uses only aggregate shares and limited diversion data, the parametric models include all machine-week observations for all products (a few years worth of machine-week specific data). Our preferred specifications of  $I = 2$  or  $I = 3$  tend to outperform comparable parametric methods for both MAD and RMSE. Figure 7 also reports a variant of model that penalizes the  $\sum_i \pi^2$ . This tends to make the weights on types more equal, but otherwise has little impact. We produce a full version of these results in Appendix E.

We provide a visual comparison of the estimated diversion matrices for  $I = \{2, 3, 4\}$  Figure 8. Each subfigure is a heatmap of the transposed diversion matrix  $D^T(\theta)$ . The first quadrant shows the observed data from our six experiments (now as columns). The remaining quadrants show how we “complete the matrix” of substitution patterns with different choices for  $I$ . The bottom right plot for  $I = 4$  shows signs of “overfitting” (it goes from looking smooth to more pixelated).

## 5.3. Comparison of implied diversion

Table 6 and Table 7 provide the in-sample nonparametric diversion estimates ( $\mathcal{D}_{jk}$ ) for two of the ten experiment arms, and compare them to diversion implied by parametric estimates  $D_{jk}(\theta)$  and

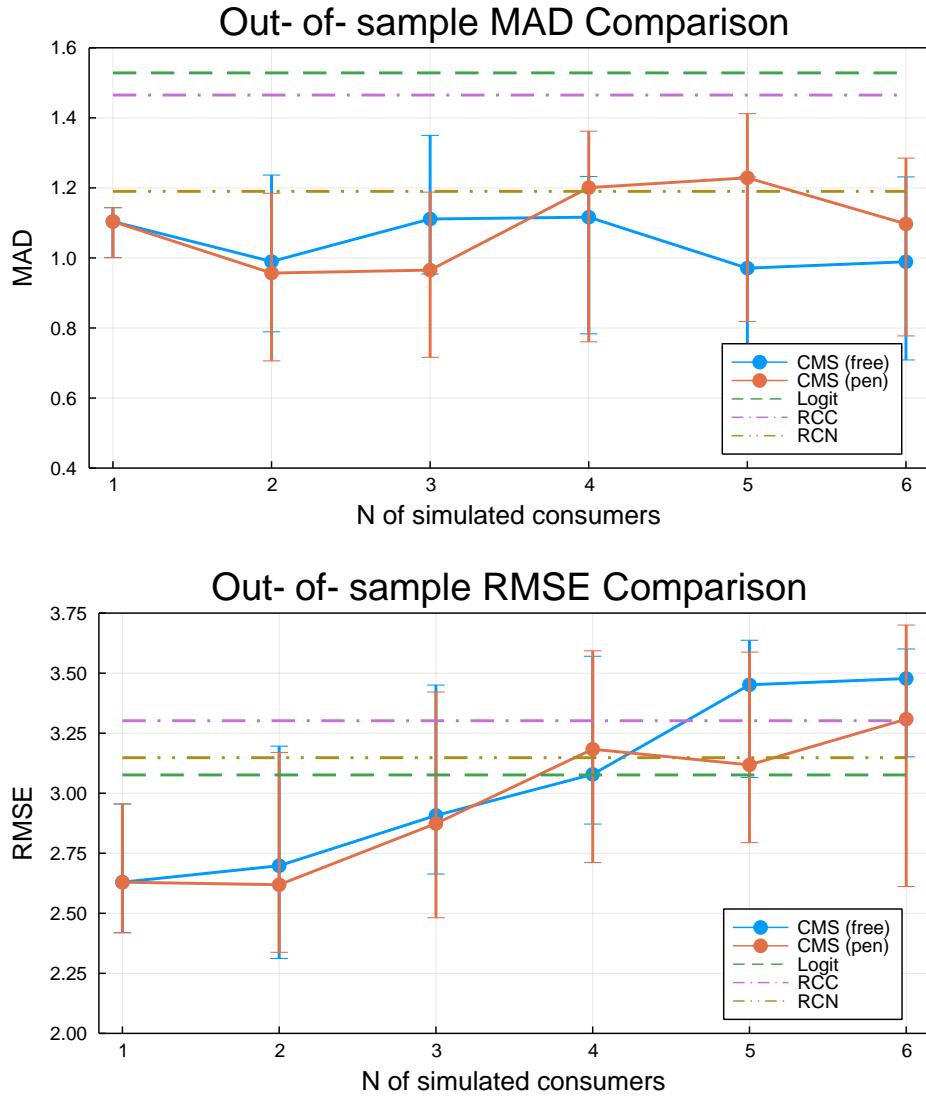


Figure 7: Comparison of RMSE and MAD

Notes: Parametric Models are estimated on full sample.

Semiparametric model (CMS) is estimated on shares and all but one observed column of diversion matrix. (Out of sample fit reported).

RCC: Independent Normal Coefficients on (Sugar, Salt, Nut Content).

RCN: Independent Normal on Category Dummies (Salty Snack, Chocolate Candy, Non-Chocolate Candy, Cookies, Pastry, Other). [See Table X].

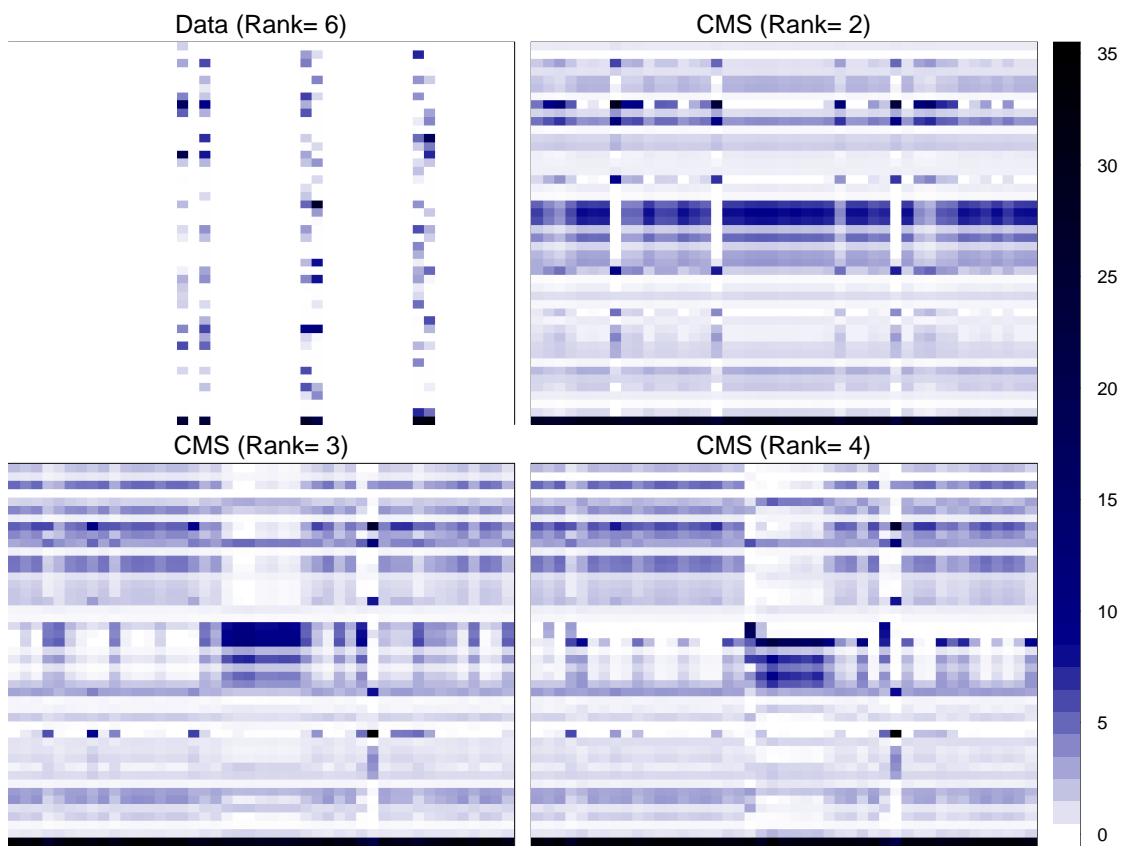


Figure 8: Predicted Diversion Matrices: Semiparametric Models

Products are sorted by category, and then by descending share within each category.

our preferred semiparametric specifications ( $I \in \{2, 3\}$ ). The most notable feature of our model is the ability to assign high diversion ratios to both the outside good and top substitutes, while parametric models tend to have flatter substitution profiles as well as underestimated outside good diversion.

Product	Shares	Nonparam	Logit	RCC	RCN	CMS(I=2)	CMS(I=3)
Outside Good	30.12	23.86	22.93	23.37	20.2	29.27	25.44
M&M Peanut 1.74 oz	4.14	9.72	3.34	3.74	2.43	5.94	7.31
Twix Caramel	2.42	7.81	2.44	2.91	1.79	6.72	9.0
Snickers 2.07oz	3.96	6.93	3.19	3.53	2.33	7.35	8.08
Planters (Con)	1.92	6.15	2.19	1.85	1.23	3.99	5.21
Choc Chip Famous Amos	2.05	5.99	1.66	1.87	8.23	0.31	4.74
Rold Gold (Con)	2.56	5.15	3.01	1.29	1.68	2.16	2.08
Choc Herhsey (Con)	0.22	3.68	1.31	1.63	1.0	2.07	2.42
Rice Krispies Treats 1.7oz	0.27	3.63	0.99	1.07	0.69	2.49	1.74
Baked (Con)	2.39	3.04	2.09	1.82	1.17	2.32	1.15
Popcorn (Con)	0.42	2.81	1.0	0.71	0.56	0.58	0.56

Table 6: Top Substitutes: Zoo Animal Crackers Removal

Product	Shares	Nonparam	Logit	RCC	RCN	CMS(I=2)	CMS(I=3)
Outside Good	30.12	36.33	24.02	24.64	27.99	34.78	33.62
Twix Caramel	2.42	11.64	2.56	2.93	5.31	8.56	11.97
M&M Milk Chocolate	1.16	8.33	2.09	2.63	4.43	5.73	7.1
Choc Mars (Con)	1.11	6.79	1.76	2.11	3.68	1.71	2.22
Reeses Peanut Butter Cups	0.59	6.57	1.84	2.78	3.83	3.48	5.12
Butterfinger	0.5	5.22	1.22	1.71	2.49	3.75	4.2
Raisinets	1.6	3.28	1.67	2.12	3.48	2.38	2.92
Nonchoc Other (Con)	0.78	2.63	1.45	1.81	1.33	0.69	0.74
Choc Chip Famous Amos	2.05	2.48	1.74	1.79	1.23	0.0	0.21
Choc Herhsey (Con)	0.22	2.16	1.37	1.69	2.98	1.85	2.15
Planters (Con)	1.92	2.02	2.3	4.16	1.52	4.1	5.13

Table 7: Top Substitutes: Snickers and MMs Peanut Removal

#### 5.4. Semiparametric model estimates

In Figure 9 we report the individual shares  $s_{ij}$  and weights  $\pi_i$  for our semiparametric estimator with different levels of complexity (rank)  $I \in \{1, 2, 3, 4\}$ . These shares and weights correspond to the estimates of our model. As we increase the rank of our semiparametric approximation, we see how the types change. The  $I = 1$  approximation is an augmented logit, as it also fits the second-choice data in addition to the aggregate shares. The  $I = 2$  case starts to provide clear separation between individuals who like salty snacks and individuals who like chocolate candy. When we add a third

“type” to the model, we see that the types become increasingly sparse (0s in the table) and have more extreme preferences for particular goods (for example Fritos, Reeses Peanut Butter Cups). When we add a fourth type, we see that they have even more sparse preferences.

Finally, in Figure 10, we present our estimated  $(J + 1) \times (J + 1)$  diversion matrix using our rank  $I = 2$  approximation, and compare it with the underlying nonparametric diversion as well as comparable parametric models. We shade in pairs of products with high diversion ratios and include diversion to the outside good. Clusters of similar products emerge visually (Chips and Salty Snacks, Chocolate Candy, etc.). We are able to infer relatively large amounts of substitution between similar products such as Twix and Raisinets even though neither product is removed during our experiment, and we don’t use any information about the characteristics of either product. Additionally, even though products like Snickers and Twix are popular products overall, we estimate little to no diversion from salty snacks such as Fritos or Lay’s Potato chips to these products. This substitution pattern would be possible in a nested logit, but we haven’t provided any information about product categories either.

## 6. Extensions

### 6.1. Multiple Product Removals

In some settings, we might have data on diversion ratios from multiple product removals. For example, in the vending machine experiment described in this paper, two treatment arms removed two products simultaneously: Snickers and M&M’s in one arm and Doritos and Cheetos in the other. We show in Appendix A.1 that the diversion ratio to remaining product  $k$  from removing multiple products  $j \in \mathcal{J}_{\text{rem}}$  with total choice probability:  $s_{i,\mathcal{J}_{\text{rem}}}(\mathbf{x}) \equiv \sum_{j \in \mathcal{J}_{\text{rem}}} s_{ij}(\mathbf{x})$  has a structure similar to (5), and is given by:

$$D_{\mathcal{J}_{\text{rem}},k}(\mathbf{x}) = \frac{1}{s_{\mathcal{J}_{\text{rem}}}(\mathbf{x})} \sum_{i=1}^I \pi_i s_{ik}(\mathbf{x}) \left( \frac{s_{i,\mathcal{J}_{\text{rem}}}(\mathbf{x})}{1 - s_{i,\mathcal{J}_{\text{rem}}}(\mathbf{x})} \right). \quad (15)$$

Our semiparametric estimator (6) can allow for these cases simply by minimizing the squared difference between observed and predicted substitution with multiple product removals.

	Model/Rank:		I = 1		I = 2		I = 3			I = 4			
	Weight on individual:		100.0%	81.2%	18.8%	62.8%	31.4%	5.8%	73.5%	24.1%	2.4%	0.02%	
	Product	Logit Sj	i = 1	i = 1	i = 2	i = 1	i = 2	i = 3	i = 1	i = 2	i = 3	i = 4	
SALTY SNACKS	Snyders (Con)	2.21	0.7	0.64	0.7	0	0	2.17	0	2.38	0	0.23	
	Cheetos	2.52	0.5	0	0	0.24	0.25	0.8	0.31	0	0	0	
	Ruffles (Con)	0.98	1.88	0.98	4.82	0	2.53	4.84	0.03	5.05	2.54	0	
	Dorito Nacho	2.05	0.98	0.9	1.22	0	0	0	0	0	0	0	
	Rold Gold (Con)	0.94	1.86	2.35	0	2.5	0.13	1.4	0.65	1.95	0.14	4.41	
	Baked (Con)	1.99	2.08	2.49	0.4	1.36	0	3.94	2.35	3.76	0	0.43	
	Salty Other (Con)	2.78	0.22	0.29	0	0.05	0	0.59	0	0.71	0	0.3	
	Sun Chip	1.81	4.76	0	22.96	0	19.29	5.19	0.09	5.71	18.7	0	
	Cheez-It	1.77	1.47	1.06	2.67	0	0.91	3.9	0	4	0.93	0.43	
	Jays (Con)	1.48	0.17	0.23	0	0.04	0	0.47	0	0.57	0	0.24	
	Frito	2.03	1.41	1.28	1.65	0	0	4.55	0	4.7	0	0.04	
	FritoLay (Con)	1.49	1.71	1.57	1.94	0.15	0.36	4.49	0.53	4.5	0.45	0.34	
	Smartfood	1.51	0.56	0.67	0	0.31	0.16	1.02	0	1.17	0.18	0.69	
	Lays	1.45	0.54	0.56	0.31	0	0	1.62	0	1.68	0	0.3	
	Cheetos Flamin	0.96	0.55	0.5	0.55	0	0	1.83	0	1.92	0.03	0	
	Dorito Blazin	1.45	1.47	0.01	6.47	0	5.77	1.82	0	2.06	5.62	0	
	Popcorn (Con)	2.06	0.51	0.63	0	0.61	0.33	0.34	0.22	0.54	0.36	0.8	
	Ritz Bits	0.51	0.16	0.22	0	0.17	0	0.23	0	0.32	0.03	0.26	
CHOCOLATE CANDY	M&M Peanut	3.21	4.78	6.46	0	8.95	0	1.18	16.65	0	0	0	
	Snickers	3.53	6.08	8	0	9.75	0.7	0	14.91	0	0.11	0.13	
	Twix Caramel	2.29	5.19	7.32	0	11.03	0	0	4.04	0	0	18.6	
	Raisinets	1.47	1.47	2.03	0	2.67	0	0.21	2.74	0.03	0.05	2.06	
	M&M Milk Choc	1.8	3.63	4.9	0	6.47	0	0.67	5.31	0.36	0.04	7.24	
	Choc Mars (Con)	2.13	1.03	1.46	0	2.03	0	0.11	0	0.31	0	4.69	
	Reeses PB Cups	1.68	1.62	2.98	0	4.68	0	0.36	2.95	0.25	0	7.38	
	Butterfinger	1.1	2.72	3.21	0.8	3.69	1.13	1.67	2.06	1.73	1.16	5.62	
	Choc Herhsey (Con)	1.22	2.87	1.58	7.2	1.64	5.92	2.92	0.91	3.23	5.77	2.53	
NONHOC. CANDY	Skittles Original	1.03	0.12	0.18	0	0.03	0	0.36	0	0.43	0	0.17	
	Nonchoc Other (Con)	1.06	0.41	0.59	0	0.65	0	0.32	0	0.41	0	1.54	
	Twizzlers	1.66	1.16	1.2	0.86	0.9	0.59	1.71	1.96	1.41	0.66	0	
COOKIES	ZAnimal Cracker	1.9	0.29	0.35	0.14	0.33	0.39	0	0.16	0	0	0	
	CC Fam Amos	1.58	1.57	0	3.66	0.04	26.42	0	0.23	0	28.73	0	
	Ruger Wafer (Con)	1.6	0.54	0.68	0	0.46	0	0.94	0	1.07	0	1.31	
	Grandmas CC	1.15	0.84	0.46	2.07	0.34	2.21	0.89	0.47	0.92	2.25	0.05	
	Rasbry Knotts	0.68	1.1	0.47	3.18	0.45	2.89	1.12	0.76	1.19	2.87	0	
	Choc Fam Amos	0.91	1.35	1.09	2.12	1.47	2.65	0.38	1.19	0.52	2.63	1.35	
	Nabisco (Con)	1.23	1.44	1.22	2.04	1.24	2.28	1.11	0.99	1.2	2.21	1.44	
PASTRY	Pop-Tarts (Con)	2.42	0.27	0.39	0	0.34	0	0.36	0	0.46	0	0.87	
	Rice K Treats	0.85	2.25	2.64	0.8	2.06	0.16	3.22	2.78	3.03	0.24	1.59	
OTHER	Nature Valley (Con)	2.13	1.42	1.47	1.02	0.42	0	3.54	1.63	3.22	0	0	
	Planters (Con)	1.63	4.81	3.51	9.13	4.39	8.99	2.79	4.91	2.79	8.74	3.1	
	KarNuts (Con)	1.65	1.25	1.68	0	1.71	0	1.05	2.51	0.87	0.01	0.36	
	Farleys Fruit Snax	0.99	0.58	0.57	0.45	0.04	0	1.69	0.16	1.69	0	0.08	
	Cherry Fruit Snax	0.52	0.09	0.14	0	0.05	0	0.25	0	0.3	0	0.12	
	Cliff (Con)	3.91	1.03	1.29	0	1.3	0.51	0.67	0.62	0.82	0.48	2.03	
	Outside Good	25.34	28.58	29.75	22.86	27.43	15.42	33.28	27.87	32.77	15.06	29.27	

Figure 9: Estimated individual Shares  $s_{ij}$

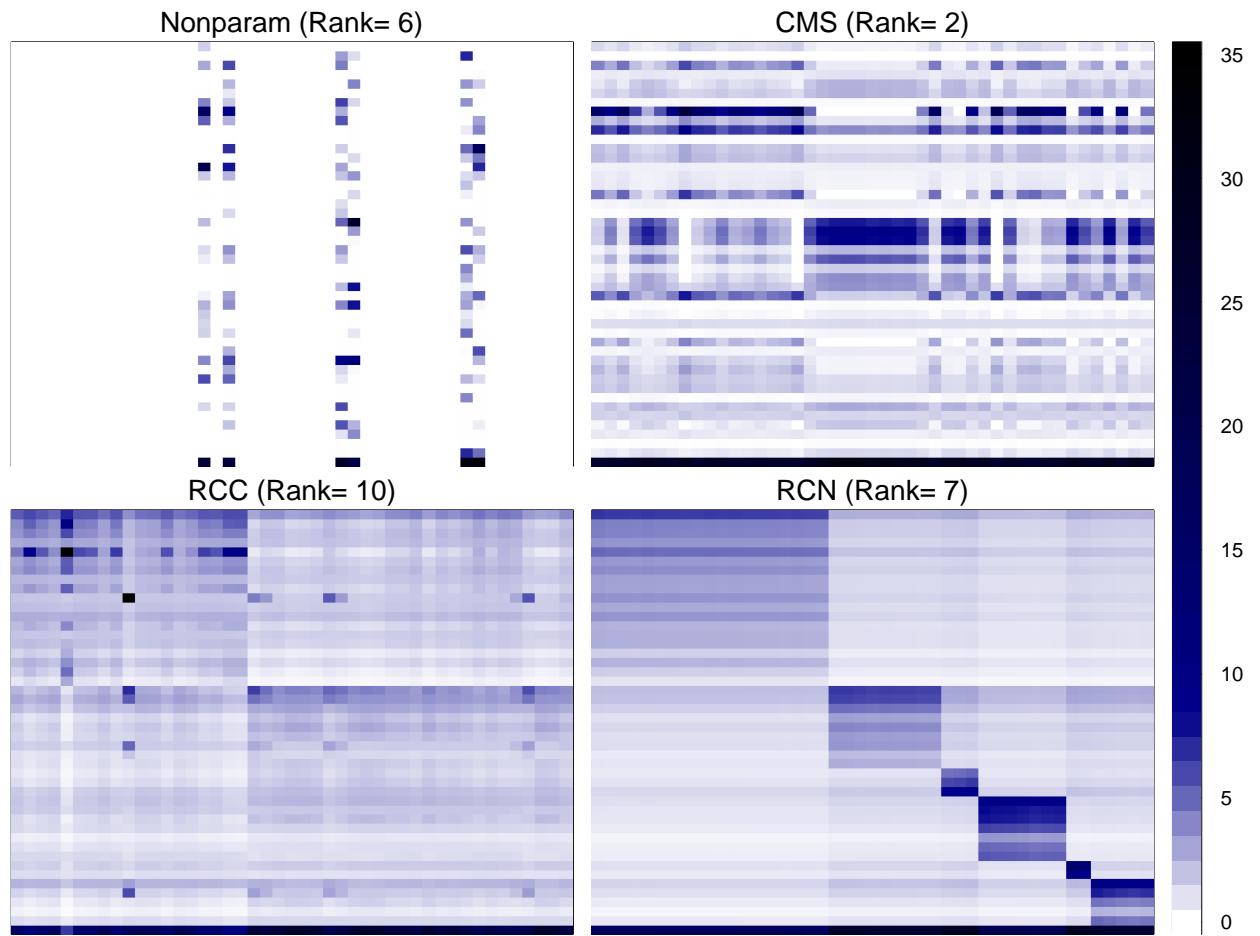


Figure 10: Estimated Diversion Matrices (all models)

## 6.2. Adding Parametric Restrictions

There may be cases where we are interested not only in substitution among existing products, but also what might happen if we were to introduce a new product, or if prices or characteristics were to change. In this case we might want to construct a prediction for  $\mathcal{D}(\mathbf{x}')$  or  $\mathcal{S}(\mathbf{x}')$  at some  $\mathbf{x}' \neq \mathbf{x}$ .

In this case having a parametric structure on characteristics like  $V_{ij}(x_j) = \beta_i x_j + d_j$  would be helpful. We can simply impose an additional set of constraints on (6), and search over  $\beta_i$  and  $d_j$  as well as  $\pi_i$  and  $s_{ij}$ :

$$s_{ij}(\mathbf{x}) = \frac{e^{\beta_i x_j + d_j}}{1 + \sum_k e^{\beta_i x_k + d_k}}.$$

This creates some challenges: the bound  $\text{rank}(D) \leq I$  may no longer be informative; for any fixed  $I$  the fit of the model is likely to be worse. However, we've effectively reduced the number of free parameters from  $(J+1) \times I$  to  $J + (K+1) \times I$  where  $\dim(\beta) = K$ .

The advantage of these additional restrictions is that it enables us to test whether or not the characteristics  $x_j$  span the space of substitution patterns  $\mathcal{D}$ . An obvious choice would be a non-nested model comparison like Rivers and Vuong (2002) which compares the fit of the model with and without parametric restrictions and adjusted for the degrees of freedom.

## 6.3. Other Variation: Prices, Quality, and Characteristics

We could derive a similar result where instead of second-choice diversion, we consider diversion with respect to an infinitesimal change in some characteristic (such as price) where  $\frac{\partial V_{ij}}{\partial z_j} = \beta_i^z$ <sup>21</sup>

$$\frac{\partial s_j}{\partial z_k}(\mathbf{x}) = \sum_{i=1}^I \beta_i^z \cdot \pi_i \cdot s_{ik}(\mathbf{x}) \cdot (1[j=k] - s_{ij}(\mathbf{x})) \quad (16)$$

$$\mathbf{D}_z(\mathbf{x}) = \text{diag}(\mathbf{s}_z)^{-1} \cdot \sum_{i=1}^I \pi_i \cdot \beta_i^z \cdot \mathbf{s}_i(\mathbf{x}) \cdot \mathbf{s}_i^T(\mathbf{x}). \quad (17)$$

---

<sup>21</sup>Here  $\text{diag}(\mathbf{s}_z)^{-1}$  is a diagonal matrix where entries are given by  $\left(\frac{\partial s_j}{\partial z_j}\right)^{-1}$  for a particular characteristic  $z$ . This term row-normalizes the matrix so that  $\sum_{j \neq k} D_{kj} = 1$ .

There are two ways to view (17). The first is that in order to predict elements of  $\mathbf{D}_z(\mathbf{x})$ , we now need an estimate of  $\beta_i^z$  for each type  $i$ . This means that counterfactuals which depend on  $\mathbf{D}_z(\mathbf{x})$  or  $\frac{\partial s_j}{\partial z_j}(\mathbf{x})$  are not identified from our simple estimator in (6) alone. The second is that if we observe part or all of  $\mathcal{D}_z(\mathbf{x})$ , we can use this variation to estimate  $\beta_i^z$ .

#### 6.4. Multiple Markets

Our estimator in (6) is conditioned on a single set of observables  $\mathbf{x}$  so that we observe aggregate shares  $\mathcal{S}(\mathbf{x})$  and some elements of  $\mathcal{D}(\mathbf{x})$  to recover  $\mathbf{s}_i(\mathbf{x})$  and  $\pi_i$ . Implicitly, this means everything is conditional on  $\mathbf{x}$ . If we observed data from multiple markets  $t = 1, \dots, T$  with different  $\mathbf{x}_t$ , we could either estimate separately market by market, or parameterize  $V_{ij}(\mathbf{x}_t)$  and use the parametric structure to pool parameters across markets.<sup>22</sup>

#### 6.5. Identification and Inference

There are different ways to think about asymptotic inference in (6). One approach would be to take the observed elements of  $\dim(\mathcal{D}_{\text{OBS}}) \rightarrow \infty$  and implicitly  $J \rightarrow \infty$  and treat (6) as a GMM estimator. This would be similar to the approach taken in Berry et al. (2004b).

A more practical approach might be to think about (6) as a minimum distance estimator where the observed second-choice diversion ratios are themselves estimated from some sample. This would be the case if we were using a survey of  $n$  individuals to estimate the second-choice diversion ratios. In that case we would need:  $\|\mathcal{D}_{jk}^n - \mathcal{D}_{jk}^0\| \xrightarrow{P} 0$  (the sample estimates converge in probability to the true population second-choice diversion ratios.)

The standard conditions for minimum distance estimators are straightforward to verify for our constrained least squares problem. Compactness of  $\theta \in \Theta$  would be guaranteed by the constraints in (6), all parameters are constrained to the unit interval. We provide the derivatives  $\frac{\partial D_{jk}}{\partial \theta}$  in Appendix A.3. For the derivatives to be bounded we need that  $\mathbf{s}_i$  is non-degenerate and has at least two nonzero elements for each  $i$ . The objective in (6) is quadratic and all other constraints are linear or quadratic in parameters, which should guarantee that  $Q(\theta)$  is twice continuously

---

<sup>22</sup>The former in applications like hospital demand. The latter is effectively the identification approach in much of the rest of the IO literature Berry et al. (1995); Nevo (2001), etc.

differentiable with bounded derivatives.

Identification of the model in (6) is straightforward so long the rank of the observed diversion matrix is sufficiently large:  $\text{rank}(\mathcal{D}) > I$ . This guarantees that the objective function is never equal to zero. We must also have more observed elements of  $\mathcal{D}$  than unknowns  $(J+1) \cdot I$ . Because  $D_{jk}(\theta)$  is a non-convex function of the parameters, at best we can only hope to establish local identification at some values of  $(\mathbf{s}_i, \pi_i)$ . We need a rank condition on the Jacobian (with respect to parameters  $\theta = [\mathbf{s}_i, \pi_i]$ ) for there to be a unique solution in the least-squares sense. A necessary though not sufficient condition is that the vectors  $\mathbf{s}_i$  be linearly independent and  $\pi_i > 0$  strictly for all  $i$ . We provide the derivatives  $\frac{\partial D_{jk}}{\partial \theta}$  in Appendix A.3.

The easiest way to construct confidence intervals for  $D_{jk}(\hat{\theta})$  (or other outputs) is to bootstrap the underlying survey data used to construct  $\mathcal{D}^n$  and re-estimate the model in (6). Even on relatively large problems, estimation only takes a few seconds in our application from Section 5.

## 6.6. Is this true

Consider an  $I \times I$  diagonal matrix  $\Pi_I$  with entries  $\pi_i$  can we write:

$$\mathbf{D}_I = \mathbf{S}_I \cdot \Pi_I \cdot \mathbf{S}_I^{*\top} \text{ where } S^* \text{ has columns } \mathbf{s}_i^* = \frac{\mathbf{s}_i}{1 - \mathbf{s}_i} \text{ and } \mathbf{S} \text{ has columns } \mathbf{s}_i.$$

## 7. Conclusion

We develop a semi-parametric estimator of the full matrix of diversion ratios based on matrix completion methods commonly used in computer science. Our approach uses data on aggregate market shares and (potentially partially observed) second-choice diversion ratios, requires no information on product characteristics, and is computationally easy to estimate. We demonstrate the approach in Monte Carlo simulations and compare it to commonly-used but potentially misspecified parametric models. We also apply the method to two settings: the US automobile market, for which we have diversion on all products in a setting with more than 300 products, and the snack food industry, for which we have experimental second-choice data for a small subset of products in a setting with about 40 products.

## References

- ABADIE, A. AND G. W. IMBENS (2011): “Bias-Corrected Matching Estimators for Average Treatment Effects,” *Journal of Business & Economic Statistics*, 29, 1–11.
- ABALUCK, J. AND A. ADAMS-PRASSL (2021): “What do Consumers Consider Before They Choose? Identification from Asymmetric Demand Responses\*,” *The Quarterly Journal of Economics*, 136, 1611–1663.
- ALLENBY, G. M., N. HARDT, AND P. E. ROSSI (2019): “Chapter 3 - Economic foundations of conjoint analysis,” in *Handbook of the Economics of Marketing*, ed. by J.-P. Dubé and P. E. Rossi, North-Holland, vol. 1 of *Handbook of the Economics of Marketing, Volume 1*, 151–192.
- BARSEGHYAN, L., M. COUGHLIN, F. MOLINARI, AND J. C. TEITELBAUM (2021a): “Heterogeneous Choice Sets and Preferences,” *Econometrica*, 89, 2015–2048.
- BARSEGHYAN, L., F. MOLINARI, AND M. THIRKETTLE (2021b): “Discrete Choice under Risk with Limited Consideration,” *American Economic Review*, 111, 1972–2006.
- BAYER, P., F. FERREIRA, AND R. McMILLAN (2007): “A Unified Framework for Measuring Preferences for Schools and Neighborhoods,” *Journal of Political Economy*, 115, 588–638.
- BAYER, P. AND C. TIMMINS (2007): “Estimating Equilibrium Models Of Sorting Across Locations,” *The Economic Journal*, 117, 353–374.
- BERRY, S. (1994): “Estimating discrete-choice models of product differentiation,” *RAND Journal of Economics*, 25, 242–261.
- BERRY, S. AND P. JIA (2010): “Tracing the Woes: An Empirical Analysis of the Airline Industry,” *American Economic Journal: Microeconomics*, 2, 1–43.
- BERRY, S., J. LEVINSOHN, AND A. PAKES (1995): “Automobile Prices in Market Equilibrium,” *Econometrica*, 63, 841–890.
- (1999): “Voluntary Export Restraints on Automobiles: Evaluating a Trade Policy,” *American Economic Review*, 89, 400–430.
- (2004a): “Differentiated Products Demand Systems from a Combination of Micro and Macro Data: The New Car Market,” *Journal of Political Economy*, 112, 68–105.
- BERRY, S., O. B. LINTON, AND A. PAKES (2004b): “Limit Theorems for Estimating the Parameters of Differentiated Product Demand Systems,” *The Review of Economic Studies*, 71, 613–654.
- BERRY, S. AND A. PAKES (2007): “The Pure Characteristics Demand Model,” *International Economic Review*, 48, 1193–1225.
- BERRY, S. T. AND P. A. HAILE (2014): “Identification in Differentiated Products Markets Using Market Level Data,” *Econometrica*, 82, 1749–1797.
- (2022): “Nonparametric Identification of Differentiated Products Demand Using Micro Data,” Tech. Rep. arXiv:2204.06637, arXiv.

- BLEI, D. M., A. Y. NG, AND M. I. JORDAN (2003): “Latent dirichlet allocation,” *The Journal of Machine Learning Research*, 3, 993–1022.
- CARDELL, N. S. (1997): “Variance Component Structures for the Extreme Value and Logistic Distributions,” *Econometric Theory*, 13, 185–213.
- CHANDRA, A., A. FINKELSTEIN, A. SACARNY, AND C. SYVERSON (2016): “Health Care Exceptionalism? Performance and Allocation in the US Health Care Sector,” *American Economic Review*, 106, 2110–2144.
- CHEN, Y., Y. CHI, J. FAN, AND C. MA (2021): “Spectral Methods for Data Science: A Statistical Perspective,” *Foundations and Trends® in Machine Learning*, 14, 566–806, citation:SpectralMethodsBook2021.
- CHETTY, R., J. N. FRIEDMAN, AND J. E. ROCKOFF (2014): “Measuring the Impacts of Teachers I: Evaluating Bias in Teacher Value-Added Estimates,” *American Economic Review*, 104, 2593–2632.
- CONLON, C. AND J. GORTMAKER (2020): “Best practices for differentiated products demand estimation with PyBLP,” *The RAND Journal of Economics*, 51, 1108–1161.
- (2023): “Incorporating Micro Data into Differentiated Products Demand Estimation with PyBLP,” .
- CONLON, C. AND J. H. MORTIMER (2021a): “Empirical properties of diversion ratios,” *The RAND Journal of Economics*, 52, 693–726.
- CONLON, C., J. H. MORTIMER, P. SARKIS, AND R. GONZALEZ VALDENEGRO (2023): “Effects of Product Availability: Experimental Evidence,” .
- CONLON, C. T. AND J. H. MORTIMER (2013): “Demand Estimation under Incomplete Product Availability,” *American Economic Journal: Microeconomics*, 5, 1–30.
- (2021b): “Efficiency and Foreclosure Effects of Vertical Rebates: Empirical Evidence,” *Journal of Political Economy*, 129, 3357–3404, publisher: The University of Chicago Press.
- DARDANONI, V., P. MANZINI, M. MARIOTTI, AND C. J. TYSON (2020): “Inferring Cognitive Heterogeneity From Aggregate Choices,” *Econometrica*, 88, 1269–1296.
- DEATON, A. AND J. MUELLBAUER (1980): “An Almost Ideal Demand System,” *The American Economic Review*, 70, 312–326, publisher: American Economic Association.
- EFRON, B. AND C. MORRIS (1973): “Stein’s Estimation Rule and its Competitors—An Empirical Bayes Approach,” *Journal of the American Statistical Association*, 68, 117–130.
- FARRELL, J. AND C. SHAPIRO (2010): “Antitrust Evaluation of Horizontal Mergers: An Economic Alternative to Market Definition,” *The B.E. Journal of Theoretical Economics*, 10, publisher: De Gruyter.
- FOX, J. T., K. I. I. KIM, S. P. RYAN, AND P. BAJARI (2011): “A simple estimator for the distribution of random coefficients,” *Quantitative Economics*, 2.

- GOEREE, M. S. (2008): “Limited Information and Advertising in the U.S. Personal Computer Industry,” *Econometrica*, 76, 1017–1074.
- GOOLSBEE, A. AND A. PETRIN (2004): “The Consumer Gains from Direct Broadcast Satellites and the Competition with Cable TV,” *Econometrica*, 72, 351–381.
- GREENE, W. H. AND D. A. HENSHER (2003): “A latent class model for discrete choice analysis: contrasts with mixed logit,” *Transportation Research Part B: Methodological*, 37, 681–698.
- GRIECO, P. L. E., C. MURRY, J. PINKSE, AND S. SAGL (2023): “Conformant and Efficient Estimation of Discrete Choice Demand Models,” .
- GRIECO, P. L. E., C. MURRY, AND A. YURUKOGLU (2021): “The Evolution of Market Power in the US Auto Industry,” .
- HEISS, F., S. HETZENECKER, AND M. OSTERHAUS (2022): “Nonparametric estimation of the random coefficients model: An elastic net approach,” *Journal of Econometrics*, 229, 299–321.
- KANE, T. J. AND D. O. STAIGER (2008): “Estimating Teacher Impacts on Student Achievement: An Experimental Evaluation,” Working Paper 14607, National Bureau of Economic Research, series: Working Paper Series.
- KE, S., J. L. MONTIEL OLEA, AND J. NESBIT (2022): “Robust Machine Learning Algorithms for Text Analysis,” .
- KINGMA, D. P. AND J. BA (2017): “Adam: A Method for Stochastic Optimization,” .
- MAGNOLFI, L., J. MCCLURE, AND A. T. SORENSEN (2022): “Triplet Embeddings for Demand Estimation,” .
- MANZINI, P. AND M. MARIOTTI (2014): “Stochastic Choice and Consideration Sets,” *Econometrica*, 82, 1153–1176.
- MASATLIOGLU, Y., D. NAKAJIMA, AND E. Y. OZBAY (2012): “Revealed Attention,” *American Economic Review*, 102, 2183–2205.
- MATĚJKA, F. AND A. MCKAY (2015): “Rational Inattention to Discrete Choices: A New Foundation for the Multinomial Logit Model,” *American Economic Review*, 105, 272–298.
- MCFADDEN, D. (1974): “Conditional Logit Analysis of Qualitative Choice Behavior,” in *Frontiers in Econometrics*, ed. by P. Zarembka, New York: Academic Press, 105 – 142.
- (1978): “Modelling the Choice of Residential Location,” in *Spatial Interaction Theory and Planning Models*, ed. by A. Karlqvist, L. Lundqvist, F. Snickars, and J. Weibull, North-Holland.
- MCFADDEN, D. AND K. TRAIN (2000): “Mixed MNL models for discrete response,” *Journal of Applied Econometrics*, 15, 447–470.
- MORRIS, C. N. (1983): “Parametric Empirical Bayes Inference: Theory and Applications,” *Journal of the American Statistical Association*, 78, 47–55, publisher: [American Statistical Association, Taylor & Francis, Ltd.].

- NEVO, A. (2000): “Mergers with Differentiated Products: The Case of the Ready-to-Eat Cereal Industry,” *The RAND Journal of Economics*, 31, 395–421.
- (2001): “Measuring Market Power in the Ready-to-Eat Cereal Industry,” *Econometrica*, 69, 307–342.
- PETRIN, A. (2002): “Quantifying the Benefits of New Products: The Case of the Minivan,” *Journal of Political Economy*, 110, 705–729, publisher: The University of Chicago Press.
- QIU, J., M. SAWADA, AND G. SHEU (2021): “Win/Loss Data and Consumer Switching Costs: Measuring Diversion Ratios and the Impact of Mergers,” SSRN Scholarly Paper 3957662, Social Science Research Network, Rochester, NY.
- RAVAL, D., T. ROSENBAUM, AND S. A. TENN (2017): “A Semiparametric Discrete Choice Model: An Application to Hospital Mergers,” *Economic Inquiry*, 55, 1919–1944.
- RAVAL, D., T. ROSENBAUM, AND N. E. WILSON (2022): “Using Disaster Induced Closures to Evaluate Discrete Choice Models of Hospital Demand,” 88.
- REYNOLDS, G. AND C. WALTERS (2008): “The Use of Customer Surveys for Market Definition and the Competitive Assessment of Horizontal Mergers,” *Journal of Competition Law and Economics*, 4, 411–431, publisher: Oxford University Press.
- RIVERS, D. AND Q. VUONG (2002): “Model selection tests for nonlinear dynamic models,” *The Econometrics Journal*, 5, 1–39.
- TRAIN, K. E. (2009): *Discrete Choice Methods with Simulation*, Cambridge: Cambridge University Press, 2 ed.
- UDELL, M. AND A. TOWNSEND (2019): “Why Are Big Data Matrices Approximately Low Rank?” *SIAM Journal on Mathematics of Data Science*, 1, 144–160.
- WU, L., Y. YANG, AND H. LIU (2014): “Nonnegative-lasso and application in index tracking,” *Computational Statistics & Data Analysis*, 70, 116–126.

# Appendices

## A. Theoretical Results

### A.1. Multiple Product Removals

From ?, we know that, in a random coefficients logit demand framework, diversion from good  $j$  to good  $k$  when  $j$  is no longer available is:

$$D_{jk} = \frac{s_k(\mathcal{J}, x) - s_k(\mathcal{J} \setminus j, x)}{s_j(\mathcal{J}, x)} = -\frac{1}{s_j(\mathcal{J}, x)} \int \frac{s_{ij}(\mathcal{J}, x)s_{ik}(\mathcal{J}, x)}{(1 - s_{ij}(\mathcal{J}, x))}$$

For clarity, we use the subscript  $\setminus j$  as an equivalent for writing  $(\mathcal{J} \setminus j)$  and we remove the  $x$ . We can rewrite the previous expression to get the share of  $k$  when  $j$  is not available:

$$s_{k \setminus j} = s_k + \int \frac{s_{ij}s_{ik}}{(1 - s_{ij})} \quad (\text{A1})$$

In addition, we know that for each individual type within the mixed logit model, diversion is written as:

$$D_{ijk} = \frac{s_{ik} - s_{ik \setminus j}}{s_{ij}} = -\frac{s_{ik}}{(1 - s_{ij})}$$

which we can rewrite to find individual  $i$ 's share of  $k$  when  $j$  is not available:

$$s_{ik \setminus j} = s_{ik} + \frac{s_{ik}s_{ij}}{(1 - s_{ij})} = s_{ik} \cdot \left(1 + \frac{s_{ij}}{(1 - s_{ij})}\right) = \frac{s_{ik}}{(1 - s_{ij})} \quad (\text{A2})$$

For the multiple product removals case, we want to prove that the individual share of good  $k$  when all products  $j \in \mathcal{J}_{\text{rem}}$  are removed is:

$$s_{ik \setminus \mathcal{J}_{\text{rem}}} = \frac{s_{ik}}{1 - \sum_{j \in \mathcal{J}_{\text{rem}}} s_{ij}}$$

We can proceed by induction. Equation (2) is the our first case, when removing only one product. Assume that after removing  $p$  products we have

$$s_{ik \setminus 1, \dots, p} = \frac{s_{ik}}{1 - \sum_{j=1}^p s_{ij}},$$

then when we remove  $p+1$  products, we have

$$\begin{aligned} s_{ik \setminus 1, \dots, p+1} &= \frac{s_{ik \setminus 1, \dots, p}}{1 - s_{ip+1 \setminus 1, \dots, p}} = \frac{s_{ik}}{1 - \sum_{j=1}^p s_{ij}} \cdot \frac{1}{1 - \frac{s_{ip+1}}{1 - \sum_{j=1}^p s_{ij}}} = \frac{s_{ik}}{1 - \sum_{j=1}^p s_{ij}} \cdot \frac{1 - \sum_{j=1}^p s_{ij}}{1 - \sum_{j=1}^{p+1} s_{ij}} \\ &= \frac{s_{ik}}{1 - \sum_{j=1}^{p+1} s_{ij}} \end{aligned}$$

Therefore, for a particular instance of IIA logit  $i$ , we have:

$$s_{ik \setminus \mathcal{J}_{\text{rem}}} = \frac{s_{ik}}{1 - \sum_{j \in \mathcal{J}_{\text{rem}}} s_{ij}}$$

For the mixed logit, we want to prove that:

$$s_{k \setminus \mathcal{J}_{\text{rem}}} = s_k + \int s_{ik} \cdot \frac{\sum_{j \in \mathcal{J}_{\text{rem}}} s_{ij}}{1 - \sum_{j \in \mathcal{J}_{\text{rem}}} s_{ij}}$$

We already know that this is the case for a single removed product, as expressed in (1). And, if after removing  $p$  goods we have:

$$s_{0 \setminus 1, \dots, p} = s_k + \int s_{ik} \cdot \frac{\sum_{j=1}^p s_{ij}}{1 - \sum_{j=1}^p s_{ij}}$$

then for  $p+1$ , we have

$$\begin{aligned} s_{k \setminus 1, \dots, p+1} &= s_{k \setminus 1, \dots, p} + \int \frac{s_{ik \setminus 1, \dots, p} \cdot s_{ip+1 \setminus 1, \dots, p}}{1 - s_{ip+1 \setminus 1, \dots, p}} \\ &= s_k + \int s_{ik} \cdot \frac{\sum_{j=1}^p s_{ij}}{1 - \sum_{j=1}^p s_{ij}} + \int \frac{\frac{s_{ik}}{1 - \sum_{j=1}^p s_{ij}} \cdot \frac{s_{ip+1}}{1 - \sum_{j=1}^p s_{ij}}}{1 - \frac{s_{ip+1}}{1 - \sum_{j=1}^p s_{ij}}} \\ &= s_k + \int s_{ik} \cdot \frac{\sum_{j=1}^p s_{ij}}{1 - \sum_{j=1}^p s_{ij}} + \int s_{ik} \cdot \frac{s_{ip+1}}{(1 - \sum_{j=1}^p s_{ij})(1 - \sum_{k=1}^{p+1} s_{ik})} \\ &= s_k + \int s_{ik} \cdot \left[ \frac{\sum_{j=1}^p s_{ij}}{1 - \sum_{j=1}^p s_{ij}} + \frac{s_{ip+1}}{(1 - \sum_{j=1}^p s_{ij})(1 - \sum_{k=1}^{p+1} s_{ik})} \right] \\ &= s_k + \int s_{ik} \cdot \left[ \frac{(1 - \sum_{k=1}^{p+1} s_{ik}) \cdot (\sum_{j=1}^p s_{ij}) + s_{ip+1}}{(1 - \sum_{j=1}^p s_{ij})(1 - \sum_{k=1}^{p+1} s_{ik})} \right] \\ &= s_k + \int s_{ik} \cdot \left[ \frac{(1 - \sum_{k=1}^p s_{ik}) \cdot (\sum_{j=1}^p s_{ij}) + (1 - \sum_{k=1}^p s_{ik}) \cdot s_{ip+1}}{(1 - \sum_{j=1}^p s_{ij})(1 - \sum_{k=1}^{p+1} s_{ik})} \right] \\ &= s_k + \int s_{ik} \cdot \left[ \frac{\sum_{k=1}^{p+1} s_{ik}}{(1 - \sum_{k=1}^{p+1} s_{ik})} \right] \end{aligned}$$

Therefore,

$$s_{k \setminus \mathcal{J}_{\text{rem}}} = s_k + \int s_{ik} \cdot \frac{\sum_{j \in \mathcal{J}_{\text{rem}}} s_{ij}}{1 - \sum_{j \in \mathcal{J}_{\text{rem}}} s_{ij}}$$

and we can write diversion from multiple products removed to  $k$  as:

$$D_{\mathcal{J}_{\text{rem}}, k} = \frac{s_{k \setminus \mathcal{J}_{\text{rem}}} - s_k}{\sum_{j \in \mathcal{J}_{\text{rem}}} s_j} = \frac{1}{\sum_{j \in \mathcal{J}_{\text{rem}}} s_j} \int \frac{s_{ik} \sum_{j \in \mathcal{J}_{\text{rem}}} s_{ij}}{1 - \sum_{j \in \mathcal{J}_{\text{rem}}} s_{ij}}$$

## A.2. Nested Logit Details

We use the Cardell (1997); Berry (1994) parameterization of the nested logit model. (This is not the same as the Train (2009); McFadden (1978) version).

A consumer  $i$  purchasing product  $j$  in a market where it is available obtains utility given by:

$$u_{ij} = \delta_j + \zeta_{ig}(\rho) + (1 - \rho)\varepsilon_{ij}$$

where  $\delta_j$  and  $\rho$  are parameters,  $\varepsilon_{ij}$  is i.i.d. Type-I Extreme-Value and  $\zeta_{ig}$  is the idiosyncratic nest preference, such that  $\zeta_{ig}(\rho) + (1 - \rho)\varepsilon_{ij}$  is also extreme value. Thus, if we denote  $\mathcal{J}_g$  as the set of products in nest  $g$ , the logit inclusive value  $IV_g$  for nest  $g$  is given by:

$$IV_g = \sum_{k \in \mathcal{J}_g} \exp\left(\frac{\delta_k}{1 - \rho}\right)$$

The choice probabilities can be written as the product of the logit for choice  $j$  conditional on category choice  $g$  and the logit probability of choosing category  $g$ :

$$s_{j|g} = \frac{\exp\left(\frac{\delta_j}{1 - \rho}\right)}{IV_g}, \quad s_g = \frac{IV_g^{(1-\rho)}}{\sum_{g'} IV_{g'}^{(1-\rho)}} \quad s_j = s_{j|g} \cdot s_g.$$

Following the results in the appendix of ?, and defining  $Z(\rho, s_g) = [\rho + (1 - \rho)s_g] \in (0, 1]$ , we get two formulas for diversion from product  $j$  to product  $k$ :

$$\begin{aligned} (\text{Same Nest } g:) D_{jk} &= \frac{s_{k|g}}{Z^{-1}(\rho, s_g) - s_{j|g}} \\ (\text{Different Nests:}) D_{jk} &= \frac{s_k \cdot (1 - \rho)}{1 - Z(\rho, s_g) \cdot s_{j|g}}. \end{aligned}$$

It is helpful to define the  $J \times 1$  vector  $\mathbf{s}_{|\mathbf{g}}$  as having entries  $s_{j|g}$  if  $j$  is a member of nest  $g$  and 0 if it is not. This allows us to write the transposed diversion matrix in terms of:

$$\begin{aligned} \mathbf{D}_{cross}^T &= \sum_{g=1}^G (1 - \rho) \mathbf{s} \cdot \left[ \frac{Z^{-1}(\rho, s_g)}{Z^{-1}(\rho, s_g) - \mathbf{s}_{|\mathbf{g}}} \right]^T \\ \mathbf{D}_{same}^T &= \sum_{g=1}^G \mathbf{s}_{|\mathbf{g}} \cdot \left[ \frac{1}{Z^{-1}(\rho, s_g) - \mathbf{s}_{|\mathbf{g}}} \right]^T. \end{aligned}$$

We can combine both so that:

$$\mathbf{D}^T = \sum_{g=1}^G [(1 - \rho) Z^{-1}(\rho, s_g) \mathbf{s} + \mathbf{s}_{|\mathbf{g}}] \cdot \left[ \frac{1}{Z^{-1}(\rho, s_g) - \mathbf{s}_{|\mathbf{g}}} \right]^T$$

Because the term inside the summation can be written as the product of two vectors this diversion matrix will have at most rank  $G$ , the number of nests.

### A.3. Jacobian

Our objective in (6) is  $Q(\theta) = \sum_{(j,k) \in \text{OBS}} (\mathcal{D}_{jk} - D_{jk})^2$ . We assume that  $\bar{q}_0$  is known and set  $\theta = [\mathbf{s}_i, \pi_i, \mathbf{s}]$ . We can write the Jacobian as :

$$\frac{\partial Q(\theta)}{\partial \theta} = \sum_{(j,k) \in \text{OBS}} 2 \cdot (\mathcal{D}_{jk} - D_{jk}) \frac{\partial D_{jk}}{\partial \theta}$$

We can look at this element-by-element for  $j \neq k \neq l$ :

$$\begin{aligned} \frac{\partial D_{jk}}{\partial s_{il}} &= 0, \quad \frac{\partial D_{jk}}{\partial s_k} = 0 \\ \frac{\partial D_{jk}}{\partial s_{ik}} &= \sum_{i=1}^I \frac{\pi_i}{s_j} \cdot \frac{s_{ij}}{1-s_{ij}} = \sum_{i=1}^I \pi_i \cdot D_{jk,i} \cdot \frac{1}{s_{ik}} \\ \frac{\partial D_{jk}}{\partial s_{ij}} &= \sum_{i=1}^I \frac{\pi_i}{s_j} \cdot \frac{s_{ik}}{(1-s_{ij})^2} = \sum_{i=1}^I \pi_i \cdot D_{jk,i} \cdot \frac{1}{s_{ij}(1-s_{ij})} \\ \frac{\partial D_{jk}}{\partial s_j} &= \sum_{i=1}^I -\frac{\pi_i}{s_j^2} \cdot \frac{s_{ik} s_{ij}}{1-s_{ij}} = -\frac{1}{s_j} \sum_{i=1}^I \pi_i \cdot D_{jk,i} \\ \frac{\partial D_{jk}}{\partial \pi_i} &= \frac{s_{ij}}{s_j} \cdot \frac{s_{ik}}{1-s_{ij}} = D_{jk,i} \end{aligned}$$

Each of these derivatives are bounded because  $(\pi_i, s_{ij}, s_j) \in [0, 1]$ . Since we expect  $s_j > 0$  from the data, the only risk is that  $s_{ij} = 1$  (sparsity  $s_{ij} = 0$  does not cause any issues). In this case it is sufficient that  $D_{jk} > 0$  for all  $j$  and at least one  $k$  so that  $\mathbf{s}_i$  has at least two nonzero elements.

## B. Additional Monte Carlo Results

### C. Description of Vending Data

#### C.1. Description of Data and Industry

Globally, the snack foods industry is a \$300 billion market annually, composed of a number of large, well-known firms and some of the most heavily-advertised global brands. Mars Incorporated reported over \$50 billion in revenue in 2010, and represents the third-largest privately-held firm in the US. Other substantial players include Hershey, Nestle, Kraft, Kellogg, and the Frito-Lay division of PepsiCo. While the snack-food industry as a whole might not appear highly concentrated, sales within product categories can be very concentrated. For example, Frito-Lay comprises around 40% of all savory snack sales in the United States, and reported over \$13 billion in US revenues last year, but its sales outside the salty-snack category are minimal, coming mostly through parent PepsiCo's Quaker Oats brand and the sales of *Quaker Chewy Granola Bars*.<sup>23</sup> We report HHI's at both the category level and for all vending products in table C1 from the midwest region of the U.S. If the relevant market is defined at the category level, all categories are considered highly concentrated, with HHIs in the range of roughly 4500-6300. If the relevant market is defined as all products sold

---

<sup>23</sup>Most analysts believe Pepsi's acquisition of Quaker Oats in 2001 was unrelated to its namesake business but rather for Quaker Oats' ownership of Gatorade, a close competitor in the soft drink business.

Manufacturer:	Category:			
	Salty Snack	Cookie	Confection	Total
PepsiCo	78.82	9.00	0.00	37.81
Mars	0.00	0.00	58.79	25.07
Hershey	0.00	0.00	30.40	12.96
Nestle	0.00	0.00	10.81	4.61
Kellogg's	7.75	76.94	0.00	11.78
Nabisco	0.00	14.06	0.00	1.49
General Mills	5.29	0.00	0.00	2.47
Snyder's	1.47	0.00	0.00	0.69
ConAgra	1.42	0.00	0.00	0.67
TGIFriday	5.25	0.00	0.00	2.46
Total	100.00	100.00	100.00	100.00
HHI	6332.02	6198.67	4497.54	2401.41

Table C1: Manufacturer Market Shares and HHI's by Category and Total

Source: IRM Brandshare FY 2006 and Frito-Lay Direct Sales For Vending Machines Data, Heartland Region, 50 best-selling products. ([http://www.vending.com/Vending\\_Affiliates/Pepsico/Heartland\\_Sales\\_Data](http://www.vending.com/Vending_Affiliates/Pepsico/Heartland_Sales_Data))

in a snack-food vending machine, the HHI is below the critical threshold of 2500. Any evaluation of a merger in this industry would hinge on the closeness of competition.

Over the last 25 years, the industry has been characterized by a large amount of merger and acquisition activity, both on the level of individual brands and entire firms. For example, the *Famous Amos* cookie brand was owned by at least seven firms between 1985 and 2001, including the Keebler Cookie Company (acquired by Kellogg in 2001), and the Presidential Baking Company (acquired by Keebler in 1998). *Zoo Animal Crackers* have a similarly complicated history, having been owned by Austin Quality Foods before they too were acquired by the Keebler Cookie Co. (which in turn was acquired by Kellogg).<sup>24</sup>

Our study measures diversion through the lens of a single medium-sized retail vending operator in the Chicago metropolitan area, Mark Vend Company. Each of Mark Vend's machines internally records price and quantity information. The data track total vends and revenues since the last service visit on an item-level basis, but do not include time-stamps for each sale. Any given machine can carry roughly 35 products at one time, depending on configuration.

We observe retail and wholesale prices for each product at each service visit during a 38-month panel that runs from January 2006 to February 2009. There is relatively little price variation within a site, and almost no price variation within a category (e.g., chocolate candy) at a site. This is helpful from an experimental design perspective, but can pose a challenge to structural demand estimation. Very few “natural” stock-outs occur at our set of machines.<sup>25</sup> Most changes to the set of products available to consumers are a result of product rotations, new product introductions, and

<sup>24</sup>Snack foods have an important historic role in market definition. A landmark case was brought by *Tastykake* in 1987 in an attempt to block the acquisition of *Drake* (the maker of Ring-Dings) by *Ralston-Purina's Hostess* brand (the maker of Twinkies). That case established the importance of geographically significant markets, as Drake's had only a 2% market share nationwide, but a much larger share in the Northeast (including 50% of the New York market). Tastykake successfully argued that the relevant market was single-serving snack cakes rather than a broad category of snack foods involving cookies and candy bars. [Tasty Baking Co. v. Ralston Purina, Inc., 653 F. Supp. 1250 - Dist. Court, ED Pennsylvania 1987.]

<sup>25</sup>Mark Vend commits to a low level of stock-out events in its service contracts.

product retirements. Over all sites and months, we observe 185 unique products. Some products have very low levels of sales and we consolidate them with similar products within a category produced by the same manufacturer, until we are left with 73 ‘products’ that form the basis of the rest of our exercise.<sup>26</sup>

## C.2. Experimental Design

We implemented ten exogenous product removals with the help of Mark Vend Company. These product removals have been used in two other projects including . Our experiment uses 66 snack machines located in professional office buildings and serviced by Mark Vend. Most of the customers at these sites are employees of law firms and insurance companies. Our goal in selecting the machines was to choose machines that could be analyzed together, in order to be able to run each product removal over a shorter period of time across more machines.<sup>27</sup> These machines were also located on routes that were staffed by experienced drivers, which maximized the chance that the product removal would be successfully implemented. The 66 machines used for each treatment are distributed across five of Mark Vend’s clients, which had between 3 and 21 machines each.<sup>28</sup>

For each treatment, we remove a product from all machines at a client site for a period of 2.5 to 3 weeks. The four products that we remove are the two best-selling products from either (a) confections seller Mars Incorporated (Snickers and Peanut M&Ms) or (b) cookie seller Kellogg’s (Famous Amos Chocolate Chip Cookies and Zoo Animal Crackers). We refer to exogenously-removed products as the *focal products* throughout our analysis.<sup>29</sup>. Whenever a product was exogenously removed, poster-card announcements were placed at the front of the empty product column.<sup>30</sup> The dates of the interventions range from June 2007 to September 2008, with all removals run during the months of May - October. We collected data for all machines for just over three years, from January of 2006 until February of 2009. Although data are recorded at the level of a service visit, it is more convenient to organize observations by week, because different visits occur on different days of the week.<sup>31</sup> The cost of implementing the experiment consisted primarily of drivers’ time.<sup>32</sup>

---

<sup>26</sup>For example, we combine Milky Way Midnight with Milky Way, Ruffles Original with Ruffles Sour Cream & Cheddar, and various flavors of Pop-Tarts together.

<sup>27</sup>Many high-volume machines are located in public areas (e.g., museums or hospitals), and feature demand patterns (and populations) that vary enormously from one day to the next, so we did not use machines of this nature. In contrast, the work-force populations at our experimental sites have relatively stable demand patterns.

<sup>28</sup>The largest client had two sets of floors serviced on different days, and we divided this client into two sites. Generally, each site is spread across multiple floors in a single high-rise office building, with machines located on each floor.

<sup>29</sup>Not reported here are two treatment arms removing best-selling products from Pepsi’s Frito Lay Division, which we omit for space considerations, and because Pepsi’s products already dominate the salty snack category (which makes merger analysis less relevant). We also ran two additional treatments in which we removed two products at once; again we omit those for space considerations and because they don’t speak to our diversion ratio example. These are analyzed in Conlon and Mortimer (2021b)

<sup>30</sup>The announcements read: “This product is temporarily unavailable. We apologize for any inconvenience.” The purpose of the card was two-fold: first, we wanted to avoid dynamic effects on sales as much as possible, and second, Mark Vend wanted to minimize the number of phone calls received in response to the stock-out events.

<sup>31</sup>During each 2-3 week experimental period, most machines receive service visits about three times. However, the length of service visits varies across machines, with some machines visited more frequently than others. In order to define weekly observations, we assume that sales are distributed uniformly among the business days in a service interval, and assign sales to weeks. We allow our definition of when weeks start and end to depend on the client site and experiment, because different experimental treatments start on different days of the week. At some site-experiment pairs, weeks run Tuesday to Monday, while others run Thursday to Wednesday.

<sup>32</sup>Drivers had to spend extra time removing and reintroducing products to machines, and the driver dispatcher had

Our experiment differs somewhat from an ideal experiment. Ideally, we would be able to randomize the choice set for each individual consumer. Technologically, of course, that is difficult in both vending and traditional brick and mortar contexts.<sup>33</sup> Additionally, because we remove all of the products at an entire client site for a period of 2.5 to 3 weeks, we lack a contemporaneous “same-side” group of untreated machines. We chose this design, rather than randomly staggering the product removals, because we (and the participating clients) were afraid consumers might travel from floor to floor searching for stocked-out products. This design consideration prevents us from using contemporaneous control machines in the same building, and makes it more difficult to capture weekly variation in sales due to unrelated factors, such as a client location hitting a busy period that temporarily induces long work hours and higher vending sales. Conversely, the design has the benefit that we can aggregate over all machines at a client site, and treat the entire site as if it were a single machine. Despite the imperfections of field experiments in general, these are often the kinds of tests run by firms in their regular course of business, and may most closely approximate the type of experimental information that a firm may already have available at the time when a proposed merger is initially screened.

### C.3. Description of Experimental Data

We summarize the data generated by our product removals in Table 4. Across our four treatments and 66 machines, we observe between 161-223 treated machine-weeks. In the untreated group, we observe 8,525 machine-weeks and more than 700,000 units sold. Each treatment week exposes around 2,700-3,500 individuals to the product removal, of which around 134-274 would have purchased the focal product in an average week. Each treatment lasts 2.5-3 weeks, and between approximately 14,000-19,000 sales are recorded during the treated periods. The treated group consists of the 400-1,200 individuals who would have purchased the focal product had it been available for each treatment. This highlights one of the main challenges of measuring diversion experimentally: for the purposes of measuring the treatment effect, only individuals who would have purchased the focal product, had it been available, are considered “treated,” yet we must expose many more individuals to the product removal, knowing that many of them were not interested in the focal product in the first place.

In general, we see that the overall sales per machine-week are higher during the treatment period (between 83.3-89.4) than during the control period (82.2).<sup>34</sup> This illustrates a second challenge, which is that there is a large amount of variation in overall sales at the weekly level, independent of our product removals. This weekly variation in overall sales is common in many retail environments. We often observe week-over-week sales that vary by more than 20%.

In our particular setting, many of the product removals were implemented during the summer

---

to spend time instructing the drivers, tracking the dates of each product removal, and reviewing the data as they were collected. Drivers are generally paid a small commission on the sales on their routes, so if sales levels fell dramatically as a result of the product removals, their commissions could be affected. Tracking commissions and extra minutes on each route for each driver would have been prohibitively expensive to do, and so drivers were provided with \$25 gift cards for gasoline during each week in which a product was removed on their route to compensate them for the extra time and the potential for lower commissions.

<sup>33</sup>This leaves our design susceptible to contamination if for example, Kraft runs a large advertising campaign for Planters Peanuts that corresponds to the timing of one of our product removals.

<sup>34</sup>The per-week sales can be misleading because not all machines are measured in every week during the treatment period. This is because the product removals have slightly different start dates at different client site locations. This leads to a somewhat liberal definition of “treatment week” as only one or two machines might be treated in the final week.

of 2007, which was a high point in demand at several sites, most likely due to macroeconomic conditions.

Our estimates of diversion face two major challenges: (1) the set of products may vary for non-experimental reasons across machines and time; (2) demand is volatile both at the product level, and at the aggregate level.

Using four simple assumptions motivated by economic theory, we develop an estimator for the average treatment effect version of the diversion ratio which deals with these challenges. The first two assumptions restrict the set of machine-weeks that can act as a control for a particular machine-treated week as in a *matching estimator*. The second two assumptions affect the way in which estimates of  $\widehat{\Delta q}_j, \widehat{\Delta q}_k$  are used to construct  $\overline{D}_{jk}$  and employ the principle of *Bayesian Shrinkage*. All four assumptions are implications of the economic restriction that consumers make discrete choices among substitutes.<sup>35</sup> Using the following four assumptions we demonstrate how we estimate  $\overline{D}_{jk}$  from our experimental data.

**Assumption 1. Valid Controls** For a machine-week observation to be included as a control for  $q_{k,t}$  it must: (a) have product  $k$  available; (b) be from the same vending machine; (c) not be included in any of our treatments.

**Assumption 2. Substitutes:** Removing product  $j$  can never increase the overall level of sales during a period, and cannot decrease sales by more than the sales of  $j$ .

**Assumption 3. Unit Interval:**  $D_{jk} \in [0, 1]$ .

**Assumption 4. Unit Simplex:**  $D_{jk} \in [0, 1]$  and  $\sum_{\forall k} D_{jk} = 1$ .

#### C.4. Matching Assumptions

Consider an estimate of  $\widehat{\Delta q}_k$ , where  $W = 1$  denotes the removal of product  $j$ :<sup>36</sup>

$$\widehat{\Delta q}_k = E[q_k|W = 1] - E[q_k|W = 0]$$

We adjust our calculation of the expectation to address the volatility of demand. To be explicit, one can introduce a covariate  $\xi$  (demand shock):

$$E[q_k|W = w] = \int q_k(\xi, w) f(\xi|W = w) d\xi$$

The treated and control periods have different distributions of covariates (demand shocks) because  $f(\xi|W = 1) \neq f(\xi|W = 0)$ . The typical solution involves *matching* or *balancing*, where one re-weights observations in the control period using measure  $g(\cdot)$  so that  $f(\xi|W = 1) = g(\xi|W = 0)$  and then calculates the expectation  $E_g[q_k|Z = 0]$  with respect to measure  $g$ .<sup>37</sup> For each treated week  $t$ , one can construct a set of matched control weeks within a neighborhood  $S(\xi_t)$ , where  $S(\xi_t)$

<sup>35</sup>While we estimate the ATE version of  $D_{jk}$  in our example, the procedure described in this section could be used to estimate a LATE if the treatment were, for example, a 10% price increase instead of a product removal.

<sup>36</sup>One advantage of using a product removal experiment is that  $E[q_j|W = 1] = 0$  by construction (consumers cannot purchase products that are unavailable). This also helps rule out one set of potential defiers. The second set of defiers, those that purchase  $k$  only when  $j$  is available are ruled out if  $(j, k)$  are substitutes rather than complements.

<sup>37</sup>We omit the usual discussion of the conditional independence assumption because we have randomized assignment of  $W$ .

is the set of control weeks that correspond to treated week  $t$ , and  $\xi_t$  is an unobserved demand shock. Having chosen  $S(\xi_t)$ , the change in sales for the chosen control weeks is given as:

$$\Delta q_{k,t}(\xi_t) = q_{k,t}(\xi_t) - \frac{1}{|\#s \in S(\xi_t)|} \sum_{s \in S(\xi_t)} q_{k,s} \quad \text{with} \quad \widehat{\Delta q_k} = \sum_t \Delta q_{k,t}(\xi_t) \quad (\text{C3})$$

Our first two assumptions tell us how to choose  $S(\xi_t)$ . Assumption 1 is straightforward: it controls for unobserved machine-level heterogeneity by restricting potential controls to different (untreated) weeks at the same machine. If  $\xi_t$  were observed, one could employ conventional matching estimators (such as  $k$ -nearest neighbor or local-linear regression Abadie and Imbens (2011)). However,  $\xi_t$  is unobserved, so we rely on Assumption 2 (removing a product cannot increase total sales, and cannot reduce sales by more than the sales of the product removed) instead.

We implement Assumption 2 as follows. We let  $Q_t$  denote the sales of all products during the treated machine-week, and  $Q_s$  denote the overall sales of a potential control machine-week. Given a treated machine-week  $t$ , we look for the corresponding set of control periods that satisfy Assumption 1 and further restrict them to satisfy Assumption 2:

$$\{s : Q_s - Q_t \in [0, q_{js}]\} \quad (\text{C4})$$

The problem with a direct implementation of ((C4)) is that periods with (unexpectedly) higher sales of the focal product  $q_{js}$  are more likely to be included as a control, which would lead us to underestimate the diversion ratio. We propose a modification of ((C4)) that is unbiased. We replace  $q_{js}$  with  $\widehat{q}_{js} = E[q_{js}|Q_s, W = 0]$ . An easy way to obtain the expectation is to run an OLS regression of  $q_{js}$  on  $Q_s$  using data only from untreated machine-weeks that satisfy Assumption 1:

$$S_t \equiv \{s : Q_s^0 - Q_t^1 \in [0, \widehat{b}_0 + \widehat{b}_1 Q_s^0]\} \quad (\text{C5})$$

Thus ((C5)) defines the set of control periods  $S_t$  that correspond to treatment period  $t$  under Assumption 1 and Assumption 2.<sup>38</sup> Plugging this into equation ((C3)) gives estimates of  $\widehat{\Delta q_k}$  and  $\widehat{\Delta q_j}$ .<sup>39</sup>

### C.5. Bayesian Shrinkage Assumptions

Our Assumption 3 and Assumption 4 place restrictions on how we calculate the diversion ratio given our estimates of  $\widehat{\Delta q_k}$  and  $\widehat{\Delta q_j}$ . The idea is that there may be better estimates of  $\overline{D}_{jk}$  than the simple ratio  $\frac{\widehat{\Delta q_k}}{\widehat{\Delta q_j}}$ . For example, we might find large but noisy estimates of diversion to a substitute product based on only a few observations and a better estimate might adjust for that uncertainty.<sup>40</sup>

<sup>38</sup>The economic implication of Assumption 2 is that if all treatment and control weeks faced an identical set of substitute products, the sum of the diversion ratios from  $j$  to all other products would lie between zero and one (for each  $t$ ):  $\sum_{k \neq j} D_{jk,t} \in [0, 100\%]$ .

<sup>39</sup>There are stronger assumptions one could make in order to implement a more traditional matching or balancing estimator in the spirit of Abadie and Imbens (2011). Suppose a third product  $k'$  was similarly affected by the demand shock  $x$  but we knew ex-ante that  $D_{jk'} = 0$ . If so, one could match on similar sales levels of  $q_{k'}$ . For our vending example this might involve using sales at a nearby soft drink machine to control for overall demand at the snack machine, or it using sales of chips to control for sales of candy bars. We find that when all four assumptions are used, additional matching criteria have no appreciable effect on our estimates.

<sup>40</sup>A baseball analogy is apt. If one hitter has 150 hits in 500 at bats and another has 2 hits in 5 at bats, one would likely believe the hitter with the .300 average is a better batter, even though he has a lower batting average. This is

We can see how these assumptions work by writing the diversion ratio as the probability of a binomial with  $\Delta q_j$  trials and  $\Delta q_k$  successes:

$$\Delta q_k | \Delta q_j, D_{jk} \sim Bin(n = \Delta q_j, p = D_{jk}) \quad (C6)$$

This is considered a nonparametric estimator as long as we estimate a separate binomial probability  $D_{jk}$  for each  $(j, k)$ .

We implement Assumption 3 by placing a prior on  $D_{jk}$  that restricts all of the mass to the unit interval  $D_{jk} | \mu_{jk}, m_{jk} \sim Beta(\mu_{jk}, m_{jk})$ . Assumption 4 goes further and restricts the vector  $\mathbf{D}_j$  to the unit simplex, which we implement with the prior  $\mathbf{D}_j \sim Dirichlet(\mu_{j0}, \mu_{j1}, \dots, \mu_{jK}, m_{jk})$ . This has the effect of using information about  $D_{jk}$  to inform our estimates for  $D_{jk'}$ .

There are two ways to parametrize the Beta (and Dirichlet) distributions. In the traditional  $Beta(\beta_1, \beta_2)$  formulation  $\beta_1$  denotes the number of prior successes and  $\beta_2$  denotes the number of prior failures (observed before any experimental observations). Under an alternative formulation,  $Beta(\mu, m)$ :  $\mu = \frac{\beta_1}{\beta_1 + \beta_2}$  denotes the prior mean and  $m$  denotes the number of “psuedo-observations”  $m = \beta_1 + \beta_2$ . We work with the latter formulation for both the Beta and Dirichlet distributions.<sup>41</sup> This formula makes it easy to express the posterior mean (under Assumption 3) as a *shrinkage estimator* that combines our prior information with our experimental data:

$$\widehat{D}_{jk} = \lambda \cdot \mu_{jk} + (1 - \lambda) \frac{\Delta q_k}{\Delta q_j}, \quad \lambda = \frac{m_{jk}}{m_{jk} + \Delta q_j} \quad (C7)$$

The weight put on our prior mean is denoted by  $\lambda$ , and directly depends on how many “pseudo-observations” we observe from our prior before observing experimental outcomes. One reason this estimator is referred to as a “shrinkage” estimator, is because as  $\Delta q_j$  becomes smaller (and our experimental outcomes are less informative),  $\widehat{D}_{jk}$  shrinks towards  $\mu_{jk}$  (from either direction). Thus, when our product removals provide lots of information about diversion from  $j$  to  $k$  we rely on the experimental outcomes, but when our experimental variation is less informative, we rely more on our prior information.<sup>42</sup> This has the desirable property of taking extreme but imprecisely estimated parameters and pushing them towards the prior mean.

Our remaining challenge is how to specify the prior  $(\mu_{jk}, m_{jk})$ . Ideally, the location of the prior  $\mu$  would be largely irrelevant while the prior strength  $m$  would be as small as possible.<sup>43</sup> An uniform or uninformative prior might be to let  $\mu_{jk} = \frac{1}{K+1}$  where  $K$  is the number of substitutes. An informative prior centered on the plain IIA logit estimates would let  $\mu_{jk} = \frac{s_k}{1-s_j}$  so that (prior) diversion is proportional to market shares.<sup>44</sup>

---

because there is a much greater chance that the second hitter was merely lucky.

<sup>41</sup>The Dirichlet is a generalization of the Beta to the unit simplex. The mean parameters  $[\mu_0, \mu_1, \dots, \mu_K]$  form a unit simplex while  $m$  denotes the number of pseudo-observations.

<sup>42</sup>We cannot provide a similar closed-form characterization under Assumption 4. Though there is a conjugacy relationship between the Dirichlet and the Multinomial, there is no conjugacy relationship between the Dirichlet and the Binomial except under the special case where the same number of treated individuals  $\Delta q_j$  are observed for each substitute  $k$ .

<sup>43</sup>Indeed, with all four assumptions this is true. We use only a small number of prior observations  $m < 4$ , and the location of the prior is almost completely irrelevant.

<sup>44</sup>When  $\mu_{jk}$  is chosen as a function of the same observed dataset (including from estimated demand parameters) this is a form of an Empirical Bayes estimator. The development of Empirical Bayes shrinkage is attributed to Morris (1983); Efron and Morris (1973) and has been widely used in applied microeconomics to shrink outliers from a distribution of fixed effects in teacher value added Chetty et al. (2014); Kane and Staiger (2008) or hospital quality

We use the IIA logit prior not because it is the best estimate of the diversion ratio absent experimental data, but rather because assuming diversion proportional to market share is commonplace among practitioners in the absence of better data.<sup>45</sup> An advantage of the shrinkage estimator is that it allows us to nest the parametric estimate of diversion currently used in practice and the experimental outcomes, depending on our choice of  $m$ . A smaller  $m$  implies a *weaker prior* and more weight on the observed data.

We report the results of our nonparametric empirical Bayes estimates below in Table 5. These include the aggregate share  $S_j$  as well as the estimated diversion ratios  $D_{jk}$  under the Dirichlet prior with  $m = 4$  pseudo-observations. In the Appendix we document that even under this very weak prior, the confidence intervals for our second-choice diversion estimates are relatively precise.

## D. Other Product Removals

In this section we present our diversion estimates and compare them with parametric methods for all treatment arms.

Product	Shares	Nonparam	Logit	RCC	RCN	CMS(I=2)	CMS(I=3)
Outside Good	30.12	24.42	23.21	23.39	29.16	32.34	30.2
Twix Caramel	2.42	20.25	2.47	2.91	2.31	7.95	12.06
M&M Peanut 1.74 oz	4.14	16.44	3.38	3.74	3.14	7.02	9.79
Rold Gold (Con)	2.56	5.93	3.04	1.28	2.14	2.56	2.74
M&M Milk Chocolate	1.16	5.3	2.02	2.5	1.92	5.33	7.08
Butterfinger	0.5	4.13	1.17	1.29	1.08	3.49	4.05
Planters (Con)	1.92	3.96	2.22	1.84	1.57	3.81	4.92
Choc Chip Famous Amos	2.05	2.87	1.68	1.86	1.3	0.0	0.39
Choc Herhsey (Con)	0.22	2.57	1.32	1.63	1.29	1.72	1.87
Sun Chip LSS	2.12	2.23	1.93	1.96	1.35	0.0	0.25
Choc SandFamous Amos	0.4	1.75	1.0	1.12	0.8	1.18	1.64

Table D2: Top Substitutes: Snickers Removal

Product	Shares	Nonparam	Logit	RCC	RCN	CMS(I=2)	CMS(I=3)
Outside Good	30.12	23.86	22.93	23.37	20.2	29.27	25.44
M&M Peanut 1.74 oz	4.14	9.72	3.34	3.74	2.43	5.94	7.31
Twix Caramel	2.42	7.81	2.44	2.91	1.79	6.72	9.0
Snickers 2.07oz	3.96	6.93	3.19	3.53	2.33	7.35	8.08
Planters (Con)	1.92	6.15	2.19	1.85	1.23	3.99	5.21
Choc Chip Famous Amos	2.05	5.99	1.66	1.87	8.23	0.31	4.74
Rold Gold (Con)	2.56	5.15	3.01	1.29	1.68	2.16	2.08
Choc Herhsey (Con)	0.22	3.68	1.31	1.63	1.0	2.07	2.42
Rice Krispies Treats 1.7oz	0.27	3.63	0.99	1.07	0.69	2.49	1.74
Baked (Con)	2.39	3.04	2.09	1.82	1.17	2.32	1.15
Popcorn (Con)	0.42	2.81	1.0	0.71	0.56	0.58	0.56

Table D3: Top Substitutes: Zoo Animal Crackers Removal

---

Chandra et al. (2016)

<sup>45</sup>If we had estimates from a random coefficients demand model, we could use those estimates of the diversion ratio instead. In practice, we find that under Assumption 4 the choice of  $\mu_{jk}$  becomes irrelevant.

Product	Shares	Nonparam	Logit	RCC	RCN	CMS(I=2)	CMS(I=3)
Outside Good	30.12	27.13	22.98	12.58	17.82	28.39	26.8
Planters (Con)	1.92	9.33	2.2	1.99	3.59	4.9	4.85
Butterfinger	0.5	6.28	1.16	0.88	0.76	2.66	2.44
Cliff (Con)	0.28	5.94	1.16	1.0	0.72	0.99	0.91
Nabisco (Con)	0.39	5.71	1.1	1.3	0.62	1.43	1.43
M&M Milk Chocolate	1.16	5.67	2.0	1.17	1.34	3.77	3.06
Sun Chip LSS	2.12	5.13	1.91	2.43	3.1	5.53	6.02
Choc Chip Famous Amos	2.05	4.33	1.66	1.61	0.9	0.88	5.82
Choc Herhsey (Con)	0.22	3.98	1.31	0.83	0.9	2.95	3.02
Rice Krispies Treats 1.7oz	0.27	3.3	0.99	1.13	0.7	2.22	2.04
Cheez-It Original SS	2.0	3.26	1.89	3.42	3.07	1.46	1.54

Table D4: Top Substitutes: Doritos Removal

Product	Shares	Nonparam	Logit	RCC	RCN	CMS(I=2)	CMS(I=3)
Outside Good	30.12	33.83	23.25	25.33	27.48	31.81	30.87
Snickers 2.07oz	3.96	17.58	3.24	4.93	6.76	8.55	8.47
M&M Milk Chocolate	1.16	10.69	2.02	2.73	4.32	5.24	5.77
Planters (Con)	1.92	8.99	2.22	4.7	1.48	3.75	4.4
Butterfinger	0.5	5.41	1.18	1.81	2.43	3.43	3.56
Raisinets	1.6	5.21	1.62	2.22	3.4	2.17	2.36
Twix Caramel	2.42	4.96	2.48	2.95	5.19	7.82	9.58
Sun Chip LSS	2.12	1.98	1.93	1.66	1.28	0.0	1.1
Cliff (Con)	0.28	1.68	1.17	1.42	1.03	1.38	1.27
Nabisco (Con)	0.39	1.63	1.12	1.0	0.82	1.31	1.31
Choc Herhsey (Con)	0.22	1.45	1.32	1.74	2.91	1.69	2.05

Table D5: Top Substitutes: M&M's Peanut Removal

Product	Shares	Nonparam	Logit	RCC	RCN	CMS(I=2)	CMS(I=3)
Outside Good	30.12	27.96	23.07	19.33	17.61	27.05	32.32
Sun Chip LSS	2.12	13.16	1.92	2.2	3.04	8.99	4.9
Planters (Con)	1.92	8.3	2.21	1.96	3.52	5.71	3.17
M&M Peanut 1.74 oz	4.14	8.27	3.36	2.97	2.17	3.93	2.17
Ruffles (Con)	2.8	6.98	2.57	2.97	4.07	2.49	4.19
Choc Herhsey (Con)	0.22	6.32	1.31	1.32	0.89	3.78	2.85
Dorito Nacho LSS	2.7	5.9	2.19	2.53	3.47	1.03	0.0
Snickers 2.07oz	3.96	4.58	3.21	3.21	2.08	4.86	1.29
Frito LSS	1.86	4.17	2.19	2.5	3.45	1.42	3.88
M&M Milk Chocolate	1.16	2.57	2.01	1.97	1.33	2.98	1.42
Cheetos Flaming Hot LSS	0.62	2.55	0.97	1.12	1.55	0.52	1.56

Table D6: Top Substitutes: Cheetos Removal

	<b>Product</b>	<b>Shares</b>	<b>Nonparam</b>	<b>Logit</b>	<b>RCC</b>	<b>RCN</b>	<b>CMS(I=2)</b>	<b>CMS(I=3)</b>
Dorito Blazin Buffalo Ranch LSS	Sun Chip LSS	2.12	26.81	1.9	2.05	1.04	23.83	25.99
	Outside Good	30.12	21.41	22.86	21.17	19.64	23.73	21.02
	Planters (Con)	1.92	11.57	2.19	1.79	1.21	9.48	12.15
	Choc Herhsey (Con)	0.22	7.87	1.3	1.46	0.96	7.47	7.99
	Rasbry Knotts	0.7	4.07	0.71	0.78	3.42	3.3	3.9
	Ruffles (Con)	2.8	3.85	2.54	2.55	1.4	5.0	3.41
	Choc SandFamous Amos	0.4	3.54	0.98	1.09	4.9	2.2	3.58
	Grandmas Choc Chip	1.09	3.37	1.26	1.4	6.16	2.15	2.98
	Nabisco (Con)	0.39	2.52	1.1	1.19	5.45	2.12	3.08
	Butterfinger	0.5	1.33	1.16	1.2	0.8	0.83	1.56

Table D7: Top Substitutes: Chocolate Chip Famous Amos Removal

	<b>Product</b>	<b>Shares</b>	<b>Nonparam</b>	<b>Logit</b>	<b>RCC</b>	<b>RCN</b>	<b>CMS(I=2)</b>	<b>CMS(I=3)</b>
Farleys Mixed Fruit Snacks	Outside Good	30.12	31.21	23.07	12.41	17.54	27.05	32.32
	Cheez-It Original SS	2.0	7.25	1.9	3.38	3.05	1.69	3.35
	Nature Valley (Con)	2.0	7.11	2.28	2.33	1.44	1.29	3.07
	Rice Krispies Treats 1.7oz	0.27	5.85	1.0	1.11	0.7	1.92	3.02
	Snickers 2.07oz	3.96	4.9	3.21	2.39	2.02	4.86	1.29
	FritoLay (Con)	1.48	4.8	1.58	2.67	2.53	1.72	3.86
	Snyders (Con)	4.15	4.01	4.07	6.01	6.53	0.67	1.85
	Butterfinger	0.5	3.67	1.17	0.85	0.73	2.26	1.93
	Sun Chip LSS	2.12	3.62	1.92	2.39	3.09	8.99	4.9
	Rasbry Knotts	0.7	3.54	0.71	0.53	0.38	1.53	1.09

Table D8: Top Substitutes: Cheetos (bis) Removal

	<b>Product</b>	<b><math>\mathcal{S}</math></b>	<b><math>\mathcal{D}_{jk}</math></b>	<b>Logit</b>	<b>RCC</b>	<b>RCN</b>	<b>CMS(I=2)</b>	<b>CMS(I=3)</b>
Reeses Peanut Butter Cups	Outside Good	30.12	33.21	23.25	24.74	26.30	31.81	30.87
	Snickers	3.96	17.74	3.24	4.69	6.29	8.55	8.47
	Twix Caramel	2.42	7.23	2.48	2.85	4.83	7.82	9.58
	KarNuts (Con)	1.70	5.91	1.74	3.13	1.51	1.80	1.71
	Baked (Con)	2.39	5.04	2.12	1.30	1.37	2.66	2.02
	Nature Valley (Con)	2.00	4.78	2.30	2.55	2.01	1.57	1.12
	Twizzlers	0.33	4.37	1.86	2.07	1.74	1.28	1.15
	M&M Milk Chocolate	0.59	4.03	1.78	2.78	3.48	3.19	4.14
	Rice Krispies Treats	0.27	2.62	1.01	0.90	0.94	2.82	2.47
	Planters (Con)	1.92	2.36	2.22	4.33	1.45	3.75	4.40

Table D9: Top Substitutes: M&M's Peanut (bis) Removal

<b>Product</b>	$\mathcal{S}$	$\mathcal{D}_{jk}$	<b>Logit</b>	RCC	RCN	<b>CMS(I=2)</b>	<b>CMS(I=3)</b>
Outside Good	30.12	35.57	23.59	14.99	18.03	28.39	32.32
Baked (Con)	2.39	8.85	2.15	3.34	3.54	2.01	3.54
FritoLay (Con)	1.48	7.74	1.61	2.47	2.64	1.68	3.86
Snickers	3.96	7.52	3.28	2.72	2.08	6.15	1.29
Ruffles (Con)	2.80	6.77	2.62	3.96	4.31	1.92	4.19
Frito	1.86	6.74	2.23	2.75	3.66	1.38	3.88
Dorito Blazin Buffalo Ranch	0.61	3.07	1.56	2.40	2.59	1.57	1.69
Cheez-It Original	2.00	2.70	1.94	3.09	3.19	1.46	3.35
Rold Gold (Con)	2.56	2.66	3.09	6.82	5.09	1.81	1.52
Rice Krispies Treats	0.27	2.59	1.02	1.16	0.72	2.22	3.02
KarNuts (Con)	1.70	2.49	1.76	1.34	1.10	1.29	1.12

Table D10: Top Substitutes: Cheetos and Doritos (double) Removal

<b>Product</b>	<b>Shares</b>	<b>Nonparam</b>	<b>Logit</b>	RCC	RCN	<b>CMS(I=2)</b>	<b>CMS(I=3)</b>
Outside Good	30.12	36.33	24.02	24.64	27.99	34.78	33.62
Twix Caramel	2.42	11.64	2.56	2.93	5.31	8.56	11.97
M&M Milk Chocolate	1.16	8.33	2.09	2.63	4.43	5.73	7.1
Choc Mars (Con)	1.11	6.79	1.76	2.11	3.68	1.71	2.22
Reeses Peanut Butter Cups	0.59	6.57	1.84	2.78	3.83	3.48	5.12
Butterfinger	0.5	5.22	1.22	1.71	2.49	3.75	4.2
Raisinets	1.6	3.28	1.67	2.12	3.48	2.38	2.92
Nonchoc Other (Con)	0.78	2.63	1.45	1.81	1.33	0.69	0.74
Choc Chip Famous Amos	2.05	2.48	1.74	1.79	1.23	0.0	0.21
Choc Herhsey (Con)	0.22	2.16	1.37	1.69	2.98	1.85	2.15
Planters (Con)	1.92	2.02	2.3	4.16	1.52	4.1	5.13

Table D11: Top Substitutes: Snickers and M&M's Peanut (double) Removal

## E. Additional Vending Results

We present alternate results with  $L_2$  penalty  $\left(\sum_{i=1}^I \pi_i^2\right)$  so that all things equal we prefer models with smaller outside good shares.<sup>46</sup>

---

<sup>46</sup>This prevents trivial fits by driving the outside good share to one, and all inside shares to zero. Though the inclusion of observed elements of  $\mathcal{D}$  should also prevent this in our main specification without the penalty.

	Model/Rank:		I = 1		I = 2		I = 3			I = 4			
	Weight on individual:		100.0%	50.8%	49.2%	33.6%	34.1%	32.3%	25.8%	25.1%	24.7%	24.4%	
	Product	Logit Sj	i = 1	i = 1	i = 2	i = 1	i = 2	i = 3	i = 1	i = 2	i = 3	i = 4	
SALTY SNACKS	Snyders (Con)	2.21	0.7	0.34	1.5	0	2.61	0.33	0	2.44	0.26	1.24	
	Cheetos	2.52	0.5	0	0	0.19	0.76	0.29	0.39	1.03	0.07	0.01	
	Ruffles (Con)	0.98	1.88	0.61	4.91	0	5.62	3.06	0	6.41	2.97	0	
	Dorito Nacho	2.05	0.98	0.93	0.51	0	0	0	0	0	0	0	
	Rold Gold (Con)	0.94	1.86	2.35	0.41	2.48	1.37	0.67	0	2.38	0.25	8.44	
	Baked (Con)	1.99	2.08	2.39	1.08	1.28	5.64	0	2.41	4.62	0	0	
	Salty Other (Con)	2.78	0.22	0.17	0.25	0.09	0.42	0.25	0	0.38	0.16	0.93	
	Sun Chip	1.81	4.76	0	18.6	0	0	21.83	0	5.81	20.63	0	
	Cheez-It	1.77	1.47	0.76	3.09	0	4.7	1.47	0.49	4.65	1.35	0	
	Jays (Con)	1.48	0.17	0.16	0.13	0.08	0.33	0.17	0	0.3	0.1	0.73	
	Frito	2.03	1.41	0.98	2.33	0	6.36	0	0	6.26	0	0.62	
	FritoLay (Con)	1.49	1.71	1.35	2.41	0.03	6.25	0.48	0.11	6.16	0.41	1.59	
	Smartfood	1.51	0.56	0.53	0.49	0.34	0.87	0.61	0.36	0.85	0.49	0.58	
	Lays	1.45	0.54	0.4	0.77	0	2.1	0.05	0	1.98	0.01	0.8	
	Cheetos Flamin	0.96	0.55	0.35	0.91	0	2.43	0.06	0	2.44	0.06	0.11	
	Dorito Blazin	1.45	1.47	0	5.49	0	0	6.8	0	1.82	6.37	0	
	Popcorn (Con)	2.06	0.51	0.58	0.22	0.63	0	0.69	0.41	0.12	0.58	0.95	
	Ritz Bits	0.51	0.16	0.21	0	0.19	0.09	0.11	0	0.17	0.04	0.63	
CHOCOLATE CANDY	M&M Peanut	3.21	4.78	7.08	0	8.74	1.98	0.02	10.19	0	0	2.22	
	Snickers	3.53	6.08	8.63	0	9.94	0	0.53	9.52	0.68	0.46	9.1	
	Twix Caramel	2.29	5.19	8.1	0	10.88	0	0	0	0	0	32.6	
	Raisinets	1.47	1.47	2.19	0	2.74	0	0	3.44	0	0	0	
	M&M Milk Choc	1.8	3.63	5.32	0	6.58	0.41	0.02	6.91	0	0	3.42	
	Choc Mars (Con)	2.13	1.03	1.57	0	2.04	0	0	0	0.35	0	6.49	
	Reeses PB Cups	1.68	1.62	3.5	0	4.73	0	0	3.58	0	0	3.96	
	Butterfinger	1.1	2.72	3.33	1.01	3.73	1.18	1.8	3.62	0.99	1.5	3.14	
	Choc Herhsey (Con)	1.22	2.87	1.46	6.28	1.61	1.34	7.09	1.26	3.34	6.48	2.22	
NONCHOC. CANDY	Skittles Original	1.03	0.12	0.14	0.01	0.07	0.23	0.11	0	0.22	0.04	0.51	
	Nonchoc Other (Con)	1.06	0.41	0.6	0	0.69	0.17	0.03	0	0.27	0	2.31	
	Twizzlers	1.66	1.16	1.14	1.08	0.91	2.1	0.61	1.75	1.27	0.67	0	
COOKIES	ZAnimal Cracker	1.9	0.29	0.37	0.05	0.62	0.02	0.13	0.53	0	0.17	0.25	
	CC Fam Amos	1.58	1.57	0	3.35	0.01	0	9.61	0.05	0	16.91	0	
	Ruger Wafer (Con)	1.6	0.54	0.62	0.25	0.47	1.05	0.08	0	1.18	0	2.2	
	Grandmas CC	1.15	0.84	0.39	1.9	0.36	0.18	2.58	0.55	0.67	2.43	0	
	Rasbry Knotts	0.68	1.1	0.4	2.75	0.44	0.34	3.3	0.63	1.15	3.08	0	
	Choc Fam Amos	0.91	1.35	1.1	1.86	1.37	0	2.95	1.39	0	2.94	1	
	Nabisco (Con)	1.23	1.44	1.14	2.08	1.29	0	3.14	1.97	0.05	2.86	0	
PASTRY	Pop-Tarts (Con)	2.42	0.27	0.37	0	0.38	0.21	0.09	0	0.24	0.02	1.48	
	Rice K Treats	0.85	2.25	2.61	1.24	2.06	4.3	0.33	3.44	3.19	0.25	0	
OTHER	Nature Valley (Con)	2.13	1.42	1.29	1.58	0.38	5	0	1.4	3.94	0	0	
	Planters (Con)	1.63	4.81	3.54	7.77	4.37	0	10.57	6.08	0.85	9.94	0	
	KarNuts (Con)	1.65	1.25	1.77	0	1.73	1.36	0	2.56	0.41	0	0	
	Farleys Fruit Snax	0.99	0.58	0.45	0.8	0.03	2.27	0	0.27	2.02	0	0.18	
	Cherry Fruit Snax	0.52	0.09	0.12	0	0.07	0.15	0.05	0	0.16	0	0.34	
	Cliff (Con)	3.91	1.03	1.22	0.41	1.35	0	1.26	1.84	0	1.04	0	
	Outside Good	25.34	28.58	29.44	24.49	27.09	38.12	18.85	34.84	31.17	17.44	11.96	

Figure E1: Estimated individual Shares  $s_{ij}$