

# Estimating Demand for Differentiated Products with Zeroes in Market Share Data\*

Amit Gandhi  
Airbnb

Zhentong Lu  
Bank of Canada

Xiaoxia Shi †  
UW-Madison

July 9, 2022

## Abstract

In this paper we introduce a new approach to estimating differentiated product demand systems that allows for products with zero sales in the data. Zeroes in demand are a common problem in differentiated product markets, but fall outside the scope of existing demand estimation techniques. We show that with a lower bound imposed on the expected sales quantities, we can construct upper and lower bounds for the conditional expectation of the inverse demand. These bounds can be translated into moment inequalities that are shown to yield consistent and asymptotically normal point estimators for demand parameters under natural conditions. In Monte Carlo simulations, we demonstrate that the new approach works well even when the fraction of zeroes is as high as 95%. We apply our estimator to supermarket scanner data and find that correcting the bias caused by zeroes has important empirical implications, e.g., price elasticities become twice as large when zeroes are properly controlled.

**Keywords:** Demand Estimation, Differentiated Products, Measurement Error, Moment Inequality, Zero

---

\*Previous version of this paper was circulated under the title “Estimating Demand for Differentiated Products with Error in Market Shares.”

†Corresponding author: Xiaoxia Shi, xshi@ssc.wisc.edu

‡We are thankful to Steven Berry, Jean-Pierre Dubé, Philip Haile, Bruce Hansen, Ulrich Müller, Aviv Nevo, Jack Porter, Dan Quint, and Chris Taber for insightful discussions and suggestions; We would also like to thank the participants at the 2011 MIT Econometrics of Demand Conference, Chicago-Booth Marketing Lunch, the 2012 Northwestern Conference on “Junior Festival on New Developments in Microeconometrics,” the 2012 Cowles Foundation Conference on “Structural Empirical Microeconomic Models,” 3rd Cornell - Penn State Econometrics & Industrial Organization Workshop (2016), the IAAE Annual Conference in 2021, as well as seminar participants at Wisconsin-Madison, Wisconsin-Milwaukee, Cornell, Indiana, Princeton, NYU, Penn, the Federal Trade Commission, and UCSD for their many helpful comments and questions. Xiaoxia Shi dedicates this paper to the loving memory of her mother Guizhi Ma who passed away in a car crash in December 2019.

# 1 Introduction

In this paper we introduce a new approach to differentiated product demand estimation that allows for zeroes in empirical market share data. Such zeroes are a highly prevalent feature of demand in a variety of empirical settings, ranging from workhorse retail scanner data, to data as diverse as homicide rates and international trade flows (we discuss these examples in further depth below). Zeroes naturally arise in “big data” applications which allow for increasingly granular views of consumers, products, and markets (see for example Quan and Williams (2015), Nurski and Verboven (2016)). Unfortunately, the standard estimation procedures using inverse demand function following the seminal Berry, Levinsohn, and Pakes (1995) (BLP for short) cannot be used in the presence of zero empirical shares - the inverse demand is simply not well defined at zeroes. Furthermore, ad hoc fixes to market zeroes that are sometimes used in practice, such as dropping zeroes from the data or replacing them with small positive numbers, are subject to biases which can be quite large (because the slope of the inverse demand is arbitrarily large around zero). This has left empirical work on demand for differentiated products without satisfying solutions to the zero shares problem, and often force researchers to aggregate their rich data on naturally defined products to crude artificial products which limits the type of questions that can be answered. This is the key problem that our paper aims to solve.

In this paper we provide an approach to estimating differentiated product demand models that provides consistency and asymptotic normality for demand parameters despite a possibly large presence of zero market shares in the data. We start by noting that the zeroes are caused by the wedge between the empirical shares ( $s_{jt}$ ) and the true choice probabilities ( $\pi_{jt}$ ): while the latter is always positive, the former can be zero because of sample noise. We show how the zeroes in empirical shares may not simply be a data anomaly, but an essential feature of markets with a rich product variety, even if the number of consumers ( $n_t$ ) is large. By market design, expected sales ( $n_t\pi_{jt}$ ) of some products do not increase with  $n_t$ , and as a result, their empirical market shares are zero with non-vanishing probabilities. We then show that by imposing a lower bound to the expected sales we can construct upper and lower bounds for the conditional expectation of the inverse demand. The bounds are used to construct a set of moment *inequalities* which are valid in the presence of the zeroes, and

more generally in the presence of sampling error in market shares.<sup>1</sup>

The moment inequalities can be directly used for parameter inference with the help of set inference methods in the econometrics literature. But for computational reasons, we give a point identification condition and propose a point estimator instead. We show that our point estimator is consistent so long as  $n_t$  is large and there is an exogenous product or market characteristic, or a group of them, that can identify a positive mass of observations whose latent choice probabilities are bounded sufficiently away from zero, e.g., product-market pairs for whom the observed market shares are not likely to be zero. This is natural in many applications (as illustrated in Section 5), and strictly generalizes the restrictions on choice probabilities for consistency under the traditional approach. Asymptotic normality then follows by similar arguments as those for censored regression models by Kahn and Tamer (2009).

Computationally, our estimator closely resembles the traditional approach with only a slight adjustment in how the empirical moments are constructed. In particular it is no more burdensome than the usual estimation procedures for BLP and can be implemented using either the standard nested fixed point method of the original BLP, or the MPEC method as advocated more recently by Dubé, Fox, and Su (2012).

We investigate the finite sample performance of the approach in a variety of mixed logit examples. We find that our estimator works well even when the the fraction of zeros is as high as 95%, while the standard procedure with the observations with zeroes deleted yields severely biased estimators even with mild or moderate fractions of zeroes.

We apply our bounds approach to widely used scanner data from the Dominicks Finer Foods (DFF) retail chain. In particular, we estimate demand for the tuna category as previously studied by Chevalier, Kashyap, and Rossi (2003) and continued by Nevo and Hatzitaskos (2006) in the context of testing the loss leader hypothesis of retail sales. We find that controlling for products with zero demand using our approach gives demand estimates that can be more than twice as elastic than standard estimates that select out the zeroes. We also show that the estimated price elasticities increase substantially during Lent (a high demand period for this product category) after we control for the zeroes. Both of these

---

<sup>1</sup>In the last couple of years, new aggregate demand models have been considered that accommodate zeroes in market share data in Dube, Hortacsu, and Joo (2020) and Lima (2021). Dube et. al. model the products with zero market shares as ones that are not in any consumer’s consideration set. Lima’s model rationalizes the zeros in market shares by restricting the support of the idiosyncratic taste shock. Neither paper deals with the sample noise issue in observed market shares. Since Dube et. al., Lima, and our paper’s methods rely on nonnested assumptions on the source of zeros, in practice, knowing the true source of zero is important for choosing the appropriate method. When in doubt, it is advisable to implement multiple methods and compare the results. A potentially interesting direction for future research is to combine those methods into a more generally applicable solution to the problem of zero market shares.

findings have implications for reconciling the loss-leader hypothesis with the data.

The plan of the paper is the following. In Section 2, we illustrate the stylized empirical pattern of Zipf’s law where market zeroes naturally arise. In Section 3, we describe our solution to the zeroes problem using a simple logit setup without random coefficients to make the essential matters transparent. In Section 4, we extend the moment inequality construction and our estimator to general discrete choice model possibly with random coefficients. Section 5 discusses the point identification condition. Sections 6 and 7 present the theoretical properties of the proposed estimator. Section 8 present results of Monte Carlo simulations and Section 9 presents the application to the DFF data, respectively. Section 10 concludes.

## 2 Market Zeroes

In this section, we highlight the empirical pattern of zeroes. Here we primarily use workhorse store level scanner data to illustrate these patterns. It is the same data that will also be used for our application. However we emphasize that our focus here on scanner data is only for the sake of a concrete illustration of the market zeroes problem - the key patterns we highlight in scanner data are also present in many other economic settings where demand estimation techniques are used (discussed further below and illustrated in the Appendix).

We employ here a widely studied store level scanner data set from the Dominick’s Finer Foods grocery chain, which is a public data set that has been used by many researchers.<sup>2</sup> The data comprise 93 Dominick’s Finer Foods stores in the Chicago metropolitan area over the years from 1989 to 1997. Like other store level scanner data sets, this data set provides demand information (price, sales, marketing) at store/week/UPC level, where a UPC (universal product code) is a unique bar code that identifies a natural product<sup>3</sup>.

Table 1 presents information on the resulting product variety across the different product categories in the data. The first column shows the number of products in an average store/week - the number of UPC’s can be seen varying from roughly 50 (e.g., bath tissue) to over four hundred (e.g., soft drinks) within even these narrowly defined categories. Thus there is considerable product variety in the data. The next two columns illustrate an important aspect of this large product variety: there are often just a few UPC’s that dominate each product category whereas most UPC’s are not frequently chosen. The second column

---

<sup>2</sup>For a complete list of papers using this data set, see the website of Dominick’s Database: <http://research.chicagobooth.edu/marketing/databases/dominicks/index.aspx>

<sup>3</sup>Store level scanner data can often be augmented with a panel of household level purchases (available, for example, through IRI or Nielsen). Although the DFF data do not contain this micro level data, the main points of our analysis are equally applicable to the case where household level data is available. Store level purchase data can be viewed as a special case household level data where all households are observationally identical (no observable individual level characteristics).

illustrates this pattern by showing the well known “80/20” rule that prevails in our data: we see that roughly 80 percent of the total quantity purchased in each category is driven by the top 20 percent of the UPC’s in the category. In contrast to these “top sellers”, the other 80 percent of UPC’s contain relatively “sparse sellers” that share the remaining 20 percent of the total volume in the category. The third column shows an important consequence of this sparsity: many UPC’s in a given week at a store simply do not sell. In particular, we see that the fraction of observations with zero sales can even be nearly 60% for some categories.

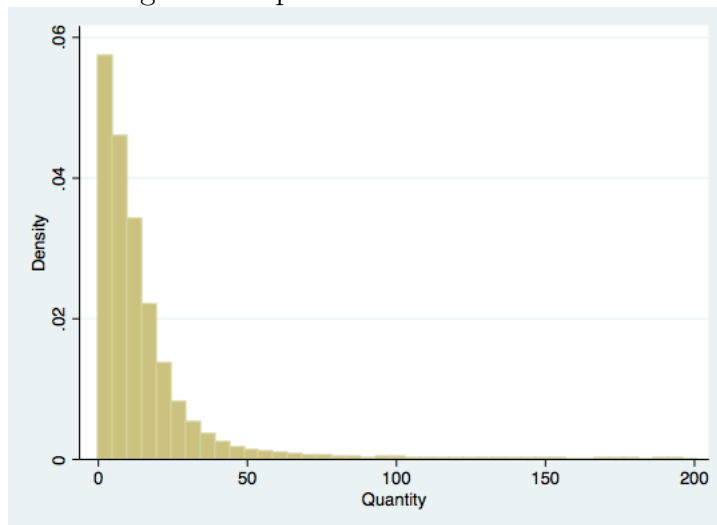
Table 1: Selected Product Categories in the Dominick’s Database

Category	Average Number of UPC’s in a Store/Week Pair	Percent of Total Sale of the Top 20% UPC’s	Percent of Zero Sales
Beer	179	87.18%	50.45%
Cereals	212	72.08%	27.14%
Crackers	112	81.63%	37.33%
Dish Detergent	115	69.04%	42.39%
Frozen Dinners	123	66.53%	38.32%
Frozen Juices	94	75.16%	23.54%
Laundry Detergents	200	65.52%	50.46%
Paper Towels	56	83.56%	48.27%
Refrigerated Juices	91	83.18%	27.83%
Soft Drinks	537	91.21%	38.54%
Snack Crackers	166	76.39%	34.53%
Soaps	140	77.26%	44.39%
Toothbrushes	137	73.69%	58.63%
Canned Tuna	118	82.74%	35.34%
Bathroom Tissues	50	84.06%	28.14%

We can visualize this situation in another way by fixing a product category (here we use canned tuna) and simply plotting the histogram of the volume sold for each week/UPC realization for a single store in the data. This frequency plot is given in *Figure 1*. As can be seen there is a sharp decay in the empirical frequency as the purchase quantity becomes larger, with a long thin tail.<sup>4</sup> In particular the bulk of UPC’s in the store have small purchase volume: the median UPC sells less than 10 units a week, which is less than 1.5%

<sup>4</sup>We plot the long tail pattern differently from a commonly seen illustration of power law using rank-size distribution (“size against rank or popularity”), but the difference is only cosmetic (basically flipping the x- and y-axis); the two ways of plotting convey the same information.

Figure 1: Zipf’s Law in Scanner Data



of the median volume of Tuna the store sells in a week. The mode of the frequency plot is a zero share.

This power-law decay in the frequency of product demand is often associated with “Zipf’s law” or the “the long tail”, which has a long history in empirical economics.<sup>5</sup> We present further illustrations of this long-tail demand pattern found in international trade flows as well as cross-county homicide rates in Appendix A, which provides a sense of the generality of these stylized facts.

The key takeaway from these illustrations is that the presence of market zeroes in the data is closely intertwined to the prevalence of power-law patterns of demand. We will exploit this relationship to place structure on the data generating process that underlies market zeroes.

### 3 A First Pass Through Logit Demand

Why do zero shares create a problem for demand estimation? In this section, we use the workhorse multinomial logit model to explain the zeroes problem and our solution. The general case is treated in the next section. In both cases, we assume that the econometrician observes a data set of  $\{(n_t, s_{jt}, x_{jt}) : j = 1, \dots, J_t, t = 1, \dots, T\}$ , where  $n_t$  is the number of potential consumers in market  $t$ ,  $s_{jt}$  is the fraction of those consumers choosing product  $j$ , and  $x_{jt}$  is the vector of observed characteristics of the product  $j$  and/or market  $t$  that often includes price,  $J_t$  is the number of inside products in market  $t$  and  $T$  is the number of markets. We focus on the case where there are many markets.

---

<sup>5</sup>See Anderson (2006) for a historical summary of Zipf’s law and many examples from the social and natural sciences. See Gabaix (1999) for an application of Zipf’s law to the economics literature.

### 3.1 Making Sense of The Zeroes

Consider a multinomial logit model for the demand of  $J_t$  products ( $j = 1, \dots, J_t$ ) and an outside option ( $j = 0$ ). A consumer  $i$  derives utility  $u_{ijt} = \delta_{jt} + \epsilon_{ijt}$  from product  $j$  in market  $t$ , where  $\delta_{jt}$  is the mean-utility of product  $j$  in market  $t$ , and  $\epsilon_{ijt}$  is the idiosyncratic taste shock that follows the type-I extreme value distribution. As is standard, the mean-utility  $\delta_{jt}$  of product  $j > 0$  is modeled as

$$\delta_{jt} = x'_{jt}\beta_0 + \xi_{jt}, \quad (1)$$

where  $\xi_{jt}$  is the unobserved characteristic. The outside good  $j = 0$  has mean utility normalized to  $\delta_{0t} = 0$ . The parameter of interest is  $\beta_0$ .

Each consumer chooses the product that yields the highest utility:

$$s_{ijt} = 1\{u_{jt} \geq u_{j't} \ \forall j' = 0, 1, \dots, J_t\}, \text{ for } j = 0, 1, \dots, J_t. \quad (2)$$

Aggregating consumers' choices, we obtain the true choice probability of product  $j$  in market  $t$ , denoted as

$$\pi_{jt} = \Pr(\text{product } j \text{ is chosen in market } t) = E[s_{ijt} | \delta_{1t}, \dots, \delta_{J_t t}].$$

The standard approach introduced by Berry (1994) for estimating  $\beta_0$  is to combine demand system inversion and instrumental variables.

First, for demand inversion, one uses the logit structure to find that

$$\delta_{jt} = \log(\pi_{jt}) - \log(\pi_{0t}), \text{ for } j = 1, \dots, J_t. \quad (3)$$

To handle the potential endogeneity of  $x_{jt}$  (i.e., its correlation with  $\xi_{jt}$ ), one finds some excluded instruments which along with the exogenous controls in  $x_{jt}$  form  $z_{jt}$  such that

$$E[\xi_{jt} | z_{jt}] = 0. \quad (4)$$

Then two stage least squares with  $\delta_{jt}$  defined in (3) as the dependent variable becomes the identification strategy for  $\beta_0$ .

Unfortunately  $\pi_{jt}$  is not observed as data - it is a theoretical choice probability defined by the model but only indirectly revealed through actual consumer choices. The standard approach to this following Berry (1994), Berry, Levinsohn, and Pakes (1995), and many subsequent papers in the literature has been to substitute  $s_{jt}$  the empirical market share for

$\pi_{jt}$ , where

$$s_{jt} = n_t^{-1} \sum_{i=1}^{n_t} s_{ijt} \text{ for } j = 0, 1, \dots, J_t, \quad (5)$$

and run a two-stage least square with  $\log(s_{jt}) - \log(s_{0t})$  as dependent variable,  $x_{jt}$  as covariates, and  $z_{jt}$  as instruments to obtain estimates for  $\beta_0$ . The theoretical justification used in the literature assume that  $n_t$  is large and importantly,  $\pi_{jt}$  either is bounded away from zero or converges to zero at a slower rate than  $1/n_t$ . Under these assumptions, Berry, Linton, and Pakes (2004) and Freyberger (2015) show that plugging in  $s_{jt}$  for  $\pi_{jt}$  at worst leads to a correctible bias.

However, for data sets with the power law pattern described in Section 2, a large proportion of the  $s_{jt}$ 's are zeroes. Substituting  $s_{jt}$  for  $\pi_{jt}$  is no longer feasible, and the theoretical assumptions used to justify that practice are no longer compatible with the data. The former is because  $\log(0)$  is not finite, and the later is because under the assumption that  $\pi_{jt}$  approaches zero at a slower rate than  $1/n_t$ , we have  $\Pr(s_{jt} = 0) \rightarrow 0$ , which is not consistent with the large number of zeroes in the data.

We rationalize the large number of zeros in  $s_{jt}$  at seemingly large  $n_t$  by allowing  $\pi_{jt}$  to approach zero at the rate of  $1/n_t$ . When  $\pi_{jt}$  approaches zero at this rate, for example,  $\pi_{jt} = c/n_t$  for a constant  $c > 0$ , we have

$$\lim_{n_t \rightarrow \infty} \Pr(s_{jt} = 0) = \lim_{n_t \rightarrow \infty} (1 - c/n_t)^{n_t} = \exp(-c). \quad (6)$$

Thus, zeroes arise naturally in this framework. In our bound construction below, we will assume a much weaker lower bound on  $\pi_{jt}$  than the existing literature:  $\pi_{jt} \geq \underline{\varepsilon}_1/n_t$  for some fixed constant  $\underline{\varepsilon}_1$ .

There is a simple supply side explanation for why the choice probability of some products should approach zero at the exact rate of  $1/n_t$  and why there may be a lower bound for  $n_t\pi_{jt}$ . A market with the power-law feature described in Section 2 may be thought of as one with a few dominant products that coexist with a competitive fringe (see e.g. Shimomura and Thisse (2012)). The fringe products enjoy free entry and exit and are subject to a fixed cost, denoted  $f_{jt}$ . The free entry and exit drives their expected profit to zero:

$$n_t\pi_{jt}m_{jt} - f_{jt} = 0, \quad (7)$$

where  $m_{jt}$  is the average mark-up. Then  $n_t\pi_{jt} = f_{jt}/m_{jt}$ . And  $\pi_{jt} \geq \underline{\varepsilon}_1/n_t$  holds for some  $\underline{\varepsilon}_1$



if there are a lower bound for  $f_{jt}$  and an upper bound for  $m_{jt}$ .<sup>6</sup> If there are also an upper bound for  $f_{jt}$  and a lower bound for  $m_{jt}$ , then  $\pi_{jt}$  approaches zero at the rate of  $1/n_t$ . The existence of such bounds are reasonable in differentiated product markets.<sup>7</sup>

### 3.2 Estimation Problem with Zeroes

As mentioned above, the zeroes pose an immediate challenge to estimation:  $\log(s_{jt})$  is  $-\infty$  when  $s_{jt} = 0$ . This makes the standard BLP estimator ill-defined. A common workaround is to ignore the  $(j,t)$ 's with  $s_{jt} = 0$ , effectively lumping those  $j$ 's into the outside option in market  $t$ . This however leads to a selection problem. To see this, suppose  $s_{jt} = 0$  for some  $(j, t)$  and one drops these observations from the analysis - effectively one is using a selected sample where the selection criterion is  $s_{jt} > 0$ . In this selected sample, the conditional mean of  $\xi_{jt}$  is no longer a constant. This is the well-known selection-on-unobservables problem and with such sample selection, an attenuation bias ensues.<sup>8</sup> The attenuation bias generally leads to demand estimates that appear to be too inelastic.<sup>9</sup>

Another commonly adopted empirical “trick” is to add a small positive number  $\epsilon > 0$  to the  $s_{jt}$ 's that are zero, and use the resulting modified shares  $s_{jt}^\epsilon > 0$  in place of  $\pi_{jt}$ .<sup>10</sup> However, this trick only treats the symptom, i.e.,  $s_{jt} = 0$ , but overlooks the nature of the problem: the true choice probability  $\pi_{jt}$  is small. And in this case, small estimation error in any estimator  $\hat{\pi}_{jt}$  of  $\pi_{jt}$  would lead to large error in the plugged-in version of  $\delta_{jt}$  and the estimation of  $\beta_0$ . This problem manifests itself directly because the estimate  $\hat{\beta}$  can be incredibly sensitive to the particular choice of the small number being added and there is little guidance on what is the “right” choice of the small number. In general, like selecting

---

<sup>6</sup>The calculation assumes single-product firms. Multi-product firms stop putting out new products sooner because they internalize the business stealing effect of new products on their existing products.

<sup>7</sup>The only bound that might be disputable is the lower bound for the average markup because markup is endogenous. But even that has some supporting evidence in the literature: Armstrong (2016) shows that the markup converges to a positive constant rather than zero when the number of firms grows to infinity.

<sup>8</sup>To see why  $E[\xi_{jt}|x_{jt}, s_{jt} > 0]$  is not a constant, consider two values of  $x_{jt}$ :  $x, x^*$  such that  $x'\beta > x^*\beta$ , and consider the homoskedastic case for simplicity. For each given value of  $x_{jt}$ , the criterion  $s_{jt} > 0$  selects high values of  $\xi_{jt}$  and leaves out low values of  $\xi_{jt}$ . Moreover, the selection is more severe for  $x^*$  than for  $x$  because the unobservable (to econometricians) needs to be more appealing to induce a positive observed market share when the observable characteristic is less appealing.

Thus, we should have

$$E[\xi_{jt}|x_{jt} = x^*, s_{jt} > 0] > E[\xi_{jt}|x_{jt} = x, s_{jt} > 0], \quad (8)$$

and clearly,  $E[\xi_{jt}|x_{jt}, s_{jt} > 0]$  is not a constant.

<sup>9</sup>It is easy to see that the selection bias is of the same direction if the selection criterion is instead  $s_{jt} > 0$  for all  $t$ , as one is effectively doing when focusing on a few top sellers that never demonstrate zero sales in the data. The reason is that the event  $s_{jt} > 0$  for all  $t$  contains the event  $s_{jt} > 0$  for a particular  $t$ . If the markets ( $\xi_{jt}$ 's) are independent, the particular  $t$  part of the selection dominates.

<sup>10</sup>Berry, Linton, and Pakes (2004) and Freyberger (2015) study the biasing effect of plugging in  $s_{jt}$  for  $\pi_{jt}$ . Their bias corrections do not apply when there are zeroes in the empirical shares.

away the zeroes, the “adding a small number trick” is also a biased estimator for  $\beta_0$ . We illustrate both biases in the Monte Carlo section (Section 8).

Despite their failure as general solutions, these “ad hoc zero fixes” have in them what could be a useful idea – Perhaps the variation among the non-zero share observations can be used to estimate the model parameters, while at the same time the presence of zeroes is controlled in such a way that avoids bias. We will present a new estimator that formalizes this possibility by using moment *inequalities* to control for the zeroes in the data while using the variation in the remaining part of the data to estimate the demand parameters.

### 3.3 Constructing Moment Inequalities

Our approach builds on two estimators of  $\log(\pi_{jt})$ . We refer to them as the upper and lower bounds of  $\log(\pi_{jt})$  because they bound  $\log(\pi_{jt})$  from above and below *on average* in the sense discussed below. These bounds are:

$$\log((n_t s_{jt} + \iota_u)/n_t) \text{ and } \log((n_t s_{jt} + \iota_\ell)/n_t) \quad (9)$$

where  $\iota_u$  and  $\iota_\ell$  are two positive numbers that we now construct.

To construct  $\iota_u$  and  $\iota_\ell$ , note that  $n_t s_{jt}$  follows a binomial distribution given  $n_t$  and  $\pi_{jt}$ :  $Bin(n_t, \pi_{jt})$ .<sup>11</sup> For each fixed  $n$  and  $\pi$ , and  $\iota \geq 0$ , define the function

$$f(\iota; n, n\pi) := E[\log(n_t s_{jt} + \iota) - \log(n_t \pi_{jt}) | n_t = n, \pi_{jt} = \pi].$$

The function  $f$  is negative infinity at  $\iota = 0$  (because  $s_{jt}$  can be 0 with positive probability), strictly increasing with  $\iota$ , and approaches positive infinity as  $\iota \rightarrow \infty$ . Therefore, at each  $n$  and  $\pi$ , the function crosses zero once and only once. We let the point of crossing be denoted  $\iota^*(n, n\pi)$ , which is defined implicitly by the equation:

$$f(\iota^*(n, n\pi); n, n\pi) = 0. \quad (10)$$

This quantity can be calculated because the function  $f(\iota; n, n\pi)$  (i.e., the expectation) can be calculated using the binomial distribution.

As explained in Section 3.1 above, we assume that  $n_t \pi_{jt}$  is bounded below by a small

---

<sup>11</sup>Here we maintain the standard assumption that in each given market, consumers’ choices are independent and identically distributed.

constant  $\underline{\varepsilon}_1 > 0$ , then we can define

$$\underline{\iota}_u := \sup_{n, \pi: n\pi \geq \underline{\varepsilon}_1} \iota^*(n, n\pi) \text{ and } \bar{\iota}_\ell := \inf_{n, \pi: n\pi \geq \underline{\varepsilon}_1} \iota^*(n, n\pi). \quad (11)$$

Furthermore, suppose that  $\underline{\iota}_u$  and  $\bar{\iota}_\ell$  are known and  $\underline{\iota}_u < \infty, \bar{\iota}_\ell > 0$  for now, which we will discuss shortly below. Then, if we let  $\iota_u$  and  $\iota_\ell$  be any finite number satisfying  $\iota_u \geq \underline{\iota}_u$  and  $0 < \iota_\ell \leq \bar{\iota}_\ell$ , we will have

$$E[\log((n_t s_{jt} + \iota_u)/n_t) - \log(\pi_{jt}) | z_{jt}] \geq 0 \text{ and } E[\log((n_t s_{jt} + \iota_\ell)/n_t) - \log(\pi_{jt}) | z_{jt}] \leq 0. \quad (12)$$

Combining this with the orthogonality condition  $E[\xi_{jt} | z_{jt}] = 0$ , we obtain a set of conditional moment inequalities

$$\begin{aligned} E[\log(n_t s_{jt} + \iota_u)/n_t - \log(\pi_{0t}) - x'_{jt} \beta_0 | z_{jt}] &\geq 0 \\ E[\log(n_t s_{jt} + \iota_\ell)/n_t - \log(\pi_{0t}) - x'_{jt} \beta_0 | z_{jt}] &\leq 0. \end{aligned} \quad (13)$$

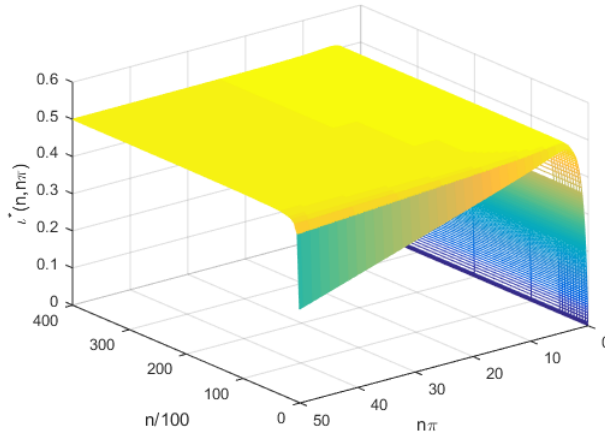
The piece  $\log(\pi_{0t})$  is easy to estimate because  $\pi_{0t}$  is typically large (sufficiently distant from zero) in most empirical work. We can plug in  $s_{0t}$  or any modification  $\tilde{s}_{0t}$  of  $s_{0t}$  for  $\pi_{0t}$ . As long as the modification is negligible relative to the estimation error in  $s_{0t}$ , standard arguments will imply  $T^{-1} \sum_{t=1}^T [\log(\tilde{s}_{0t}) - \log(\pi_{0t})] = o_p(1)$ . We specify  $\tilde{s}_{0t}$  in the general case later. For the logit case,  $\tilde{s}_{0t} = s_{0t}$  works just fine.

Now we discuss the choice of  $\iota_\ell$  and  $\iota_u$  in greater detail. The first two questions we seek to answer are whether  $\bar{\iota}_\ell$  is positive and  $\underline{\iota}_u$  is finite, and whether we know them without the knowledge of the lower bound  $\underline{\varepsilon}_1$  for  $n_t \pi_{jt}$ . The third question is how to choose  $\iota_\ell$  and  $\iota_u$  given our answers to the first two questions.

We answer the first two questions by numerically obtaining  $\iota^*(n, n\pi)$  for a large representative set of values of  $n$  and  $n\pi$  and plot them in Figure 2.<sup>12</sup> The figure shows that  $\iota^*(n, n\pi)$  varies smoothly with its two arguments, which gives us confidence that the supremum and the infimum from these discrete values are close to those of the function. Specifically, Figure 2 shows that  $\underline{\iota}_u \approx 0.5$  and it is not affected by  $\underline{\varepsilon}_1$ . For  $\bar{\iota}_\ell$ , the figure shows that it approaches zero as  $\underline{\varepsilon}_1$  approaches zero. Thus, without knowing  $\underline{\varepsilon}_1$ , we do not know  $\bar{\iota}_\ell$ . Nevertheless, the calculation that leads to Figure 2 also produces Table 2, which gives us an idea of how  $\bar{\iota}_\ell$

<sup>12</sup>We considered the values:  $n \in \{100, 500, 1000, 1500, 2000, 2500, 3000, 3500, 4000, 4500, 5000, 10000, 20000, 40000\}$  and  $n\pi \in \{0.0001 : 0.0001 : 0.01, 0.02 : .01 : 1, 1 : 0.1 : 50\}$ , where the numbers between semicolons are the step sizes.

Figure 2:  $\iota^*(n, n\pi)$  for a Range of  $n$  and  $n\pi$  Values



changes with  $\underline{\varepsilon}_1$ . As the table shows, when  $\underline{\varepsilon}_1$  is very small,  $\bar{\iota}_\ell$  is well approximated by  $\underline{\varepsilon}_1$ .<sup>13</sup>

Table 2: Computed  $\bar{\iota}_\ell$  for Various Values of  $\underline{\varepsilon}_1$

$\underline{\varepsilon}_1$	$n$	$\bar{\iota}_\ell$
$\geq 0.5$	$\in [100, 40000]$	$\geq 0.250$
0.1	$\in [100, 40000]$	0.0776
0.01	$\in [100, 40000]$	0.00955
0.001	$\in [100, 40000]$	0.000993
0.0001	$\in [100, 40000]$	0.0000998

Given what we have learned about  $\iota_u$  and  $\bar{\iota}_\ell$ , we recommend choosing  $\iota_u$  and  $\iota_\ell$  as follows. For  $\iota_u$ , any  $\iota_u > \iota_u$  works in theory, but for better finite sample property, we recommend an  $\iota_u$  a bit larger. In the Monte Carlo simulations we find that  $\iota_u = 2$  works well. Moreover, using  $\iota_u = 2$  has an added benefit: it not only satisfies the theoretical requirement for the logit model, but also satisfies the requirement for non-logit based models, as we will see in Section ??.

For  $\iota_\ell$ , one can make a guess about how small  $\underline{\varepsilon}_1$  can be based on institutional knowledge, and simply use an  $\iota_\ell$  that is smaller than this number. In practice, it sometimes is not difficult to make an educated guess of  $\underline{\varepsilon}_1$  when you realize that  $\underline{\varepsilon}_1$  is the lowest *number of units* that one expects a product to sell in a market. For example, if the market unit is week and the product is a particular yogurt, the supermarket probably will not put it on the shelf if it is

<sup>13</sup>Complete analytical investigation of the shape of  $\iota^*(n, n\pi)$  is difficult due to the lack of analytical solution to expectations of the logarithm of binomial random variables. However, we provide some partial answers by analytically deriving the limit of  $\iota^*(n, n\pi)$  as  $n\pi$  approaches infinity and that as  $n\pi$  approaches zero in Appendix E. These limits are consistent with the numerical results reported in Figure 2 and Table 2.

expected to sell less than one unit per 100 weeks. That gives us a lower bound  $\underline{\varepsilon}_1 = 0.01$ .

What if one makes a wrong guess at the lowest number of sales? Over-guessing can cause violations of the moment inequalities (12), but fortunately, under-guessing does *not*. Setting  $\iota_\ell$  at a value much lower than the actual  $\bar{\iota}_\ell$  can guarantee the validity of (12). In our Monte Carlo and application, we in fact use an extremely low  $\iota_\ell = 2^{-52}$  just to be on the safe side. As we see in the Monte Carlo and the empirical application, the estimates have good precision despite the extremely small  $\iota_\ell$  used.<sup>14</sup>

### 3.4 Point Estimation

One can use any of the inference procedures for moment inequality models on (13), for example, Andrews and Shi (2013) and Cox and Shi (2019). Point identification is not required. On the other hand, point identification can greatly reduce the computational cost because inference without point identification generally requires costly test inversion. This is especially important for more complicated demand models than multinomial logit where even standard BLP estimation is computationally nontrivial.

In later sections, we discuss conditions that guarantee point identification. Under those conditions, the inequalities in (13) hold as equalities asymptotically on a set of  $z_{jt}$  values of positive measure, and ensure point identification in the same spirit as Kahn and Tamer (2009) in the context of endogenously censored regression models. To capture the identification information provided by those  $z_{jt}$  values, we consider a countable collection  $\mathcal{G}$  of instrumental indicator functions  $g : R^{d_z} \rightarrow \{0, 1\}$ , where  $d_z$  is the dimension of  $z_{jt}$ . We adopt the collections of instrumental functions in Andrews and Shi (2013). Such collections are shown therein to preserve all the identification information in the conditional moment inequality model (13), and thus they preserve the point-identification provided by the set of  $z_{jt}$  values at which the inequalities asymptotically hold as equalities, without that set of values being known. An example of  $\mathcal{G}$  is given below.

We form the sample moments

$$\begin{aligned} \bar{m}_T^u(\beta, g) &:= (T\bar{J}_T)^{-1} \sum_{t=1}^T \sum_{j=1}^{J_t} (\hat{\delta}_{jt}^u - x'_{jt}\beta)g(z_{jt}) \text{ and} \\ \bar{m}_T^\ell(\beta, g) &:= (T\bar{J}_T)^{-1} \sum_{t=1}^T \sum_{j=1}^{J_t} (\hat{\delta}_{jt}^\ell - x'_{jt}\beta)g(z_{jt}), \end{aligned}$$

---

<sup>14</sup>We note that this is generally true if the point identification condition in Section 5 below holds and  $n_t$  is large. But if the point identification condition does not hold or  $n_t$  is small, too small an  $\iota_\ell$  can affect the precision of the inference. In that case, one should use institutional knowledge to carefully determine  $\underline{\varepsilon}_1$ —the lower bound for  $n_t\pi_{jt}$ , subsequently determine  $\underline{\iota}_\ell$  according to Table 2, and choose  $\iota_\ell = \underline{\iota}_\ell$ .

where  $\bar{J}_T = T^{-1} \sum_{t=1}^T J_t$ ,  $\hat{\delta}_{jt}^u = \log((n_t s_{jt} + \iota_u)/n_t)$  and  $\hat{\delta}_{jt}^\ell = \log((n_t s_{jt} + \iota_\ell)/n_t)$ . These moments are used to form the criterion function:

$$\widehat{Q}_T(\beta) = \sum_{g \in \mathcal{G}} \mu(g) \{ [\bar{m}_T^u(\beta, g)]_-^2 + [\bar{m}_T^\ell(\beta, g)]_+^2 \}, \quad (14)$$

where  $\mu(g) : \mathcal{G} \rightarrow [0, 1]$  is a probability mass function on  $\mathcal{G}$ ,  $[x]_- = \min\{0, x\}$  and  $[x]_+ = \max\{0, x\}$ . Finally, we define the estimator  $\widehat{\beta}_T$  to be the minimizer of  $\widehat{Q}_T(\beta)$ :

$$\widehat{\beta}_T = \arg \min_{\beta \in B} \widehat{Q}_T(\beta), \quad (15)$$

where  $B$  is the parameter space of  $\beta$ . As we can see, computation of this estimator is on par with the standard GMM estimator for the multinomial logit model.

For  $\mathcal{G}$ , we divide the instrument vector  $z_{jt}$  into discrete instruments,  $z_{d,jt}$ , and continuous instruments  $z_{c,jt}$ . Without loss of generality assume that  $z_{c,jt}$  lies in  $[0, 1]^{d_{z_c}}$ .<sup>15</sup> Let the set  $\mathcal{Z}_d$  be the discrete set of values that  $z_{d,jt}$  can take. The set  $\mathcal{G}$  is defined as

$$\begin{aligned} \mathcal{G} &= \{g_{a,r,\zeta}(z_d, z_c) = 1((z'_c, z'_d)' \in C_{a,r,\zeta}) : C_{a,r,\zeta} \in \mathcal{C}\}, \text{ where} \\ \mathcal{C} &= \{(\times_{u=1}^{d_{z_c}} ((a_u - 1)/(2r), a_u/(2r))) \times \{\zeta\} : a_u \in \{1, 2, \dots, 2r\}, \text{ for } u = 1, \dots, d_{z_c}, \\ &\quad r = r_0, r_0 + 1, \dots, \text{ and } \zeta \in \mathcal{Z}_d\}. \end{aligned} \quad (16)$$

In practice, we truncate  $r$  at a finite value  $\bar{r}_T$ .<sup>16</sup> This does not affect the first order asymptotic property of our estimator as long as  $\bar{r}_T \rightarrow \infty$  as  $T \rightarrow \infty$ . For  $\mu(\cdot)$ , we use

$$\mu(\{g_{a,r,\zeta}\}) \propto (100 + r)^{-2} (2r)^{-d_{z_c}} K_d^{-1}, \quad (17)$$

where  $K_d$  is the number of elements in  $\mathcal{Z}_d$ . The same  $\mu$  measure is used and works well in Andrews and Shi (2013).<sup>17</sup>

<sup>15</sup>If not, we can normalize it to lie in  $[0, 1]$  as suggested in Andrews and Shi (2013). For example, we can let  $\tilde{z}_{c,jt} = F_{N(0,1)}(\widehat{\Sigma}_{z_c}^{-1/2} z_{c,jt})$ , where  $F_{N(0,1)}(\cdot)$  is the standard normal cdf and  $\widehat{\Sigma}_{z_c}$  is the sample covariance matrix of  $z_{c,jt}$ , and use  $\tilde{z}_{c,jt}$  in place of  $z_{c,jt}$  to construct the instrumental functions.

<sup>16</sup>We shall show some simulation results in the Monte Carlo section that provides useful guidance on choosing  $\bar{r}_T$  (and other ways of keeping the dimension of  $\mathcal{G}$  manageable) in practice.

<sup>17</sup>Note that appropriate choices of  $\mathcal{G}$  and  $\mu$  are not unique. For other possible choices, see Andrews and Shi (2013).

## 4 The General Model

Now we extend our discussion to the general differentiated product demand model and present our parameter estimator.

### 4.1 Setup

The specification of the general model is the same as the logit model except that the consumer level shock  $\epsilon_{ijt}$  in  $u_{ijt} = \delta_{jt} + \epsilon_{ijt} \equiv x'_{jt}\beta + \xi_{jt} + \epsilon_{ijt}$  is no longer type-I extreme value distribution. Instead, we assume that

$$\epsilon_{it} = (\epsilon_{i0t}, \dots, \epsilon_{iJ_t t}) \sim F(\cdot | x_t; \lambda), \quad (18)$$

where  $x_t$  stands for  $(x'_{1t}, \dots, x'_{J_t t})'$ , and  $F(\cdot | x_t, \lambda)$  is a conditional cumulative distribution function known up to the finite dimensional unknown parameter  $\lambda$ . By allowing  $x_t$  and an unknown parameter to enter the distribution of  $\epsilon_{ijt}$ , this specification is general enough to encompass most models used in empirical work. In particular, it encompasses the random coefficient specifications  $\epsilon_{ijt} = x'_{jt}(\beta_i - \beta) + \nu_{ijt}$ , where  $\beta_i$  is a vector of random coefficients that follows a distribution (e.g., joint normal) known up to some unknown parameter,  $\nu_{ijt}$  is the idiosyncratic taste shock.<sup>18</sup>

Given the specification, the unknown parameter in the general model is  $\theta = (\beta', \lambda)'$ . For clarity, we use  $\theta_0 \equiv (\beta'_0, \lambda'_0)'$  to denote the true value of  $\theta$ . Let  $B \subseteq R^{d_\beta}$  denote the parameter space of  $\beta$ , and  $\Lambda \subseteq R^{d_\lambda}$  the parameter space of  $\lambda$ . Let  $\Theta = B \times \Lambda$  be the parameter space of  $\theta$ .

In this model, the choice probability of each product is determined by:

$$\pi_{jt} = \int 1\{\delta_{jt} + \epsilon_j \geq \max_{j'=0,1,\dots,J_t} (\delta_{j't} + \epsilon_{j'})\} dF(\epsilon_0, \epsilon_1, \dots, \epsilon_{J_t} | x_t, \lambda_0), \quad j = 0, 1, \dots, J_t. \quad (19)$$

Let  $\pi_t = (\pi_{1t}, \dots, \pi_{J_t t})'$ . This system is invertible under the connected substitute condition in Berry, Gandhi, and Haile (2013). In other words, we can define the inverse demand function  $\delta_t(\pi_t, \lambda) := (\delta_{jt}(\pi_t, \lambda))_{j=1}^{J_t}$  as the solution to the equation system

$$\pi_{jt} = \int 1\{\delta_{jt}(\pi_t, \lambda) + \epsilon_j \geq \max_{j'=0,1,\dots,J_t} (\delta_{j't}(\pi_t, \lambda) + \epsilon_{j'})\} dF(\epsilon_0, \epsilon_1, \dots, \epsilon_{J_t} | x_t, \lambda), \quad j = 1, \dots, J_t. \quad (20)$$

---

<sup>18</sup>Requiring  $F(\cdot | x_t, \lambda)$  to be known up to a finite dimensional parameter rules out the vertical model (see Berry and Pakes (2007)) because for the vertical model,  $\epsilon_{ijt}$  is a function of the unobservable product characteristics (quality).

Inverting the demand system allows for the use of instrumental variables to identify  $\theta$  based on the exclusion restriction:

$$E[\xi_{jt} | z_{jt}] = 0. \quad (21)$$

where  $z_{jt}$  is a vector of exogenous variables including exogenous components of  $x_{jt}$  and excluded instruments if there are any. This is because one can then obtain the following moment restriction:

$$E[\delta_{jt}(\pi_t, \lambda_0) - x'_{jt}\beta_0 | z_{jt}] = 0. \quad (22)$$

If  $\pi_t$  were observed, the parameter  $\theta$  in the model would be identified under standard GMM identification conditions. However, as discussed in the logit case,  $\pi_t$  is not observed. Instead only a noisy measure  $s_t := (s_{1t}, \dots, s_{J_t})'$  is, and  $s_t$  frequently contains zero elements in many commonly used data sets. As in the logit case,  $\delta_t(s_t, \lambda)$  is typically not well defined when  $s_t$  contains zero elements, and thus simply substituting  $s_t$  for  $\pi_t$  in the moment conditions (22) is problematic.

## 4.2 Bound Construction

Like in the logit case, we construct a pair of functions:  $\hat{\delta}_{jt}^u(s_t, \lambda)$  and  $\hat{\delta}_{jt}^\ell(s_t, \lambda)$ , to form bounds for  $\delta_{jt}(\pi_t, \lambda)$ . The construction is based on the bounds for the logit case but adjusts for the different functional form:

$$\begin{aligned} \hat{\delta}_{jt}^u(s_t, \lambda) &= \log((n_t s_{jt} + \iota_u)/n_t) + \delta_{jt}(\tilde{s}_t, \lambda) - \log(\tilde{s}_{jt}), \\ \hat{\delta}_{jt}^\ell(s_t, \lambda) &= \log((n_t s_{jt} + \iota_\ell)/n_t) + \delta_{jt}(\tilde{s}_t, \lambda) - \log(\tilde{s}_{jt}), \end{aligned} \quad (23)$$

where  $\iota_\ell$  and  $\iota_u$  are fixed numbers, and  $\tilde{s}_t$  is a slight modification of  $s_t$  to take it off the boundary of the probability simplex. We will require that the modification of  $\tilde{s}_{jt}$  to  $s_{jt}$  is small so that  $\|\tilde{s}_t - s_t\| = O_p(1/n_t)$ . For example  $\tilde{s}_{jt} = s_{jt} + 1/n_t$  (when  $J_t$  is bounded) or  $\tilde{s}_{jt} = s_{jt} + 1/(n_t J_t)$  (when  $J_t$  is unrestricted) for  $j = 1, \dots, J_t$ .<sup>19</sup>

To see why the construction in (23) may be valid and what new requirements we may need on  $\iota_u$  and  $\iota_\ell$  if any, consider for example, the upper bound:

$$\hat{\delta}_{jt}^u(s_t, \lambda) - \delta_{jt}(\pi_{jt}, \lambda) = [\log((n_t s_{jt} + \iota_u)/n_t) - \log(\pi_{jt})]$$

---

<sup>19</sup>We note that this implies  $\tilde{s}_{0t} = s_{0t} - J_t/n_t$  or  $\tilde{s}_{0t} = s_{0t} - 1/n_t$ . This in principle could be less than or equal to zero. But in typical data sets, this is not an issue because  $s_{0t}$  is much larger than  $J_t/n_t$ . It is not an issue asymptotically as we will assume that  $\pi_{0t}$ —the share of the outside good—is bounded away from zero.



$$+ [(\delta_{jt}(\tilde{s}_t, \lambda) - \log(\tilde{s}_{jt})) - (\delta_{jt}(\pi_t, \lambda) - \log(\pi_{jt}))].$$

We already know from the logit case that the first summand is nonnegative in expectation conditional on  $\pi_{jt}$  as long as  $\iota_u \geq \underline{\iota}_u$  for  $\underline{\iota}_u$  defined in equation (11). It is then clear that the bound  $\hat{\delta}_{jt}^u(s_t, \lambda)$  will be asymptotically valid if either (i) the second summand is asymptotically negligible, or (ii) the conditional expectation of the second summand can be bounded from above by that of the first. Next we show that the first case applies to logit-based models, while the second case applies to models where the idiosyncratic error has a thinner tail than the logistic distribution, for example, normal distributions.

### When $\delta_{jt}(\cdot, \lambda) - \log(\cdot)_j$ is Uniformly Continuous

Let  $\Delta_{J_t}^0$  denote a subset of  $\{\pi \in (0, 1)^{J_t} : \mathbf{1}'_{J_t} \pi < 1\}$  that  $\pi_t$  can take value in. Let  $\Delta_{J_t}^c$  denote an  $c$ -expansion of  $\Delta_{J_t}^0$ , that is,  $\Delta_{J_t}^c = \{\pi \in (0, 1)^{J_t} : \pi' \mathbf{1}_{J_t} < 1, \min_{p \in \Delta_{J_t}^0} \|p - \pi\|_f \leq c\}$  for  $c > 0$ , where  $\|p - \pi\|_f = \sqrt{\|p - \pi\|^2 + (\mathbf{1}'(p - \pi))^2}$ . Note that the metric  $\|\cdot\|_f$  takes into account the difference for the outside share, while the Euclidean norm on  $\{\pi \in (0, 1)^{J_t} : \mathbf{1}'_{J_t} \pi < 1\}$  only considers the shares for the inside goods.

Define the function  $\check{\delta}_t(\cdot, \lambda) = (\check{\delta}_{1t}(\cdot, \lambda), \dots, \check{\delta}_{J_t t}(\cdot, \lambda))' : \Delta_{J_t}^c \rightarrow R^{J_t}$  where

$$\check{\delta}_{jt}(\pi, \lambda) := \delta_{jt}(\pi, \lambda) - \log(\pi_j).$$

Since  $\Delta_{J_t}^c$  may contain points arbitrarily close to the boundary of the probability simplex, in general neither  $\delta_{jt}(\cdot, \lambda)$  nor  $\log(\cdot)_j$  is uniformly continuous on  $\Delta_{J_t}^c$ . Thus, neither  $\delta_{jt}(\tilde{s}_t, \lambda) - \delta_{jt}(\pi_t, \lambda)$  nor  $\log(\tilde{s}_{jt}) - \log(\pi_{jt})$  may converge to zero as  $n_t \rightarrow \infty$  and  $\pi_{jt} \rightarrow 0$  even if  $\tilde{s}_t$  is the most efficient consistent estimate of  $\pi_t$ . However, in many models used in empirical work the logit inverse demand ( $\log(\pi_j) - \log(\pi_0)$ ) is a good first-order approximation of  $\delta_{jt}(\pi, \lambda)$  when  $\pi_j$  is close to zero and this first order term is the entire reason that the inverse demand is not uniformly continuous. For such models, the following assumption is reasonable:

**Assumption 1.** (a) For some  $c > 0$ ,  $\max_{t=1, \dots, T; j=1, \dots, J_t} \sup_{\pi, \tilde{\pi} \in \Delta_{J_t}^c : \pi \neq \tilde{\pi}} \sup_{\lambda \in \Lambda} \frac{|\check{\delta}_{jt}(\tilde{\pi}, \lambda) - \check{\delta}_{jt}(\pi, \lambda)|}{\|\tilde{\pi} - \pi\|_f \sqrt{J_t}} \leq O(1)$ .

(b)  $0 < \iota_\ell \leq \bar{\iota}_\ell$  and  $\underline{\iota}_u \leq \iota_u < \infty$ , where  $\bar{\iota}_\ell$  and  $\underline{\iota}_u$  are defined in equation (11), and  $\sup_{t=1, \dots, T} n_t \|\tilde{s}_t - s_t\|_f = O_p(1)$ .

Now we give two examples where Assumption 1(a) is satisfied.

**Example 4.1.** Nested Logit. The inverse demand of the nested logit model can be written as  $\delta_{jt}(\pi_t, \lambda) = \log(\pi_{jt}/\pi_{0t}) - \lambda \log(\pi_{gt}/\pi_{0t})$  where  $\pi_{gt}$  is the aggregate share of all the products in the nest (nest  $g$ ) that  $j$  is in. In this case,  $\check{\delta}_{jt}(\pi_t, \lambda) = (\lambda - 1) \log \pi_{0t} - \lambda \log(\pi_{gt})$ . Assumption

1(a) is satisfied if  $\Delta_{J_t}^0 = \{\pi \in (0, 1)^{J_t} : 1 - \mathbf{1}'_{J_t} \pi > \underline{\varepsilon}_0, \pi_{gt} > \underline{\varepsilon}_0 \text{ for all nests } g\}$ . In fact, Assumption 1(a) holds without the  $\sqrt{J_t}$ , which is a stronger version of the assumption. The requirement that  $\pi_{0t}$  and  $\pi_{gt}$  are bounded away from zero is reasonable for data sets in which neither the outside good nor any of the nests have zero shares.

**Example 4.2.** Random Coefficient Logit. For the random coefficient logit model,  $\delta_{jt}(\pi_t; \lambda)$  is the solution to the following equation system:

$$\pi_{jt} = \exp(\delta_{jt}) \int \frac{\exp(w'_{jt}v)}{1 + \sum_{k=1}^{J_t} \exp(\delta_{kt} + w'_{kt}v)} dF(v; \lambda), \quad j = 1, \dots, J_t,$$

where  $w_{jt}$  is a vector of covariates with random coefficients, and  $F(\cdot; \lambda)$  is the distribution of the random coefficient known up to the unknown parameter  $\lambda$ . Using the definition of  $\check{\delta}_{jt}$  above, we can write

$$\exp(-\check{\delta}_{jt}(\pi_t; \lambda)) = \int \frac{\exp(w'_{jt}v)}{1 + \sum_{k=1}^{J_t} \exp(\check{\delta}_{kt}(\pi_t; \lambda) + w'_{kt}v)\pi_{kt}} dF(v; \lambda). \quad (24)$$

Assume that  $\|w_{jt}\|$  is bounded by  $\bar{w}$  and  $0 < \sup_{w: \|w\| \leq \bar{w}} \int \exp(w'v) dF(v; \lambda) < \infty$ . We can already see that  $\check{\delta}_{jt}(\pi_t; \lambda)$  is bounded away from  $-\infty$  when  $\pi_{jt} \rightarrow 0$  (in which case,  $\delta_{jt}(\pi_t; \lambda) \rightarrow -\infty$ ). With additional algebra, we can show that  $\partial \check{\delta}_{jt}(\pi_t; \lambda) / \partial \pi_t$  is bounded, which essentially guarantees Assumption 1(a). The details are given in Appendix D.

Under Assumption 1(a), the requirement for  $\iota_u$  and  $\iota_\ell$  are the same as in the logit case, which is formally stated in Assumption 1(b).

### When $\delta_{jt}(\cdot, \lambda) - \log(\cdot)$ is Not Uniformly Continuous

In some models used in empirical work, Assumption 1 can fail to hold. For example, if the model is a simple probit with  $J_t = 1$ ,  $\delta_t(\pi) = \Phi^{-1}(\pi)$ , where  $\Phi^{-1}$  is the inverse of the standard normal cdf. In this case,  $\check{\delta}_t(\pi) = \delta_t(\pi) - \log(\pi) = \Phi^{-1}(\pi) - \log \pi$ . This function approaches  $+\infty$  when  $\pi \rightarrow 0$ , and has arbitrarily large slope near zero. For such cases, an alternative assumption may be reasonable and this is given in the following Assumption.

- Assumption 2.** (a)  $\max_{j,t} E[\check{\delta}_{jt}(\tilde{s}_t, \lambda_0) - \check{\delta}_{jt}(\pi_t, \lambda_0) | \pi_t, z_t] \leq 0$ ,  
 (b)  $\min_{j,t} E[\delta_{jt}(\tilde{s}_t, \lambda_0) - \delta_{jt}(\pi_t, \lambda_0) | \pi_t, z_t] \geq 0$ , and  
 (c) for  $j = 1, \dots, J_t$ ,  $\tilde{s}_{jt} = s_{jt} + 1/n_t$ ,  $0 < \iota_\ell \leq \bar{\iota}_\ell$  and  $1 < \iota_u < \infty$ .  
 (d)  $\sup_j \sup_{\pi: \pi_j \geq (\varepsilon_1 \wedge 1)/n_t} |\delta_{jt}(\pi, \lambda_0)| \leq C_0 \log(n_t)$  for a constant  $C_0 > 0$  for all  $t$ .

**Example 4.3.** Binary Probit. For binary probit model, we verified numerically that parts (a)-(b) hold given part (c), even though we do not have a theoretical proof. The intuition is that  $\Phi^{-1}(\pi)$  decreases slower than  $\log(\pi)$  as  $\pi \rightarrow 0$ . Thus, smaller modification  $\iota$  is needed for  $E[\Phi^{-1}(s_{jt} + \iota/n_t)|\pi_t, z_t]$  to exceed  $E[\Phi^{-1}(\pi_{jt})|\pi_t, z_t]$  than for  $E[\log(s_{jt} + \iota/n_t)|\pi_t, z_t]$  to exceed  $E[\log(\pi_{jt})|\pi_t, z_t]$ . And  $\iota = 1$  is sufficient for the latter, as we discussed in the logit case. Part (d) holds simply because of the shape of  $\Phi^{-1}(\cdot)$  which increases slower than  $\log(\cdot)$  as the argument decreases to zero.

The following lemma shows that the bounds constructed in (23) are asymptotically valid:

**Lemma 1.** *Suppose that  $\min_{t=1, \dots, T} n_t \rightarrow \infty$  as  $T \rightarrow \infty$ , that  $n_t \pi_{jt} \geq \underline{\varepsilon}_1$  for  $j = 1, \dots, J_t$  and  $\underline{\varepsilon}_1$  being the positive number used in (11), and that  $E[\xi_{jt}|z_{jt}] = 0$ . If either Assumption 1 or Assumption 2(a)-(c) holds. Then, there exist random variables  $e_{jt}^u$  and  $e_{jt}^\ell$  such that  $\sup_{j=1, \dots, J_t; t=1, \dots, T} \frac{n_t^{1/2}}{T^{1/4} J_t^{1/2}} |e_{jt}^y| = O_p(1)$  for  $y = u, \ell$ , and*

$$\begin{aligned} E[\hat{\delta}_{jt}^u(s_t, \lambda_0) - x_{jt}\beta_0 + e_{jt}^u|z_{jt}] &\geq 0 \\ E[\hat{\delta}_{jt}^\ell(s_t, \lambda_0) - x_{jt}\beta_0 + e_{jt}^\ell|z_{jt}] &\leq 0. \end{aligned} \quad (25)$$

The moment inequalities (25) can be taken to the data since the term  $e_{jt}^y$  ( $y = u, \ell$ ) is ignorable provided that  $n_t$  increases at a faster rate than  $T^{1/2}J_t$ . Also note that for the multinomial logit and the nested logit case, the lemma holds without the  $J_t^{1/2}$  in the denominator because Assumption 1(a) holds without the  $J_t^{1/2}$  in  $\sup_{j=1, \dots, J_t; t=1, \dots, T} \frac{n_t^{1/2}}{T^{1/4} J_t^{1/2}} |e_{jt}^y| = O_p(1)$ . Thus for these models  $n_t$  only needs to increase faster than  $T^{1/2}$ .

### 4.3 Point Estimation

Like in the logit case, one can apply any of the inference procedures for moment inequality models on (25). Yet point identification can greatly simplify computation. Point identification conditions are given in Section 5. Under those conditions the inequalities in (25) hold as equalities asymptotically on a set of  $z_{jt}$  values with positive measure.

We define the point estimator analogous to the logit case:

$$\hat{\theta}_T := (\hat{\beta}'_T, \hat{\lambda}'_T)' = \arg \min_{\theta \in \Theta} \hat{Q}_T(\theta), \quad (26)$$

where

$$\hat{Q}_T(\theta) = \sum_{g \in \mathcal{G}} \mu(g) \{ [\bar{m}_T^u(\theta, g)]_-^2 + [\bar{m}_T^\ell(\theta, g)]_+^2 \}, \quad \text{with} \quad (27)$$

$$\bar{m}_T^u(\theta, g) := (T\bar{J}_T)^{-1} \sum_{t=1}^T \sum_{j=1}^{J_t} (\hat{\delta}_{jt}^u(s_t, \lambda) - x'_{jt}\beta)g(z_{jt}) \text{ and}$$

$$\bar{m}_T^\ell(\theta, g) := (T\bar{J}_T)^{-1} \sum_{t=1}^T \sum_{j=1}^{J_t} (\hat{\delta}_{jt}^\ell(s_t, \lambda) - x'_{jt}\beta)g(z_{jt}).$$

where  $\mu(g) : \mathcal{G} \rightarrow [0, 1]$  is a probability mass function on  $\mathcal{G}$  and  $\mathcal{G}$  is a collection of instrumental functions. Both  $\mathcal{G}$  and  $\mu(\cdot)$  have been given in Section 3.4.

## 5 Point Identification Condition

The point identification condition is motivated by the power law feature of the data demonstrated in Section 2. The feature indicates a coexistence of a few dominant products with thick demand and a large number of fringe products with thin demand. For a dominant product  $j$  in market  $t$ ,  $\pi_{jt}$  is large and  $\log(s_{jt} + \iota_u/n_t)$  and  $\log(s_{jt} + \iota_\ell/n_t)$  are close to each other and close to  $\log(\pi_{jt})$  at large  $n_t$ . If certain values of the exogenous variables  $z_{jt}$  predict such  $\pi_{jt}$ 's, then at those  $z_{jt}$  values, the conditional moment inequalities in (25) hold as equalities asymptotically. These equalities may yield point identification by standard identification arguments for BLP moment conditions.

Formally, let  $\mathcal{Z}_0$  stand for the set of values of  $z_{jt}$  that predict dominant products (those with choice probabilities that do not approach zero). We can state the assumption as follows:

**Assumption 3.** *There exists a fixed constant  $\varepsilon_0 \in (0, 1)$  and a set  $\mathcal{Z}_0 \subseteq \text{supp}(z_{jt})$  such that  $\inf_{j,t,T} \Pr(z_{jt} \in \mathcal{Z}_0) > 0$ , such that  $\Pr(\pi_{jt} \geq \varepsilon_0 | z_{jt} \in \mathcal{Z}_0) = 1$  for all  $j, t$ .*

Below we give three stylized demand-supply models that could give rise to the dominant products below and discuss what the dominant product predictors are in each case. For now, it is important to note that  $z_{jt}$  includes both the exogenous covariates in the demand model and excluded instruments (if there are any). Often in practice, it is brand or UPC dummies that predict dominant status, which usually are also included exogenous covariates.

We state the lemma that shows that the bounds collapse on  $\mathcal{Z}_0$  under Assumption 3:

**Lemma 2.** *Suppose that  $\min_{t=1,\dots,T} n_t^2/T \rightarrow \infty$  as  $T \rightarrow \infty$ , and that Assumption 3 holds. Then, we have*

$$\sup_{j=1,\dots,J_t,t=1,\dots,T} \sup_{\lambda \in \Lambda} n_t |\hat{\delta}_{jt}^u(s_t, \lambda) - \hat{\delta}_{jt}^\ell(s_t, \lambda)| 1\{z_{jt} \in \mathcal{Z}_0\} = O_p(1). \quad (28)$$

*Remark.* When the bounds collapse, the moment inequalities (25) holds as equalities on  $\mathcal{Z}_0$  asymptotically. Then the standard (point) identification considerations for BLP models

apply here, except that attention is restricted on  $\mathcal{Z}_0$ . In general, if the instruments shift price and sales sufficiently for the dominant products, the model is point identified.

*Remark.* Note that neither  $\mathcal{Z}_0$  or  $\varepsilon_0$  need to be known in order to use our estimator. This is an advantage of the moment inequality approach comparing to an alternative approach that pre-selects products that never experience zeroes. The key to this is the Andrews and Shi (2013)-type instrumental functions that ensure that asymptotically, all the information in the conditional moment inequalities (25) are preserved in forming the unconditional moments. That will guarantee that the point-identification information provided by  $\mathcal{Z}_0$  is preserved as well, even though  $\mathcal{Z}_0$  is unknown.

Next, we discuss how the dominant products may come into being. Such products or firms have been a subject of interest since the early days of industrial organization. They have been studied under the name of *incumbents*, *leaders* as well as as dominant products/firms (see e.g. Markham (1951), Chapter 8 of Tirole (1988), Gowrisankaran and Holmes (2004), Shimomura and Thisse (2012)). They are the ones that enjoy a large market share and earn a positive profit despite that there are free entry and an unlimited number of potential entrants. The literature does not agree on how they achieve their dominant status. Simple explanations include (a) the dominant products are less substitutable with the fringe products than the fringe products among themselves, (b) the dominant products are much more appealing on average possibly due to brand loyalty or technological innovations, (c) the dominant products are provided with significantly lower cost possibly due to technology advances. In all explanations, a key is that the dominant products have features that are not easily replicable, so there is no free entry of products with those features. We illustrate each using a stylized example now.<sup>20</sup> In the examples, we ignore the  $t$  subscript for notational ease.

**Example 5.1.** Consider a nested logit model with three nests:  $\{0\}, \mathcal{J}_0, \mathcal{J}_1$ , where  $\mathcal{J}_0 \cup \mathcal{J}_1 = \{1, \dots, J\}$ . Let  $J_0$  and  $J_1$  denote the number of elements in  $\mathcal{J}_0$  and  $\mathcal{J}_1$ , respectively, and suppose that  $J_0$  is fixed as  $n$  grows but  $J_1$  grows proportionally to  $n$ , say  $J_1 = cn$ . Let  $\pi_{\mathcal{J}_\ell}$  stand for the probability that a product in  $\mathcal{J}_\ell$  is chosen, for  $\ell = 0, 1$ . Consider a nested logit model that yields

$$\begin{aligned} \frac{\pi_j}{\pi_{\mathcal{J}_\ell}} &= \frac{\exp(\delta_j)}{\sum_{j' \in \mathcal{J}_\ell} \exp(\delta_{j'})} \text{ for } j \in \mathcal{J}_\ell, \\ \pi_{\mathcal{J}_\ell} &= \frac{\exp(\lambda(\mathcal{I}(\mathcal{J}_\ell) - \log(J_\ell)))}{1 + \exp(\lambda(\mathcal{I}(\mathcal{J}_0) - \log(J_0))) + \exp(\lambda(\mathcal{I}(\mathcal{J}_1) - \log(J_1)))}, \text{ for } \ell = 0, 1, \end{aligned} \quad (29)$$

---

<sup>20</sup>As we can see, in each of the examples, the dominant status indicator is a discrete random variable. It is possible to conjure up a continuous dominant status indicator, but its support would need to have a discontinuity to separate the dominant and fringe products, a feature that could be difficult to justify in practice.

where  $\mathcal{I}(\mathcal{J}_\ell) = \log \left( \sum_{j \in \mathcal{J}_\ell} \exp(\delta_j) \right)$  and  $\lambda$  is a parameter. Suppose that  $\delta_j : j = 1, \dots, J$  are bounded between  $\underline{\delta}$  and  $\bar{\delta}$ . Then, it is easy to verify that  $\mathcal{I}(\mathcal{J}_\ell) - \log(J_\ell)$  is also bounded between  $\underline{\delta}$  and  $\bar{\delta}$ . Thus  $\pi_{\mathcal{J}_\ell} \in \left[ \frac{\exp(\lambda \underline{\delta})}{1 + \exp(\lambda \underline{\delta}) + \exp(\lambda \bar{\delta})}, \frac{\exp(\lambda \bar{\delta})}{1 + \exp(\lambda \underline{\delta}) + \exp(\lambda \bar{\delta})} \right]$ , and  $\frac{\pi_j}{\pi_{\mathcal{J}_\ell}} \in J_\ell^{-1} [\exp(\underline{\delta} - \bar{\delta}), \exp(\bar{\delta} - \underline{\delta})]$ . Then we have,

$$\begin{aligned} \pi_j &\geq J_0^{-1} \exp(\underline{\delta} - \bar{\delta} + \lambda \underline{\delta}) / (1 + \exp(\lambda \underline{\delta}) + \exp(\lambda \bar{\delta})) \text{ for } j \in \mathcal{J}_0, \\ n\pi_j &\geq c^{-1} \exp(\underline{\delta} - \bar{\delta} + \lambda \underline{\delta}) / (1 + \exp(\lambda \underline{\delta}) + \exp(\lambda \bar{\delta})) \text{ for } j \in \mathcal{J}_1. \end{aligned} \quad (30)$$

That is, products in nest  $\mathcal{J}_0$  are dominant products satisfying  $\pi_j \geq \underline{\varepsilon}_0$  and those in nest  $\mathcal{J}_1$  are fringe products satisfying  $n\pi_j \geq \underline{\varepsilon}_1$  for  $\underline{\varepsilon}_0 = cJ_0^{-1}\underline{\varepsilon}_1 = J_0^{-1} \exp(\underline{\delta} - \bar{\delta} + \lambda \underline{\delta}) / (1 + \exp(\lambda \underline{\delta}) + \exp(\lambda \bar{\delta}))$  for  $j \in \mathcal{J}_1$ . Assumption 3 is satisfied if  $1\{j \in \mathcal{J}_0\}$  is part of  $z_{jt}$ .

In this example, the number of fringe products is proportional to  $n$ . This appears arbitrary, but can in fact be a natural result of the zero-profit condition under free-entry into nest  $\mathcal{J}_1$  (see the discussion at the end of Section 3.1). The dominant products are dominant because they are protected from the competition of the fringe products by the substitution pattern in product demand and barrier to entry into nest  $\mathcal{J}_0$ .

**Example 5.2.** Consider a multinomial logit model for simplicity. Normalize  $\delta_{0t} = 0$ . Let

$$\delta_j = -\alpha p_j + \sum_{k=1}^J \beta_k UPC_{kj} + \xi_j, \quad (31)$$

where  $p_j$  is the price,  $UPC_{kj}$ 's are UPC dummies ( $UPC_{kj} = 1\{k = j\}$ ),  $\beta_j = b_j$  for  $j \in \mathcal{J}_0$ , and  $\beta_j = -\log(n) + b_j$  for  $j \notin \mathcal{J}_0$  for a subset  $\mathcal{J}_0$  of  $\{1, \dots, J\}$ , and  $b_j$  are bounded constants. Let  $J$  be fixed. Let  $p_j$  and  $\xi_j$  be bounded. Then for  $j \notin \mathcal{J}_0$ ,

$$\pi_j = \frac{\exp(-\alpha p_j + b_j + \xi_j) / n}{1 + \sum_{k=1}^J \exp(-\alpha p_k + \beta_k + \xi_k)} \geq n^{-1} \frac{\exp(-\alpha \bar{p} + \underline{b} + \underline{\xi})}{1 + J \exp(-\alpha \bar{p} + \bar{b} + \bar{\xi})}, \quad (32)$$

where  $\underline{p}, \underline{b}, \underline{\xi}$  are the lower bounds of  $p_j, b_j, \xi_j$ , and  $\bar{p}, \bar{b}, \bar{\xi}$  are the upper bounds. Let  $\underline{\varepsilon}_1 = \frac{\exp(\alpha \underline{p} + \underline{b} + \underline{\xi})}{1 + J \exp(\alpha \bar{p} + \bar{b} + \bar{\xi})}$ . Then this shows that  $n\pi_j \geq \underline{\varepsilon}_1$ . For  $j \in \mathcal{J}_0$ ,

$$\pi_j = \frac{\exp(-\alpha p_j + b_j + \xi_j)}{1 + \sum_{k=1}^J \exp(-\alpha p_k + \beta_k + \xi_k)} \geq \frac{\exp(-\alpha \bar{p} + \underline{b} + \underline{\xi})}{1 + J \exp(-\alpha \underline{p} + \bar{b} + \bar{\xi})}. \quad (33)$$

Let  $\underline{\varepsilon}_0 = \underline{\varepsilon}_1$ . Then this shows that Assumption 3 holds if the UPC dummies are used as part of  $z_{jt}$ .<sup>21</sup>

<sup>21</sup>Note that letting  $\alpha$  be a random coefficient or adding other covariates would not change the essence of the example.

In this example, the mean-utility of the fringe products depends on  $n$ . This can be a natural result of the zero-profit condition under free-entry (see the discussion at the end of Section 3.1): only fringe products with such mean-utilities self-select into the market.

**Example 5.3.** Consider a multinomial logit model again. Let  $\delta_j = -\alpha p_j + \xi_j$ . Let there be constant marginal cost  $c_j = \alpha^{-1} \log(n) z_j + c_{0j}$ , where  $z_j$  is a dummy variable and  $c_{0j}$  is a bounded constant. Suppose for simplicity that the products are supplied by single-product firms maximizing profit. Then it is easy to see that the optimal price is

$$p_j = c_j + \frac{1}{\alpha(1 - \pi_j)} \quad (34)$$

Let  $J$  be fixed and  $\xi_j$  be bounded. Then, for every  $j = 1, \dots, J$

$$\begin{aligned} \pi_j &= \frac{\exp(-\log(n)z_j - \alpha c_{0j} - (1 - \pi_j)^{-1})}{1 + \sum_{k=1}^J \exp(-\log(n)z_k - \alpha c_{0k} - (1 - \pi_k)^{-1})} \\ &\leq \exp(-\alpha \underline{c} - 1), \end{aligned} \quad (35)$$

where  $\underline{c}$  is the lower bound for  $c_{0j}$ . Let  $\bar{\pi} = \exp(-\alpha \underline{c} - 1)$ . Then for  $j$ 's with  $z_j=1$ , we have

$$\pi_j \geq n^{-1} \frac{\exp(-\alpha \bar{c} - (1 - \bar{\pi})^{-1})}{1 + \sum_{k=1}^J \exp(-\alpha \underline{c} - 1)}, \quad (36)$$

where  $\bar{c}$  is the upper bound for  $c_{0j}$ . Let  $\underline{\varepsilon}_1 = \frac{\exp(-\alpha \bar{c} - (1 - \bar{\pi})^{-1})}{1 + \sum_{k=1}^J \exp(-\alpha \underline{c} - 1)}$ , then this shows that  $n\pi_j \geq \underline{\varepsilon}_1$ . Similarly, we can show that for  $j$ 's with  $z_j = 0$ ,  $\pi_j \geq \underline{\varepsilon}_0 := \underline{\varepsilon}_1$ , verifying Assumption 3.

In this example, the cost of the fringe products depends on  $n$ . This can be a natural result of the zero-profit condition under free-entry (see the discussion at the end of Section 3.1): only fringe products with such costs self-select into the market.

As we see above, the point identification assumption is natural in many situations with dominant products. Nevertheless in settings where these Assumptions are questionable, we can still use (25) as a basis for partial identification and inference. For example, one can use the method developed in Andrews and Shi (2013) to construct a joint confidence set for the full vector  $\theta_0$ . This confidence set is constructed by inverting an Anderson-Rubin test:  $CS = \{\theta : T(\theta) \leq c(\theta)\}$  for some test statistic  $T(\theta)$  and critical value  $c(\theta)$ . Computing this set amounts to computing the 0-level set of the function  $T(\theta) - c(\theta)$ , where  $c(\theta)$  typically is simulated quantiles and thus a non-smooth function of  $\theta$ . A new approach that is computationally less burdensome when  $\beta$  is high dimensional is proposed in Gandhi, Lu, and Shi (2013), which also includes Monte Carlo simulations and empirical results using the profiling approach under partial identification. When the linear coefficients of the control

variables are nuisance parameters, one can also use the approach in Cox and Shi (2019) for inference to further reduce computational burden.

## 6 Consistency

In this section, we establish the consistency of the point estimator defined in (26). We need additional assumptions.

The first set of assumptions formalize the model and the data environment. They are similar to those in Berry, Linton, and Pakes (2004) and Freyberger (2015).

**Assumption 4.** (a) *The equation system (20) uniquely defines  $\delta_t(\pi_t, \lambda)$  for all  $t$ , all  $\pi_t \in \{\pi \in (0, 1)^{J_t} : \mathbf{1}'_{J_t} \pi < 1\}$  and all  $\lambda \in \Lambda$ .*

(b) *In each market, consumers' preferences  $(\epsilon_{ijt})_{j=1}^{J_t}$  are i.i.d. draws from the known distribution  $F(\cdot | x_t; \lambda_0)$  with unknown parameter  $\lambda_0 \in \Lambda$ . Consumer choice is determined by (19).*

(c) *The moment condition (22) holds.*

(d)  *$(x_t, s_t, z_t)_{t=1}^T$  are independent across markets.*

(e) *There exists a constant  $M$  such that  $E[\xi_{jt}^{2+c}] < M$  for all  $j = 1, \dots, J_t$ , all  $t = 1, \dots, T$ , and all  $T$  for some  $c > 0$ .*

(f)  *$\sup_{t=1, \dots, T} n_t \|\tilde{s}_t - s_t\|_f = O_p(1)$  as  $T \rightarrow \infty$ .*

(g)  *$\frac{\underline{n}_T}{J_T^{\max} \sqrt{T}} \rightarrow \infty$  and  $\frac{\log(\bar{n}_T)}{\sqrt{T}} \rightarrow 0$  where  $\underline{n}_T = \min_{t=1, \dots, T} n_t$ ,  $J_T^{\max} = \max_{t=1, \dots, T} J_t$ , and  $\bar{n}_T = \max_{t=1, \dots, T} n_t$ .*

*Remark.* Part (g) requires that  $n_t$  be not too small and not too big. The not-too-big part may be surprising because larger  $n_t$  is typically considered a good thing. Here larger  $n_t$  is not purely a good thing because we allow the lowest  $\pi_{jt}$  to be inversely related to  $n_t$ .<sup>22</sup> In this framework, larger  $n_t$  also implies lower minimum  $\pi_{jt}$  which *increases* the difficulty in bounding  $\log(\pi_{jt})$ . Also note that the  $J_T^{\max}$  in part (g) is not needed for multinomial logit and nested logit models for the reason discussed in the paragraph below Lemma 1.

The next assumption formalizes the lower bound for choice probabilities for the outside and the fringe products. These bounds have been discussed in detail in Sections 3 and 5.

**Assumption 5.** (a)  *$\pi_{0t} > \underline{\varepsilon}_0$  for all  $t$ .*

(b)  *$\pi_{jt} > \underline{\varepsilon}_1/n_t$  for all  $j, t$ .*

Next we impose a Lipschitz continuity assumption on  $\delta_{jt}(\pi, \lambda)$  in  $\pi$  on the part of the  $\pi$  space for the dominant products.

---

<sup>22</sup>Recall from Section 3 that this is done to rationalize the zeroes in the data.



**Assumption 6.**  $\sup_{t=1,\dots,T} \sup_{j=1,\dots,J_t} \sup_{\lambda \in \Lambda} \sup_{\pi, \bar{\pi} \in \Delta_{J_t}^{\varepsilon_0/2} : \pi \neq \bar{\pi}, \pi_j, \bar{\pi}_j \geq \varepsilon_0/2} \frac{|\delta_{jt}(\bar{\pi}, \lambda) - \delta_{jt}(\pi, \lambda)|}{\|\bar{\pi} - \pi\|_f \sqrt{J_t}} = O(1)$ .

*Remark.* Assumption 6 is a commonly accepted assumption when all products are dominant products (ref. Freyberger (2015)). The stronger version of this assumption without  $\sqrt{J_t}$  on the denominator holds for multinomial logit models:  $\delta_{jt}(\pi, \lambda) = \log(\pi_j) - \log(\pi_0)$  because the logarithm function is uniformly continuous on the interval  $[\varepsilon_0/2, \infty)$ . This argument combined with Assumption 1 (a) implies Assumption 6 for models satisfying Assumption 1. The same argument as that for the multinomial logit also works for the binary probit model.

Finally, we strengthen the point identification condition to ensure consistency. Define

$$\mathcal{G}_0 = \{g \in \mathcal{G} : g(z) = 0 \text{ for } z \notin \mathcal{Z}_0\}. \quad (37)$$

This is the set of instrumental functions that captures the identification information provided by the dominant products. Note that the dominant status predictor(s) in  $z_{jt}$  often is (are) brand/UPC dummy(-ies), thus, elements in  $\mathcal{G}_0$  are often those dummies interacted with dummies created for other elements of  $z_{jt}$  in the Andrews and Shi (2013) style (described in Section 3.4). It is also worth noting that one does not need to know  $\mathcal{G}_0$  but only need to know that  $\mathcal{G}$  contains such a  $\mathcal{G}_0$ , the latter guaranteed by Assumption 3 and the Andrews and Shi style  $\mathcal{G}$ .

Let

$$\begin{aligned} \bar{m}_T(\theta, g) &= \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\delta_{jt}(\pi_t, \lambda) - x'_{jt} \beta) g(z_{jt}) \\ \widehat{Q}_T^*(\theta) &= \sum_{g \in \mathcal{G}_0} \mu(g) \bar{m}_T(\theta, g)^2. \end{aligned} \quad (38)$$

The moments  $\bar{m}_T(\theta, g)$  is infeasible because  $\pi_{jt}$  is not observed. But we will be able to show that they are close to  $\bar{m}_T^u(\theta, g)$  and  $\bar{m}_T^l(\theta, g)$  for  $g \in \mathcal{G}_0$ . The criterion function  $\widehat{Q}_T^*(\theta)$  aggregates the infeasible moments for the dominant products. The assumption below is the additional identification condition:

**Assumption 7.** For any  $c > 0$ , there exists  $C(c) > 0$  such that

$$\lim_{T \rightarrow \infty} \Pr \left( \inf_{\theta \in \Theta : \|\theta^s - \theta_0^s\| > c} \widehat{Q}_T^*(\theta) > C(c) \right) = 1,$$

where  $\theta^s$  is a subvector of  $\theta$  and  $\theta_0^s$  is its true value.

*Remark.* This assumption ensures that the dominant products provide enough restriction to point identify the parameter  $\theta^s$ . Only a subvector of  $\theta$  is considered in this assumption

because we want to allow (but not require) specifications with product fixed effects. The fixed effects for the fringe products are clearly not identified since the data do not contain sufficiently precise information about their inverse demand. In that case  $\theta^s$  will only contain the common parameters and the fixed effects of the dominant products. Moreover, the assumption requires that the instrumental functions in  $\mathcal{G}_0$  are able to capture the variation of the moments over  $z_{jt} \in \mathcal{Z}_0$ . This in general requires that  $E[\hat{\delta}_{jt}^u(s_t, \lambda) - x_{jt}\beta | z_{jt} = z]$  and  $E[\hat{\delta}_{jt}^\ell(s_t, \lambda) - x_{jt}\beta | z_{jt} = z]$  are continuous in the continuous components of  $z$  and the projection of  $\mathcal{Z}_0$  onto the space of  $z_{c,jt}$  (the continuous components of  $z_{jt}$ ) is zero distance to an open set. This is innocuous in most applications.

Finally, Assumption 7 also requires that the instruments shift  $x_{jt}$  and  $\pi_{jt}$  sufficiently. This requirement is a standard one for BLP instruments. Thus, all the considerations for finding instruments in BLP models still apply.

The following theorem shows the consistency of the estimator defined in (26). Note that only the identified subvector  $\theta_0^s$  can be estimated consistently.

**Theorem 1.** *Suppose that either Assumption 1 holds and  $T^{-1} \sum_{t=1}^T J_t^2 / \bar{J}_T^2$  is bounded, or Assumption 2 holds and  $\sup_{t=1, \dots, T} J_t$  is bounded. Further suppose that Assumptions 3-7 hold. Then  $\|\hat{\theta}_T^s - \theta_0^s\| \rightarrow_p 0$ .*

*Remark.* Note that for logit-based models (which satisfy Assumption 1), we do not need  $J_t$  to be bounded. We require that the  $J_t$ 's are roughly even across  $t$ , which is formalized as the boundedness of  $T^{-1} \sum_{t=1}^T J_t / \bar{J}_T^2$ . For non-logit-based models satisfying Assumption 2, we require  $\sup_{t=1, \dots, T} J_t$  to be bounded because Assumption 2(c) requires  $\tilde{s}_{jt} = s_{jt} + 1/n_t$  which is incompatible with Assumption 4(f) unless  $\max_{t=1, \dots, T} J_t$  is bounded.

*Remark.* The proof of the theorem follows two steps. First we show that at the true value  $\theta_0$ ,  $\hat{Q}_T(\theta) = o_p(1)$ . Second, we show that for points in  $\Theta$  such that  $\theta^s$  is bounded away from  $\theta_0^s$ ,  $\hat{Q}_T(\theta)$  asymptotically dominate  $\hat{Q}_T^*(\theta)$  and the latter is bounded away from zero. The proof is given in Section C.1.

## 7 Inference

In this section we discuss statistical inference based on our point estimator. We show that the estimator is asymptotically normal despite that the bounds are slack for some  $g$ 's, which

is a similar result to that in Kahn and Tamer (2009) for censored regression models.<sup>23</sup>

Since the consistency is derived only for the subvector  $\theta^s$  of  $\theta$ , the asymptotically normality also will be about the subvector. For ease of notation, we consider the particular case where  $\theta^s = (\lambda', \beta^{s'})'$  where  $\beta^s$  is a subvector of  $\beta$ . The parameters in  $\beta$  excluded from  $\beta^s$  are the coefficients of variables that are zero for  $z_{jt} \in \mathcal{Z}_0$ .

More assumptions are needed. For clarity, we divide the assumptions into two groups, the first being standard ones similar to those in Freyberger (2015) and the second being the special assumptions that are needed to account for the presence and the unknown identity of the fringe products. Let  $B_c(\lambda_0) = \{\lambda \in \Lambda : \|\lambda - \lambda_0\| \leq c\}$  and  $B_c(\pi_t) = \{\tilde{\pi}_t \in (0, 1)^{J_t} : \mathbf{1}'\tilde{\pi}_t < 1, \|\pi_t - \tilde{\pi}_t\|_f \leq c\}$ . Let  $\mathcal{G} \setminus \mathcal{G}_0$  denote the relative complement of  $\mathcal{G}_0$  in  $\mathcal{G}$ . Let  $\partial m_{jt}(\lambda)$  denote  $\left( \partial \delta_{jt}(\pi_t, \lambda) / \partial \lambda' \quad x_{jt}^{s'} \right)'$ , where  $x_{jt}^s$  is the subvector of  $x_{jt}$  that correspond to  $\beta^s$ . Let

$$\Gamma_T(g) = (T \bar{J}_T)^{-1} \sum_{t=1}^T \sum_{j=1}^{J_t} E[\partial m_{jt}(\lambda_0) g(z_{jt})]. \quad (39)$$

**Assumption 8.** (a)  $\theta_0^s$  is in the interior of  $\Theta^s := \{\theta^s : \exists \theta^r \text{ s.t. } (\theta^{s'}, \theta^{r'})' \in \Theta\}$ .

(b) The function  $\delta_{jt}(\pi, \lambda)$  is twice-continuously differentiable on  $\Delta_{J_t}^0 \times \Lambda$ , for all  $j, t$ .

(c) For some  $c > 0$  and  $M < \infty$ ,

$$\begin{aligned} \sup_{j,t} E \left[ \sup_{\tilde{\pi}_t \in B_c(\pi_t)} \sup_{\lambda \in B_c(\lambda_0)} \left\| \frac{\partial \delta_{jt}(\tilde{\pi}_t, \lambda)}{\partial \lambda} \right\| \right] &\leq M, \\ \sup_{j,t} E \left[ \sup_{\lambda \in B_c(\lambda_0)} \left\| \frac{\partial^2 \delta_{jt}(\pi_t, \lambda)}{\partial \lambda \partial \lambda'} \right\| \right] &\leq M, \end{aligned}$$

and  $\sup_{j,t} E[\|x_{jt}^s\|^2 | z_{jt} \in \mathcal{Z}_0] \leq M$ .

(d)  $\lim_{T \rightarrow \infty} \sum_{g \in \mathcal{G}_0} \mu(g) \Gamma_T(g) \Gamma_T(g)' = \Upsilon$  for a matrix  $\Upsilon$  of full rank, and

$$\lim_{T \rightarrow \infty} T^{-1} \sum_{t=1}^T \sum_{g, g^* \in \mathcal{G}_0} Cov \left( \bar{J}_T^{-1} \sum_{j=1}^{J_t} \xi_{jt} g(z_{jt}), \bar{J}_T^{-1} \sum_{j=1}^{J_t} \xi_{jt} g^*(z_{jt}) \right) \Gamma_T(g) \Gamma_T(g)' \mu(g) \mu(g^*) = V.$$

(e)  $\lim_{T \rightarrow \infty} T^{-1} \bar{n}_T^{1/2} = 0$ .

---

<sup>23</sup>However, it is worth noting some subtle differences between the identification and inference arguments in this paper and those in Kahn and Tamer. In Kahn and Tamer, the upper and lower bounds collapse in *finite sample* for covariate values that indicate no-censoring, while in this paper, the upper and lower bounds collapse only asymptotically. Kahn and Tamer consider a fixed data-generating-process asymptotic framework, while the nature of our problem calls for a triangular array asymptotic framework. These are part of the reason that our conditions look more complicated than Kahn and Tamer's.

*Remark.* Parts (a)-(b) are standard regularity conditions for extreme estimators. Part (c) imposes a uniform bound on the derivatives of  $\delta_{jt}(\cdot, \lambda)$  with respect to  $\lambda$ . This bound condition is trivially satisfied for multinomial logit models and the binary probit models because  $\delta_{jt}(\cdot, \lambda)$  does not depend on  $\lambda$ . For the nested logit model,  $|\partial\delta_{jt}(\tilde{\pi}_t, \lambda)/\partial\lambda| = |\log(\tilde{\pi}_{gt}/\tilde{\pi}_{0t})| \leq 2|\log(\underline{\varepsilon}_0 - c)|$  as long as  $\pi_{gt}, \pi_{0t} > \underline{\varepsilon}_{0t}$  and  $c < \underline{\varepsilon}_0$ . And  $\partial^2\delta_{jt}(\pi_t, \lambda)/\partial\lambda\partial\lambda' = 0$ . Thus part (c) holds if the share of each nest is bounded from zero. For mixed logit models, one can verify part (c) following similar arguments as those for Lemma 9 in Supplemental Appendix D, under the additional assumptions that the covariates with random coefficients are bounded. Part (d) of the assumption is needed because we allow the data generating process to drift as  $T \rightarrow \infty$ . It regulates the limit of the drift in our asymptotic thought experiment. The only restriction it imposes on the data itself is that the Jacobian of the moment conditions has full-rank, which is standard for moment-based estimation and rules out perfect multicollinearity in  $x_{jt}^s$ .

**Assumption 9.** (a) *There exists a constant  $\eta > 0$  such that for all sufficiently small  $c > 0$  and all  $T$ , we have*

$$\begin{aligned} \sum_{g \in \mathcal{G} \setminus \mathcal{G}_0: (T\bar{J}_T)^{-1} \sum_{t=1}^T \sum_{j=1}^{J_t} E[(\log(s_{jt} + \iota_u/n_t) - \log(\pi_{jt}))g(z_{jt})] \leq c} \mu(g) &< c^\eta, \\ \sum_{g \in \mathcal{G} \setminus \mathcal{G}_0: (T\bar{J}_T)^{-1} \sum_{t=1}^T \sum_{j=1}^{J_t} E[(\log(s_{jt} + \bar{\iota}_\ell/n_t) - \log(\pi_{jt}))g(z_{jt})] \geq -c} \mu(g) &< c^\eta, \\ \sum_{g \in \mathcal{G} \setminus \mathcal{G}_0: (T\bar{J}_T)^{-1} \sum_{t=1}^T \sum_{j=1}^{J_t} E[g(z_{jt})(n_t s_{jt} + \iota_u)^{-1}] \leq c} \mu(g) &< c^\eta. \end{aligned}$$

(b) *Case 1: When Assumption 1 holds, assume that*

$$\sup_{j,t=1,\dots,T} E \left[ \left\| \frac{\partial \check{\delta}_{jt}(\pi_t, \lambda_0)}{\partial \pi} \right\|^2 \right] = O(J_T^{\max}) \quad \text{and} \quad \sup_{j,t=1,\dots,T} \sup_{\pi \in B_c(\pi_t)} \left\| \frac{\partial^2 \check{\delta}_{jt}(\pi, \lambda_0)}{\partial \pi \partial \pi'} \right\| = O_p(J_T^{\max})$$

for some  $c > 0$ .

*Case 2: When Assumption 2 holds, assume that*

$$\sup_{j,t=1,\dots,T} E \left[ \left\| \frac{\partial \delta_{jt}(\pi_t, \lambda_0)}{\partial \pi} \right\|^2 1(z_{jt} \in \mathcal{Z}_0) \right] = O(1), \quad \text{and}$$

$$\sup_{j,t=1,\dots,T} \sup_{\pi: \|\pi - \pi_t\| \leq c} \left\| \frac{\partial^2 \delta_{jt}(\pi, \lambda_0)}{\partial \pi \partial \pi'} 1(z_{jt} \in \mathcal{Z}_0) \right\| = O_p(1)$$

for some  $c > 0$ .

*Remark.* Part (a) of Assumption 9 is needed to show that the moments inequalities are slack enough for the fringe products to not interfere with the convergence rate and the asymptotic distribution of the bound estimator. It is satisfied for the  $\mu(\cdot)$  and  $\mathcal{G}$  that we propose if the exogenous variables that signal dominant products are discrete and  $\pi_{jt}$ 's for the fringe products converge to zero at the rate  $n_t^{-1}$  so that  $E[\log(s_{jt} + \underline{l}_u/n_t) - \log \pi_{jt} | z_{jt}]$  and  $E[\log(s_{jt} + \bar{l}_\ell/n_t) - \log \pi_{jt} | z_{jt}]$  are bounded away from zero. Part (b) strengthens the requirements of Assumptions 1 and 2 to ensure convergence rate of our estimator. The case 1 part of Assumption 9(b) implies the case 2 part of this assumption, thus is stronger. The weaker assumption is sufficient for case 2 because of the additional conditions in Assumption 2. Case 1 of Assumption 9(b) may be verified in a similar fashion as Assumption 1(a).<sup>24</sup>

**Theorem 2.** *Suppose that either Assumption 1 holds and  $T^{-1} \sum_{t=1}^T J_t^2 / \bar{J}_T^2$  is bounded, or Assumption 2 holds and  $\sup_{t=1, \dots, T} J_t$  is bounded. Further suppose that Assumptions 3-9 hold. Then we have*

$$\sqrt{T}(\hat{\theta}_T^s - \theta_0^s) \rightarrow_d N(0, \Upsilon^{-1} V \Upsilon^{-1}).$$

*Remark 1.* Note that  $\Upsilon$  and  $V$  depend on  $\mathcal{G}_0$  which in turn depends on the unknown set  $\mathcal{Z}_0$ . Thus, estimating the asymptotic variance covariance matrix can be difficult. Instead, following Kahn and Tamer (2009), we recommend using non-parametric bootstrap to obtain standard errors and confidence intervals. We follow this recommendation in the application in Section 9. We also evaluate the performance of bootstrap standard errors and bootstrap-based confidence intervals in our Monte Carlo experiments in Section 8.

*Remark 2.* The asymptotic variance formula also makes it clear that the choice of instrumental function set  $\mathcal{G}$  affects estimation accuracy. Potentially, one could choose  $\mathcal{G}$  to minimize the asymptotic variance, however, this does not seem to resemble the existing efficiency theory for conditional moment equalities, e.g. Chamberlain (1987), Newey (1990), and Ai and Chen (2003), mainly due to the structure that  $\mathcal{G}$  needs to take to preserve the information in the conditional moment inequalities. We thus leave this for future research.

## 8 Monte Carlo Simulations

In this section, we present three sets of Monte Carlo experiments with random coefficient logit models. The first experiment investigates the performance of our approach with moderate

---

<sup>24</sup>For multinomial logit and nested logit models part (b) is not needed. The proofs of Theorem 2 goes through with slight adjustment using the special structure of the inverse demand of such models, without using part (b). As a result, for such models the rate at which  $J_t$  increases with  $n_t$  does not need to be restricted.

fractions of zero shares, which should cover most of the empirical scenarios. In the second experiment, we test our estimator with a data generating process that produces extremely large fractions of zeros; the purpose is to further illustrate the key idea of our estimator in exploiting the long tail pattern that is naturally present in the data. In the third experiment, we use actual data from our application as the base for DGP; the purpose is to examine the performance of our estimator in a realistic setting and provide some practical guidelines regarding to the choice of instruments functions.

The first two experiments use the a random coefficient logit model, where the utility of consumer  $i$  for product  $j$  in market  $t$  is

$$u_{ijt} = \alpha_0 + x_{jt}\beta_0 + \lambda_0 x_{jt}v_i + \xi_{jt} + \epsilon_{ijt}, \quad (40)$$

where  $v_i \sim N(0, 1)$ ,  $\lambda_0$  is the standard deviation of the random coefficients on  $x_{jt}$ ,  $\epsilon_{ijt}$ 's are i.i.d. across  $i, j$  and  $t$  following Type I extreme value distribution. The parameters of interest are  $\beta_0$  and  $\lambda_0$ , while  $\alpha_0$  is a nuisance parameter. In both experiments, we set  $\lambda_0 = .5$ ,  $\beta_0 = 1$  and vary  $\alpha_0$  for different designs. We simulate  $T$  markets, each with  $J$  products. For the third experiment, we will describe the DGP in Section 8.3.<sup>25</sup>

## 8.1 Moderately Many Zeroes

In the first experiment, the observed and unobserved characteristics are generated as  $x_{jt} = \frac{j}{10} + N(0, 1)$  and  $\xi_{jt} \sim N(0, .1^2)$  for each product  $j$  in market  $t$ . Thus one feature of the design is that the  $x_{jt}$  has some persistence across markets - products with larger index tend to have higher value of  $x$  (which respects the nature of the variation in the scanner data shown in Section 2). Finally, the vector of empirical shares in market  $t$ ,  $(s_{0t}, s_{1t}, \dots, s_{Jt})$ , is generated from Multinomial  $\left(n, [\pi_{0t}, \pi_{1t}, \dots, \pi_{Jt}]'\right) / n$ , where  $n$  is the number of consumers in each market.<sup>26</sup>

With the simulated data set  $\{(s_{jt}, x_{jt}) : j = 1, \dots, J\}_{t=1}^T$ , we compute our bound esti-

<sup>25</sup>In all the three experiments, we checked the realized  $\min_{j,t} n_t \pi_{jt}$ 's in the generated data and they are all well-above our choice of  $\iota_\ell = 2^{-52} = 2.2204^{-16}$ , which approximately equals to  $\underline{\epsilon}_1$  according to our calculations shown in Table 2. Hence, the key assumption that our moment inequalities build on,  $\min_{j,t} n_t \pi_{jt} \geq \underline{\epsilon}_1$ , is satisfied easily.

<sup>26</sup>The  $\pi_t$  has no closed form solution in the random coefficient model, and thus, we compute them via simulation, i.e.,

$$\pi_{jt} = \frac{1}{s} \sum_{i=1}^s \frac{\exp(\alpha_0 + x_{jt}\beta_0 + \lambda_0 x_{jt}v_i + \xi_{jt})}{1 + \sum_{k=1}^J \exp(\alpha_0 + x_{kt}\beta_0 + \lambda_0 x_{kt}v_i + \xi_{kt})},$$

where  $s = 1000$  is the number of consumer type draws ( $v_i$ ).

imator,<sup>27</sup> the standard BLP estimator using empirical share  $s_t$  in place of  $\pi_t$  and discarding observations with  $s_{jt} = 0$ , the standard BLP estimator using Laplace shares  $s_t^L = (n_t s_t + 1)/(n_t + J_t + 1)$  in place of  $\pi_t$ .

All the estimators require simulating the market shares and solving demand systems for each trial of  $\lambda$  in optimizing the objective function for estimation. We use the same set of random draws of  $v_i$  in estimation as in the data generating process to eliminate simulation error as simulation error is not the focus of this paper. BLP contraction mapping method is employed to numerically solve the demand systems for all three estimators.

We simulate 1000 data sets  $\{(s_t^r, x_t^r) : t = 1, \dots, T\}_{r=1}^{1000}$  and implement all the estimators mentioned above on each for a repeated simulation study. For the instrumental functions, we use the countable hyper-cubes defined in (16), and set  $\bar{r}_T = 50$ . The choices of  $\iota_\ell$  and  $\iota_u$  follow Section ???. For the BLP estimator, we use  $(1, x_{jt}, x_{jt}^2 - 1, x_{jt}^3 - 3x_{jt})$  (the first three Hermite polynomials) as instruments to construct the GMM objective function. Alternative transformations of  $x_{jt}$  as instruments yield effectively the same results.

The bias and standard deviation of the estimators are presented in Table 3. As we can see from the table, The standard BLP estimator with using empirical share  $s_t$  (labeled as “ES”) shows large bias for both  $\beta$  and  $\lambda$ . Replacing the empirical share  $s_t$  with the Laplace share  $s_t^L$  (and thus not discarding the observations with  $s_{jt} = 0$ ), labeled as “LS”, increases the bias for  $\beta$  although reducing the bias for  $\lambda$ . Our bound estimator (labeled as “Bound”) is the least biased, and its bias is very small for both parameters, especially when the sample size ( $T$ ) is large.

Next, we examine the performance of our proposed bootstrap procedure and the results are reported in Table 4. We can see that bootstrap standard errors are on average slightly larger than the standard deviation of the estimators, especially for the cases with large fraction of zeros and small sample size. Also, we compute two versions of bootstrap confidence intervals and find that the “Normal CI”, based on normal quantile and bootstrap standard errors, outperforms the standard nonparametric percentile bootstrap confidence interval and gets rather close to the nominal level (95%) of coverage probability, especially for the  $\beta$ , when the sample size gets large and the fraction of zeros is not too high.

## 8.2 Extremely Many Zeroes

Next we pressure test our bound estimator by pushing the fraction of zeroes in empirical shares toward the extreme. We modify the DGP slightly to produce very high fraction of zeros. Specifically, we generate  $x_{jt}$  from the following discrete distribution

---

<sup>27</sup>We use  $\tilde{s}_{jt} = s_{jt} + 1/(n_t J_t)$  for  $j = 1, \dots, J_t$  when implementing the bound estimator for all the simulations in this section.

Table 3: Monte Carlo Results: Estimation

DGP	$T$	Ave. % of Zeros	ES		LS		Bound		
			$\lambda$	$\beta$	$\lambda$	$\beta$	$\lambda$	$\beta$	
I	25	9.52%	Bias	.3718	-.1941	.2900	-.2167	.0422	-.0432
			SD	.0337	.0160	.0221	.0115	.0477	.0352
	50	9.48%	Bias	.3712	-.1939	.2912	-.2172	.0172	-.0216
			SD	.0236	.0118	.0164	.0082	.0388	.0284
	100	9.46%	Bias	.3714	-.1941	.2900	-.2169	.0002	-.0065
			SD	.0169	.0081	.0112	.0055	.0311	.0234
II	25	18.54%	Bias	.6752	-.6115	.4023	-.4675	.0142	-.0302
			SD	.0845	.0655	.0315	.0229	.0531	.0536
	50	18.54%	Bias	.6649	-.6040	.3993	-.4657	-.0083	-.0028
			SD	.0580	.0462	.0223	.0158	.0410	.0413
	100	18.50%	Bias	.6624	-.6021	.3983	-.4651	-.0154	.0073
			SD	.0422	.0333	.0163	.0114	.0297	.0297
III	25	41.13%	Bias	.7302	-1.3220	.3868	-.9863	-.0366	.0278
			SD	.2022	.2890	.0366	.0460	.0481	.0721
	50	41.09%	Bias	.7092	-1.2947	.3830	-.9819	-.0331	.0303
			SD	.1373	.1975	.0262	.0323	.0374	.0549
	100	41.09%	Bias	.7070	-1.2935	.3809	-.9794	-.0219	.0176
			SD	.0911	.1325	.0188	.0232	.0282	.0391
IV	25	52.39%	Bias	.4013	-1.1035	.2907	-1.1412	-.0499	.0512
			SD	.1346	.2435	.0304	.0453	.0530	.0899
	50	52.35%	Bias	.3942	-1.0937	.2877	-1.1369	-.0346	.0330
			SD	.0956	.1740	.0214	.0313	.0396	.0635
	100	52.36%	Bias	.3916	-1.0901	.2862	-1.1349	-.0215	.0169
			SD	.0687	.1255	.0154	.0227	.0311	.0475

Note: 1.  $J = 50$ ,  $n = 10,000$ ,  $\beta_0 = 1$ ,  $\lambda_0 = .5$ , Number of Repetitions = 1000.

2. “ES”: Empirical Shares; “LS”: Laplace Shares.

3. DGP: I, II, III and IV correspond to  $\alpha_0 = -9, -10, -12$  and  $-13$ , respectively.

$x$	1	12	15
$\Pr(x_{jt} = x)$	.99	.005	.005

and

$$\xi_{jt} \sim 1(x_{jt} = 1) \times N(0, 2^2) + 1(x_{jt} \neq 1) \times N(0, .1^2).$$

All the other aspects of the DGP is identical to the previous simulation.

The fractions of zeroes are made very high: 82%-96% by choosing the  $\alpha_0$  parameter. With such high fractions of zeroes, the vast majority of observations are uninformative. Thus, we need larger sample size for any estimator to perform well. We consider  $T = 100, 200, 400$ . For simplicity of presentation and to reduce computational burden, here we fix  $\lambda$  at its true



Table 4: Monte Carlo Results: Bootstrap

DGP	$T$	Ave. % of Zeros	Actual SD		BS SE		CP: BS CI		CP: Normal CI	
			$\lambda$	$\beta$	$\lambda$	$\beta$	$\lambda$	$\beta$	$\lambda$	$\beta$
I	25	9.52%	.0477	.0352	.0473	.0353	.8390	.8250	.8630	.7790
	50	9.48%	.0388	.0284	.0400	.0300	.8556	.8675	.9444	.9160
	100	9.46%	.0311	.0234	.0324	.0244	.8408	.8458	.9570	.9530
II	25	18.54%	.0531	.0536	.0563	.0585	.8390	.8630	.9640	.9510
	50	18.54%	.0410	.0413	.0423	.0433	.7980	.8340	.9490	.9690
	100	18.50%	.0297	.0297	.0309	.0311	.8380	.8750	.9270	.9560
III	25	41.13%	.0481	.0721	.0537	.0840	.7700	.8310	.9040	.9680
	50	41.09%	.0374	.0549	.0388	.0581	.8360	.8760	.8690	.9360
	100	41.09%	.0282	.0391	.0290	.0417	.8740	.9250	.9000	.9450
IV	25	52.39%	.0530	.0899	.0549	.0971	.7880	.8550	.8710	.9430
	50	52.35%	.0396	.0635	.0420	.0707	.8490	.9120	.8870	.9530
	100	52.36%	.0311	.0475	.0312	.0498	.8450	.8980	.9040	.9440

Note: 1. All the settings are identical to Table 1. Bootstrap draws are taken at market level.

Bootstrap sample size is 500.

2. “BS SE” refers to average bootstrap standard error.
3. “CP: BS CI” refers to the coverage probability of the 95% nonparametric bootstrap CI.
4. “CP: Normal CI” refers to the coverage probability of the 95% normal CI with bootstrap s.e.

value, and only investigate the behavior of the estimators for  $\beta$ .

The results are reported in *Table 5*, and they are very encouraging for the bound approach. The ES estimator is severely biased toward 0, so is the LS estimator. The bound estimator is remarkably accurate in these extreme cases. The performance highlights the key idea behind our estimator: utilizing the information from the dominant products with inherently thick demand while controlling the impact of fringe products with small/zero sales on estimation.

### 8.3 Monte Carlo Simulations with Tuna Data

In this subsection, we conduct Monte Carlos simulations based on the canned tuna data set that will be used later in our application. The main purposes are two fold: 1) we want to examine performance of the bound estimator in a setting that is similar to the application; 2) we would like to understand better how the choices of instruments affect the performance of the bound estimator, especially in real empirical settings where product dummies are typically included.

To generate data, we use tuna data in one week (the week of March 30, 1995) across all the stores (there are 80 stores) as a template (a store-week as a “market” and a UPC as a “product”) and consider a random coefficient logit specification that extends (40). In

Table 5: Monte Carlo Results: Very Large Fraction of Zeros

DGP	$T$	Ave. % of Zeros		ES	LS	Bound
I	100	84.73%	Bias	-.2698	-.2643	-.0014
			SD	.0060	.0058	.0123
	200	84.68%	Bias	-.2695	-.2640	-.0011
			SD	.0042	.0040	.0094
	400	84.71%	Bias	-.2692	-.2639	-.0005
			SD	.0030	.0030	.0066
II	100	91.45%	Bias	-.3328	-.3319	-.0016
			SD	.0066	.0061	.0126
	200	91.43%	Bias	-.3324	-.3314	-.0014
			SD	.0049	.0044	.0091
	400	91.43%	Bias	-.3320	-.3313	-.0007
			SD	.0034	.0032	.0067
III	100	95.37%	Bias	-.3992	-.4028	-.0014
			SD	.0079	.0070	.0126
	200	95.35%	Bias	-.3991	-.4025	-.0014
			SD	.0056	.0049	.0093
	400	95.36%	Bias	-.3986	-.4023	-.0010
			SD	.0040	.0035	.0065

Note: 1.  $T = 100$ ,  $J = 50$ ,  $n = 10,000$ ,  $\beta_0 = 1$ ,  $\lambda_0 = .5$ .

Number of Repetitions = 1000.

2. We fix  $\lambda = \lambda_0$  (at the true value) without estimating it.

3. DGP: I, II, III correspond to  $\alpha_0 = -13, -14, -15$ .

particular, we let the price coefficient be random, i.e.,

$$u_{ijt} = a_0 + \beta_0 x_{jt} - v_i p_{jt} + \xi_{jt} + \epsilon_{ijt},$$

where  $v_i$  follows Lognormal  $(\mu_p, \sigma_p)$ . The product-market specific demand shock  $\xi_{jt}$  has a simple heteroskedasticity structure

$$\xi_{jt} = 1(\beta_0 x_{jt} \geq \text{med}(\beta_0 x_{jt})) \xi'_{jt} + 1(\beta_0 x_{jt} < \text{med}(\beta_0 x_{jt})) \xi''_{jt},$$

where  $\xi'_{jt}$  ( $\xi''_{jt}$ ) follows normal distribution  $N(0, .5^2)$  ( $N(0, 1.5^2)$ ) truncated at  $\pm 3\sigma$ . The truncation gives  $\xi_{jt}$  a finite support to ensure that Assumptions 3 and 5 hold easily. Price is generated as a linear combination of marginal cost (use the observed wholesale price from the data) and a markup term that is a function of demand shock  $\xi$ , i.e.,

$$p_{jt} = mc_{jt} + b_0 \exp(\xi_{jt}). \quad (41)$$

Note that the markup term introduces a simple endogeneity problem. The covariates  $x_{jt}$  include a continuous variable following  $N(0, 1)$  (truncated at  $\pm 3\sigma$ ) and UPC dummies from the data. The coefficient on the continuous variable is 1 and those on the UPC dummies are set to be the estimated ones (using bound estimator) in our application. Other specifications are similar to the previous DGP. And the number of consumers in each market for generating market shares is directly imported from the data.

We simulate 1000 data sets that have the same structure as the real data, with the endogenous variables, i.e., price and market shares, varying across data sets. Then we implement several estimators of interests using the data sets. To simplify the estimation, we only estimate the two parameters of the random coefficient on price and fixing other parameters (UPC fixed effects) at their true values without estimating them.

The estimation results are summarized in *Table 6*. Note that we consider three values of  $a_0$  that imply different fractions of zeros (labeled by “I”, “II” and “III”). Also, besides the baseline  $T = 80$  case with one week data (the week of March 30, 1995), we also try  $T = 160$  using two weeks’ data (the weeks of March 23 and March 30, 1995). As before, “ES” and “LS” refer to standard BLP estimator applied to empirical shares and Laplace shares, respectively. For the bound estimator, we experiment with four alternative sets of instrument functions. “Bound- $\mathcal{G}_1$ ” uses the instruments defined by (16), which includes indicators constructed from continuous variables ( $z_{jt}, mc_{jt}$ ) with  $\bar{r}_{80} = 10$  and  $\bar{r}_{160} = 15$  and UPC dummies. “Bound- $\mathcal{G}_2$ ” is the same as “Bound- $\mathcal{G}_1$ ” except with larger  $\bar{r}_T$ :  $\bar{r}_{80} = 20$  and  $\bar{r}_{160} = 30$ . “Bound- $\mathcal{G}_3$ ” (“Bound- $\mathcal{G}_4$ ”) expands the set of instruments of “Bound- $\mathcal{G}_1$ ” (“Bound- $\mathcal{G}_2$ ”) by including the interactions between indicators constructed from continuous variables (denoted by  $\mathcal{C}$  in (16)) and UPC dummies.

From the results, we can see that:

- In almost all the cases, as before, the bound estimators have much smaller biases than the ES and LS estimators do (although with slightly increased standard deviations), especially for the standard deviation of the random coefficient  $\sigma_p$ .
- By comparing Bound- $\mathcal{G}_1$  and Bound- $\mathcal{G}_2$ , we can see that increasing the tuning parameter  $\bar{r}_T$  reduces standard deviation substantially but increase biases slightly.
- Including interactions between  $\mathcal{C}$  and UPC dummies reduces standard deviations of the estimators substantially. Hence, it seems preferable to have a sufficiently large  $\bar{r}_T$  and include the interactions, and these findings guide the construction of  $\mathcal{G}$  in our empirical application.

Table 6: Monte Carlo Results: Simulation using Tuna Data

DGP	$T$	Ave. % of Zeros	Panel I: $\mu_p$						
				ES	LS	Bound- $\mathcal{G}_1$	$\mathcal{G}_2$	$\mathcal{G}_3$	$\mathcal{G}_4$
I	80	9.11%	Bias	-.0001	-.0605	-.0009	-.0007	-.0004	-.0014
			SD	.0228	.0251	.0326	.0239	.0299	.0233
	160	9.09%	Bias	-.0081	-.0678	-.0026	-.0042	-.0019	-.0037
			SD	.0158	.0185	.0296	.0175	.0276	.0173
II	80	14.29%	Bias	-.0166	-.1403	-.0098	-.0104	-.0096	-.0099
			SD	.0258	.0311	.0582	.0303	.0498	.0305
	160	14.25%	Bias	-.0246	-.1472	-.0005	-.0129	-.0037	-.0126
			SD	.0186	.0228	.0721	.0221	.0530	.0220
III	80	17.70%	Bias	-.0286	-.2205	-.0097	-.0185	-.0123	-.0180
			SD	.0269	.0318	.0767	.0344	.0635	.0344
	160	17.66%	Bias	-.0373	-.2267	.0203	-.0185	-.0013	-.0183
			SD	.0192	.0234	.1178	.0253	.0706	.0252
Panel II: $\sigma_p$									
I	80	9.11%	Bias	.2691	.4111	.0555	.0840	.0630	.0840
			SD	.0664	.0848	.1382	.0846	.1176	.0804
	160	9.09%	Bias	.2457	.3931	.0065	.0402	.0208	.0414
			SD	.0484	.0647	.1570	.0662	.1338	.0657
II	80	14.29%	Bias	.3103	.5125	.0163	.0656	.0367	.0664
			SD	.0674	.1030	.2469	.0978	.2119	.0990
	160	14.25%	Bias	.2924	.5022	.0007	.0290	.0324	.0304
			SD	.0511	.0778	.3153	.0747	.2633	.0747
III	80	17.70%	Bias	.3410	.5745	.0340	.0578	.0458	.0590
			SD	.0710	.1164	.3273	.1163	.2873	.1179
	160	17.66%	Bias	.3221	.5678	.1010	.0260	.1158	.0275
			SD	.0517	.0871	.3819	.0911	.2915	.0912

<sup>1</sup> DGP: I, II and III correspond to  $a_0 = .4, .8, 1$ , respectively.

Number of markets  $T$ : 80 and 160 correspond to one week (03/30/1995 to 04/05/1995) and two weeks (03/23/1995 to 04/05/1995) of the tuna data for all the stores, respectively.

<sup>2</sup> “E”: Empirical Shares; “LS”: Laplace Shares; “Bound- $\mathcal{G}_1$ ”:  $\bar{r}_{80} = 10, \bar{r}_{160} = 15$ ; “Bound- $\mathcal{G}_2$ ”:  $\bar{r}_{80} = 20, \bar{r}_{160} = 30$ ; “Bound- $\mathcal{G}_3$ ”:  $\bar{r}_{80} = 10, \bar{r}_{160} = 15$ , instruments in  $\mathcal{C}$  interact with product dummies; “Bound- $\mathcal{G}_4$ ”:  $\bar{r}_{80} = 20, \bar{r}_{160} = 30$ , instruments in  $\mathcal{C}$  interact with product dummies.

<sup>3</sup> True value:  $\mu_p = 0, \sigma_p = .5$ . Coefficients on product dummies are fixed at their true values without being estimated for ease of computation. Number of repetitions = 1000.

## 9 Empirical Application

In this section, we apply our estimator on the DFF scanner data previewed in Section 2.<sup>28</sup>

In particular, we focus on the canned tuna category, as previously studied by Chevalier, Kashyap, and Rossi (2003) (CKR for short) and Nevo and Hatzitaskos (2006) (NH for

<sup>28</sup>The sample period predates the price fixing conduct by the tuna cartel starting around 2011, see Miller, Remer, and Weinberg (2020) for details.

short). CKR observed using the same data discussed in Section 2 that the share-weighted average price (i.e. the price index) of tuna fell by 15 percent during Lent – a high demand period for this product. They attributed the outcome to loss-leading behavior on the part of retailers. NH on the other hand suggest that this pricing pattern in the tuna data could instead be explained by increased price sensitivity of consumers (consistent with an increase in search) which causes a re-allocation of market shares towards less expensive products in the Lent period, and hence a fall in the observed share weighted price index. They test this hypothesis directly in the data by estimating demand parameters separately in the Lent and Non-Lent periods, and find that demand becomes more elastic in the high demand (Lent) period.

Here we revisit the groundwork laid by NH to examine the difference in price elasticity between Lent and non-Lent periods. The main difference in our analysis is that we use data on all products in the analysis, while NH restrict the sample to include only the top 30 UPCs and thus automatically drop products with small/zero sales. There are two main questions we seek to address: (a) Does the selection of UPC’s with only positive shares significantly bias the estimates of price elasticity and (b) Does the difference in price elasticities between the Lent and Non-Lent period persist after properly controlling for zeroes.

To make the comparison clear, we use largely the same specification of the model used in NH. In particular we consider a logit specification

$$u_{ijt} = \alpha p_{jt} + \beta x_{jt} + \xi_{jt} + \epsilon_{ijt},$$

where the control variables  $x_{jt}$  consist of UPC fixed effects and a time trend.<sup>29</sup> The week to week variation in the product-/market-level unobserved demand shock  $\xi_{jt}$  largely captures the short-term promotional efforts, e.g., in-store advertising and shelving choices, because the UPC fixed effects control the intrinsic product quality which is likely stable over short time horizon. Since stores are likely to advertise or shelf the product in a more prominent way during weeks when the product is on a price sale, we expect a negative correlation between price and the unobservable. We construct instruments for price by inverting DFF’s data on gross margin to calculate the chain’s wholesale costs, which is the standard price

---

<sup>29</sup>Empirical market shares are constructed using quantity sales and the number of people who visited the store that week (the customer count) as the relevant market size.

instrument in the literature that has studied the DFF data.<sup>30</sup>

We implement our bound estimator defined by (26) to obtain point estimate of  $(\alpha, \beta)$  in the model.<sup>31</sup> The standard errors are obtained using nonparametric bootstrap.<sup>32</sup> The estimation results are presented in Tables 7 and 8.<sup>33</sup> Table 7 shows that standard BLP logit estimator that inverts empirical shares to recover mean utilities (and hence drops zeroes) has a significant selection bias towards zero. The UPC level elasticities for the logit model are small in economic magnitude, with the average elasticity in the data being -.572. Furthermore, over 90% of products have inelastic demand. Using our bounds approach instead to control for zeroes has a major effect on the estimated elasticities. Average demand elasticity for UPC's becomes -1.51 and less than 30% percent of observations have inelastic demand. This change in the direction of elasticities is consistent with the attenuation bias effects of dropping products with small/zero market shares.

Table 7: Demand Estimation Results

	BLP	Bound
Price Coefficient	-.39	-1.03
S.E.	(.005)	(.319)
Ave. Own Price Elasticity	-.57	-1.51
Fraction of Inelastic Products	90.04%	28.20%
No. of Obs.	862,683	959,331

Note: The S.E. for the bound approach is the bootstrap standard error (using 1000 bootstrap replications).

<sup>30</sup>The gross margin is defined as (retail price - wholesale cost)/retail price, so we get wholesale cost using retail price  $\times$  (1 - gross margin). The instrument is defensible in the store disaggregated context we consider here because it has been shown that price sales in retail price primarily reflect a reduction in retailer margins rather than a reduction in wholesale costs (see e.g., Chevalier, Kashyap, and Rossi (2003) and Hosken and Reiffen (2004)); thus sales (and hence promotions) are not being driven by the manufacturer through temporary reduction in wholesale costs. However, this instrument may be invalid if manufacturers respond to demand shocks and adjust wholesale prices accordingly. We acknowledge the potential deficiency of using this instrument but searching for a better alternative is beyond the scope of the current paper.

<sup>31</sup>The choice of  $\mathcal{G}$  is guided by the simulation results in Section 8.3: we set  $\bar{r} = 45$  when constructing instrument functions from the wholesale cost (continuous variable) and include interactions between them and the UPC dummies.

<sup>32</sup>The procedure contains the following steps: (1) draw with replacement a bootstrap sample of *markets*, denoted as  $\{t_1, \dots, t_T\}$ ; (2) compute the bound estimator  $\hat{\theta}_T^{BD*}$  using the bootstrap sample; (3) repeat (1)-(2) for  $B_T$  times and obtain  $B_T$  independent (conditional on the original sample) copies of  $\hat{\theta}_T^{BD*}$ ; (4) obtain the sample standard deviation of the  $B_T$  copies of  $\hat{\theta}_T^{BD*}$  and this is the bootstrap standard error.

<sup>33</sup>In principle we can estimate our model separately for each store, letting preferences change freely over stores depending on local preferences. These results are available upon request. Here we present for the results of demand pooling together all stores together as was done by Nevo and Hatzitaskos (2006). The store level regressions results are very similar to the pooled store regression and the latter is a more concise summary of demand behavior that we present here.

Table 8: Demand in Lent vs. Non-Lent

	BLP		Bound	
	Lent	Non-Lent	Lent	Non-Lent
Price Coefficient	-.518	-.371	-1.23	-.75
S.E.	(.018)	(.005)	(.221)	(.231)
Ave. Own Price Elasticity	-.757	-.544	-1.80	-1.10
Fraction of Inelastic Products	84.02%	92.84%	16.79%	43.94%
No. of Obs.	70,496	792,187	78,838	880,493

Note: The S.E. for the bound approach is the bootstrap standard error (using 1000 bootstrap replications).

Our second result is that demand becomes more elastic in the high demand period, as shown in Table 8. This is consistent with Nevo and Hatzitaskos (2006)’s findings that are based on the standard logit estimator with zeroes being dropped. However, the Lent effect is bigger according to our bounds estimator that controls for the zeroes. In other words, correcting the selection bias, our bound estimator brings the price coefficient and elasticity higher and the correction effect is higher for the Lent period than for the non-Lent period. Since the fractions of zeroes are remarkably close between Lent and non-Lent periods, we suspect that the difference in the correction effect is due to a difference in the distribution of the unobservable  $\xi$ .

To further investigate this, we first replicate the reduced form finding of Nevo and Hatzitaskos (2006) that suggested a change in price sensitivity in the Lent period. This is reported in Table 9, which shows that although the price index of tuna during Lent appears to be approximately 15 percent less expensive than other weeks (as previously underscored by CKR), the average price of tuna is virtually unchanged between the Lent versus non-Lent period. Hence it is a re-allocation of demand towards less expensive products during Lent that drives the change in the price index.

Table 9: Regression of Price Index on Lent

	$P$ : Price Index	$P$ : Average Price
Lent	-.150	-.009
s.e.	(.0005)	(.0003)

We take this decomposition one step further than NH, and examine the price index separately for products “on sale” and “regularly priced” during these periods.<sup>34</sup> As can be seen in Table 10, it is the sales price index that is the key driver of the aggregate price index being cheaper during Lent. However the average price of an “on-sale” product is not cheaper

<sup>34</sup>We flag an observation in the data as being on sale if that particular UPC in that particular store in that particular week has at least a 5% reduction from highest price of previous 3 weeks.

in the Lent period. This shows that it is a re-allocation towards more steeply discounted “on-sale” product during Lent that is driving change in the aggregate price index. In contrast, we do not see an analogous reallocation for “regularly priced” products.

Table 10: Regression of Sales Price Index on Lent

	$P$ : Price Index		$P$ : Average Price	
	Sale	Regular	Sale	Regular
Lent	-.199	.035	.010	.001
s.e.	(.0017)	(.0003)	(.0016)	(.0003)

This suggests a tighter coordination of promotional effort and discounting in the high demand period. In effect more steeply discounted products are receiving larger promotional effort on the part of the retailer during the high demand, which is similar in spirit to the loss-leader hypothesis originally advanced for this data by CKR. Since promotional effort in the model is largely captured through the unobservable  $\xi$ , this change in behavior of the unobservable would account for the selection effect due to dropping zeroes changing across the two periods: during Lent period, the variance of promotional effort is larger so the selection bias is worse. Hence, our results suggest that both demand and supply side effects contribute to the falling price during high demand period, which complements and reconciles the findings of NH and CKR.

## 10 Conclusion

We have shown that differentiated product demand models have enough content to construct a system of moment inequalities that can be used to consistently estimate demand parameters despite a possibly large presence of observations with zero market shares in the data. We construct a GMM-type estimator based on these moment inequalities that is consistent and asymptotically normal under assumptions that are a reasonable approximation to the DGP in many product differentiated environments. Our application to scanner data reveals that taking the market zeroes in the data into account has economically important implications for price elasticities.

## References

AI, C., AND X. CHEN (2003): “Efficient Estimation of Models with Conditional Moment Restrictions Containing Unknown Functions,” *Econometrica*, 71, 1795–1843.



- ANDERSON, C. (2006): *The Long Tail: Why the Future of Business Is Selling Less of More*. Hyperion.
- ANDREWS, D. W. K., AND X. SHI (2013): “Inference Based on Conditional Moment Inequality Models,” *Econometrica*, 81.
- ARMSTRONG, T. B. (2016): “Large Market Asymptotics for Differentiated Product Demand Estimators With Economic Models of Supply,” *Econometrica*, 84, 1961–1980.
- BERRY, S. (1994): “Estimating discrete-choice models of product differentiation,” *The RAND Journal of Economics*, pp. 242–262.
- BERRY, S., A. GANDHI, AND P. HAILE (2013): “Connected substitutes and invertibility of demand,” *Econometrica*, 81(5), 2087–2111.
- BERRY, S., J. LEVINSOHN, AND A. PAKES (1995): “Automobile prices in market equilibrium,” *Econometrica: Journal of the Econometric Society*, pp. 841–890.
- BERRY, S., J. LEVINSOHN, AND A. PAKES (2004): “Differentiated Products Demand Systems from a Combination of Micro and Macro Data: The New Vehicle Market,” *Journal of Political Economy*, 112, 68–104.
- BERRY, S., O. LINTON, AND A. PAKES (2004): “Limit theorems for estimating the parameters of differentiated product demand systems,” *Review of Economic Studies*, 71(3), 613–654.
- BERRY, S., AND A. PAKES (2007): “THE PURE CHARACTERISTICS DEMAND MODEL,” *International Economic Review*, 48(4), 1193–1225.
- CHAMBERLAIN, G. (1987): “Asymptotic Efficiency in Estimation with Conditional Moment Restrictions,” *Journal of Econometrics*, 34, 305–334.
- CHEVALIER, J. A., A. K. KASHYAP, AND P. E. ROSSI (2003): “Why Don’t Prices Rise During Periods of Peak Demand? Evidence from Scanner Data,” *American Economic Review*, 93(1), 15–37.
- COX, G., AND X. SHI (2019): “Simple Adaptive Size-Exact Testing for Full-Vector and Subvector Inference in Moment Inequality Models,” unpublished manuscript, Department of Economics, University of Wisconsin at madison.

- DUBÉ, J.-P., J. T. FOX, AND C.-L. SU (2012): “Improving the numerical performance of static and dynamic aggregate discrete choice random coefficients demand estimation,” *Econometrica*, 80(5), 2231–2267.
- DUBE, J.-P., A. HORTACSU, AND J. JOO (2020): “Random-Coefficients Logit Demand Estimation with Zero-Valued Market Shares,” Becker Friedman Institute Working Paper NO. 2020-13.
- FREYBERGER, J. (2015): “Asymptotic theory for differentiated products demand models with many markets,” *Journal of Econometrics*, 185(1), 162–181.
- GABAIX, X. (1999): “Zipf’s Law and the Growth of Cities,” *The American Economic Review, Papers and Proceedings*, 89, 129–132.
- GANDHI, A., Z. LU, AND X. SHI (2013): “Estimating Demand for Differentiated Products with Error in Market Shares,” *CeMMAP working paper*.
- GOOLSBEE, A., AND A. PETRIN (2004): “The consumer gains from direct broadcast satellites and the competition with cable TV,” *Econometrica*, 72(2), 351–381.
- GOWRISANKARAN, G., AND T. J. HOLMES (2004): “Mergers and the Evolution of Industry Concentration: Result from the Dominant Firm Model,” 35, 561–582.
- HOSKEN, D., AND D. REIFFEN (2004): “Patterns of retail price variation,” *RAND Journal of Economics*, pp. 128–146.
- KAHN, S., AND E. TAMER (2009): “Inference on Randomly Censored Regression Models Using Conditional Moment Inequalities,” *Journal of Econometrics*, 152, 104–119.
- LIMA, L. (2021): “Demand Estimation with Zeros: Moving Costs in the US,” Harvard University Working Paper.
- MARKHAM, J. W. (1951): “The nature and significance of price leadership,” *The American Economic Review*, 41, 891–905.
- MILLER, N. H., M. REMER, AND M. C. WEINBERG (2020): “The Canned Tuna Cartel,” Discussion paper.
- NEVO, A., AND K. HATZITASKOS (2006): “Why does the average price paid fall during high demand periods?,” Discussion paper, CSIO working paper.
- NEWKEY, W. K. (1990): “Efficient Instrumental Variables Estimation of Nonlinear Models,” *Econometrica*, 58, 809–837.

SHIMOMURA, K.-I., AND J.-F. THISSE (2012): “Competition among the big and the small,”  
*The RAND Journal of Economics*, 43, 329–347.

TIOLE, J. (1988): *The Theory of Industrial Organization*. The MIT Press.

# Supplemental Appendix for “Estimating Demand for Differentiated Products with Zeroes in Market Share Data”

Amit Gandhi, Zhentong Lu, Xiaoxia Shi\*

In this supplemental appendix, we present supporting materials for “Estimating Demand for Differentiated Products with Zeroes in Market Share Data” (hereafter “main text”). The supplemental appendix is organized as follows:

Section [A](#) provides further illustration of the power law pattern in homicide and international trade data sets. This section complements the illustration in Section [2](#) in the main text.

Section [B](#) gives the proofs of Lemmas [1](#) and [2](#) presented in Sections [4](#) and [5](#) in the main text, respectively. Lemma [1](#) establishes the validity of the bounds for the general model. Lemma [2](#) proves that the bounds collapse for the dominant products.

Section [C](#) proves Theorems [1](#) and [2](#) presented in Sections [6](#) and [7](#) in the main text, respectively. Theorem [1](#) establishes the consistency of our proposed estimator and Theorem [2](#) proves the asymptotic normality.

Section [D](#) proves a lemma that establishes Assumption [1](#) in Section [4](#) of the main text for the random coefficient logit model.

Section [E](#) provides analytical evidence for the bound validity in Section [3](#) in the main text. The two lemmas presented here reinforce the numerical proof given in Section [3](#).

## A Further Illustrations of Zipf’s Law

In Figure [3](#) we illustrate this regularity using data from the two other applications that were mentioned in Section [2](#): homicide rates and international trade flows. The left hand graph shows the annual murder rate (per 10,000 people) for each county in the US from 1977-1992 (for details about the data see Dezhbakhsh et al. (2003)). The right hand side graph shows the import trade flows (measured in millions of US dollars) among 160 countries that have a regional trade agreement in the year 2006 (for details about the data see Head et al. (2013)).

---

\* Corresponding author: xshi@ssc.wisc.edu

In each of these two cases we see the characteristic pattern of Zipf’s law - a sharp decay in the frequency for large outcomes and a large mass near zero (with a mode at zero in each case).

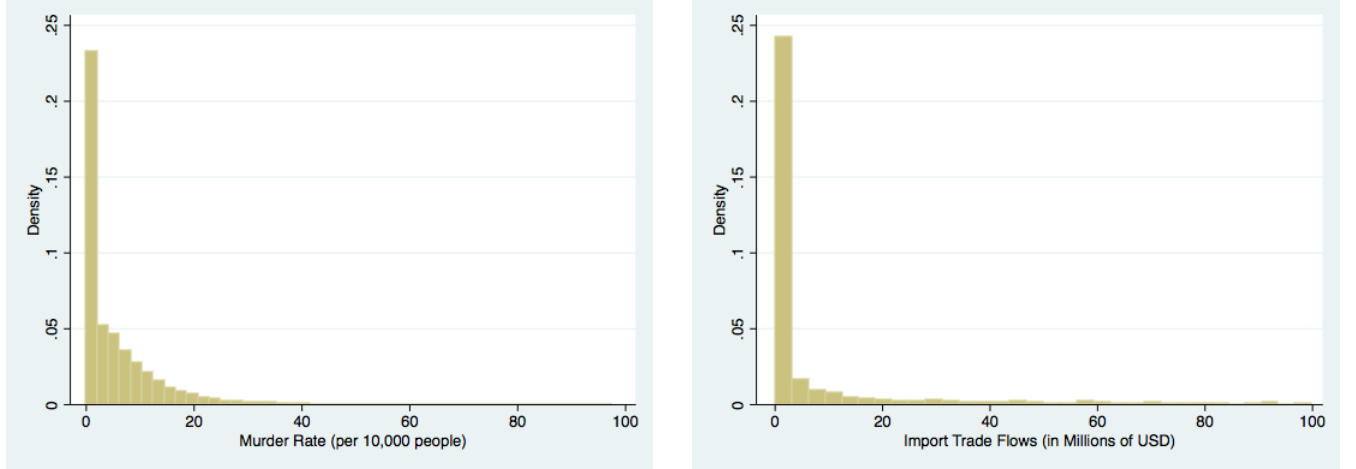


Figure 3: Zipf’s Law in Crime and Trade Data

## B Proof of Lemmas 1 and 2

*Proof of Lemma 1.* We start with the case where Assumption 1 holds. We show the argument for the upper bound only because the lower bound is analogous. Consider the derivation

$$\hat{\delta}_{jt}^u(s_t, \lambda_0) - \delta_{jt}(\pi_t, \lambda_0) = [\log((n_t s_{jt} + \iota_u)/n_t) - \log(\pi_{jt})] + [\check{\delta}_{jt}(\tilde{s}_t, \lambda_0) - \check{\delta}_{jt}(\pi_t, \lambda_0)]. \quad (42)$$

Let  $e_{jt}^u = \check{\delta}_{jt}(\tilde{s}_t, \lambda_0) - \check{\delta}_{jt}(\pi_t, \lambda_0)$ . Then by Assumption 1(b),

$$E[\hat{\delta}_{jt}^u(s_t, \lambda_0) - \delta_{jt}(\pi_t, \lambda_0) - e_{jt}^u | \pi_{jt}, z_{jt}] = E[\log((n_t s_{jt} + \iota_u)/n_t) - \log(\pi_{jt}) | \pi_{jt}, z_{jt}] \geq 0. \quad (43)$$

Since  $E[\xi_{jt} | z_{jt}] = 0$ , we have  $E[\delta_{jt}(\pi_t, \lambda_0) - x_{jt} \beta_0 | z_{jt}] = 0$ . Thus,

$$E[\hat{\delta}_{jt}^u(s_t, \lambda_0) - x_{jt} \beta_0 - e_{jt}^u | z_{jt}] \geq 0. \quad (44)$$

Let  $u_t$  stand for  $\frac{n_t^{1/2}}{T^{1/4}J_t^{1/2}}$ . Now we show that  $\sup_{j,t} u_t |e_{jt}^u| = O_p(1)$ . Consider the derivation:

$$\begin{aligned}
\sup_{j,t} u_t |e_{jt}^u| &= \sup_{j,t} u_t |\check{\delta}_{jt}(\tilde{s}_t, \lambda_0) - \check{\delta}_{jt}(\pi_t, \lambda_0)| \\
&\leq \sup_{j,t} \frac{|\check{\delta}_{jt}(\tilde{s}_t, \lambda_0) - \check{\delta}_{jt}(\pi_t, \lambda_0)|}{\|\tilde{s}_t - \pi_t\|_f \sqrt{J_t}} \sup_t \sqrt{J_t} u_t \|\tilde{s}_t - \pi_t\|_f \\
&\leq O_p(1) \left( \sup_t u_t \sqrt{J_t} \|\tilde{s}_t - s_t\|_f + \sup_{t=1} u_t \sqrt{J_t} \|s_t - \pi_t\|_f \right) \\
&= o_p(1) + O_p(1) \sup_t u_t \sqrt{J_t} \|s_t - \pi_t\|_f,
\end{aligned} \tag{45}$$

where the second inequality holds by Assumption 1(a) and the second equality holds by Assumption 1(b). To bound  $u_t \sqrt{J_t} \|s_t - \pi_t\|_f$ , note that  $n_t s_{jt}$  is a binomial random variable with parameters  $(n_t, \pi_{jt})$ . Thus,

$$\begin{aligned}
\Pr \left( \sup_{t=1, \dots, T} u_t J_t^{1/2} \|s_t - \pi_t\|_f > \varepsilon \right) &\leq \sum_{t=1}^T \Pr(u_t J_t^{1/2} \|s_t - \pi_t\|_f > \varepsilon) \\
&\leq \sum_{t=1}^T \frac{128 u_t^4 J_t^2 (3n_t^2 + n_t)}{n_t^4 \varepsilon^4} \\
&\leq \frac{512}{\varepsilon^4},
\end{aligned} \tag{46}$$

where the second inequality holds by Lemma 8. The expression in the last line does not depend on  $T$  and it can be made arbitrarily small by making  $\varepsilon$  big. This shows that

$$\sup_{t=1, \dots, T} u_t J_t^{1/2} \|s_t - \pi_t\|_f = O_p(1) \tag{47}$$

which implies  $\sup_{j,t} u_t |e_{jt}^u| = O_p(1)$  when combined with (45).

Now we move on to the case where Assumption 2 holds instead. In this case, the upper bound and the lower bound need slightly different arguments. For the upper bound, consider the derivation:

$$\begin{aligned}
\hat{\delta}_{jt}^u(s_t, \lambda_0) - \delta_{jt}(\pi_t, \lambda_0) &= \delta_{jt}(\tilde{s}_{jt}, \lambda_0) - \delta_{jt}(\pi_t, \lambda_0) + \log((n_t s_{jt} + \iota_u)/n_t) - \log(\tilde{s}_{jt}) \\
&\geq \delta_{jt}(\tilde{s}_{jt}, \lambda_0) - \delta_{jt}(\pi_t, \lambda_0),
\end{aligned} \tag{48}$$

where the inequality holds because  $\tilde{s}_{jt} = s_{jt} + 1/n_t$  and  $\iota_u > 1$  both by Assumption 2(c). Equation (48) combined with Assumption 2(b) implies the first line of (25) with  $e_{jt}^u = 0$ .

For the lower bound, consider the derivation:

$$\hat{\delta}_{jt}^\ell(s_t, \lambda_0) - \delta_{jt}(\pi_t, \lambda_0) = \check{\delta}_{jt}(\tilde{s}_t, \lambda_0) - \check{\delta}_{jt}(\pi_t, \lambda_0) + \log((n_t s_{jt} + \iota_\ell)/n_t) - \log(\pi_{jt}) \quad (49)$$

By  $\iota_\ell \leq \bar{\iota}_\ell$  (Assumption 2(c)) and the definition of  $\bar{\iota}_\ell$ , we have

$$E[\log((n_t s_{jt} + \iota_\ell)/n_t) - \log(\pi_{jt}) | n_t, \pi_t] \leq 0. \quad (50)$$

This combined with Assumption 2(a) implies the second line of (25) with  $e_{jt}^\ell = 0$ .  $\square$

*Proof of Lemma 2.* Observe that

$$\begin{aligned} n_t |\hat{\delta}_{jt}^u(s_t, \lambda) - \hat{\delta}_{jt}^\ell(s_t, \lambda)| &= n_t |\log(s_{jt} + \iota_u/n_t) - \log(s_{jt} + \iota_\ell/n_t)| \\ &\leq \frac{1}{s_{jt} + \iota_\ell/n_t} (\iota_u - \iota_\ell), \end{aligned} \quad (51)$$

using the concavity of the logarithm function. Thus

$$\begin{aligned} &\sup_{j=1, \dots, J_t; t=1, \dots, T} \sup_{\lambda} n_t |\hat{\delta}_{jt}^u(s_t, \lambda) - \hat{\delta}_{jt}^\ell(s_t, \lambda)| 1\{z_{jt} \in \mathcal{Z}_0\} \\ &\leq \sup_{j=1, \dots, J_t; t=1, \dots, T} \frac{\iota_u - \iota_\ell}{s_{jt} + \iota_\ell/n_t} 1\{z_{jt} \in \mathcal{Z}_0\} \\ &\leq \frac{\iota_u - \iota_\ell}{\inf_{j=1, \dots, J_t; t=1, \dots, T} \{(s_{jt} + \iota_\ell/n_t) 1\{z_{jt} \in \mathcal{Z}_0\} + 1\{z_{jt} \notin \mathcal{Z}_0\}\}}. \end{aligned} \quad (52)$$

The denominator of (52) is greater than or equal to

$$\inf_{j=1, \dots, J_t; t=1, \dots, T} \{\pi_{jt} 1\{z_{jt} \in \mathcal{Z}_0\} + 1\{z_{jt} \notin \mathcal{Z}_0\}\} - \sup_{j=1, \dots, J_t; t=1, \dots, T} |\pi_{jt} - s_{jt} - \iota_\ell/n_t| 1\{z_{jt} \in \mathcal{Z}_0\}. \quad (53)$$

Consider that

$$\begin{aligned} &\Pr \left( \inf_{j=1, \dots, J_t; t=1, \dots, T} \{\pi_{jt} 1\{z_{jt} \in \mathcal{Z}_0\} + 1\{z_{jt} \notin \mathcal{Z}_0\}\} < \underline{\varepsilon}_0 \right) \\ &\leq \sum_{j=1, \dots, J_t; t=1, \dots, T} \Pr(\pi_{jt} 1\{z_{jt} \in \mathcal{Z}_0\} + 1\{z_{jt} \notin \mathcal{Z}_0\} < \underline{\varepsilon}_0) \\ &= \sum_{j=1, \dots, J_t; t=1, \dots, T} \Pr(\pi_{jt} < \underline{\varepsilon}_0 | z_{jt} \in \mathcal{Z}_0) P(z_{jt} \in \mathcal{Z}_0) \\ &= 0, \end{aligned} \quad (54)$$

where the first equality holds since  $1 \geq \underline{\varepsilon}_0$ , and the second equality holds by Assumption 3.

Also consider the derivation:

$$\begin{aligned}
& \Pr \left( \sup_{j=1, \dots, J_t; t=1, \dots, T} |\pi_{jt} - s_{jt} - \iota_\ell/n_t| 1\{z_{jt} \in \mathcal{Z}_0\} > \underline{\varepsilon}_0/2 \right) \\
& \leq \Pr \left( \sup_{t=1, \dots, T} \|s_t - \pi_t\|_f > \underline{\varepsilon}_0/4 \right) + \Pr \left( \sup_{t=1, \dots, T} \iota_\ell/n_t > \underline{\varepsilon}_0/4 \right) \\
& \leq \sum_{t=1}^T \Pr(\|s_t - \pi_t\|_f > \underline{\varepsilon}_0/4) + o(1) \\
& \leq \sum_{t=1}^T \frac{128(3n_t^2 + n_t) \times 4^4}{n_t^4 \underline{\varepsilon}_0^4} + o(1) \\
& \rightarrow 0, \tag{55}
\end{aligned}$$

where the first inequality holds by the triangular inequality, the second inequality and the convergence hold Assumption 4(g), and the third inequality holds by Lemma 8.

Equations (52), (54), and (55) together imply that

$$\Pr \left( \sup_{j=1, \dots, J_t; t=1, \dots, T} \sup_{\lambda} n_t |\hat{\delta}_{jt}^u(s_t, \lambda) - \hat{\delta}_{jt}^\ell(s_t, \lambda)| 1\{z_{jt} \in \mathcal{Z}_0\} > \frac{\iota_u - \iota_\ell}{\underline{\varepsilon}_0/2} \right) \rightarrow 0. \tag{56}$$

This proves the lemma.  $\square$

## C Proofs of the Theorems

In this section, we prove the theorems that establish the consistency and the asymptotic normality of our proposed estimator.

### C.1 Proof of Theorem 1: Consistency

The proof of Theorem 1 uses the following lemma which is proved in Section C.2 below.

**Lemma 3.** *Suppose that either Assumption 1 holds and  $T^{-1} \sum_{t=1}^T J_t^2 / \bar{J}_T^2$  is bounded, or Assumption 2 holds and  $\sup_{t=1, \dots, T} J_t$  is bounded. Also suppose that Assumptions 3-6 hold. Then,*

$$\text{(i) } \sup_{\theta \in \Theta} \sup_{g \in \mathcal{G}_0} |\bar{m}_T^u(\theta, g) - \bar{m}_T(\theta, g)| = o_p(1) \text{ and } \sup_{\theta \in \Theta} \sup_{g \in \mathcal{G}_0} |\bar{m}_T^\ell(\theta, g) - \bar{m}_T(\theta, g)| = o_p(1).$$

$$\text{(ii) } \sum_{g \in \mathcal{G}} \mu(g) [\bar{m}_T^u(\theta_0, g)]_-^2 = o_p(1) \text{ and } \sum_{g \in \mathcal{G}} \mu(g) [\bar{m}_T^\ell(\theta_0, g)]_+^2 = o_p(1).$$

*Proof of Theorem 1.* First note that  $\widehat{Q}_T(\theta_0) = \sum_{g \in \mathcal{G}} \mu(g) [\bar{m}_T^u(\theta_0, g)]_-^2 + \sum_{g \in \mathcal{G}} \mu(g) [\bar{m}_T^\ell(\theta_0, g)]_+^2$ .



Thus, Lemma 3(b) implies that

$$\widehat{Q}_T(\theta_0) = o_p(1). \quad (57)$$

Define an auxiliary criterion function:

$$\widehat{Q}_{0,T}(\theta) = \sum_{g \in \mathcal{G}_0} \left\{ \left( [\bar{m}_T^u(\theta, g)]_-^2 + [\bar{m}_T^\ell(\theta, g)]_+^2 \right) \mu(g) \right\}.$$

Below we show that

$$\sup_{\theta \in \Theta} \left| \sqrt{\widehat{Q}_{0,T}(\theta)} - \sqrt{\widehat{Q}_T^*(\theta)} \right| = o_p(1). \quad (58)$$

Consider an arbitrary  $c > 0$ . The theorem is implied by the following derivation:

$$\begin{aligned} \Pr \left( \|\widehat{\theta}_T^s - \theta_0^s\| > c \right) &\leq \Pr \left( \sqrt{\widehat{Q}_T^*(\widehat{\theta}_T)} \geq \sqrt{C(c)} \right) \\ &= \Pr \left( \sqrt{\widehat{Q}_T^*(\widehat{\theta}_T)} - \sqrt{\widehat{Q}_{0,T}(\widehat{\theta}_T)} + \sqrt{\widehat{Q}_{0,T}(\widehat{\theta}_T)} \geq \sqrt{C(c)} \right) \\ &\leq \Pr \left( \sup_{\theta \in \Theta} \left| \sqrt{\widehat{Q}_{0,T}(\theta)} - \sqrt{\widehat{Q}_T^*(\theta)} \right| + \sqrt{\widehat{Q}_{0,T}(\widehat{\theta}_T)} \geq \sqrt{C(c)} \right) \\ &\leq \Pr \left( \sup_{\theta \in \Theta} \left| \sqrt{\widehat{Q}_{0,T}(\theta)} - \sqrt{\widehat{Q}_T^*(\theta)} \right| + \sqrt{\widehat{Q}_T(\widehat{\theta}_T)} \geq \sqrt{C(c)} \right) \\ &\leq \Pr \left( \sup_{\theta \in \Theta} \left| \sqrt{\widehat{Q}_{0,T}(\theta)} - \sqrt{\widehat{Q}_T^*(\theta)} \right| + \sqrt{\widehat{Q}_T(\theta_0)} \geq \sqrt{C(c)} \right) \\ &\leq \Pr \left( \sup_{\theta \in \Theta} \left| \sqrt{\widehat{Q}_{0,T}(\theta)} - \sqrt{\widehat{Q}_T^*(\theta)} \right| \geq \sqrt{C(c)}/2 \right) + \Pr \left( \widehat{Q}_T(\theta_0) \geq C(c)/4 \right) \\ &\rightarrow 0, \end{aligned} \quad (59)$$

where the first inequality holds by Assumption 7, the third inequality holds because  $\widehat{Q}_T(\widehat{\theta}_T)$  differs from  $\widehat{Q}_{0,T}(\widehat{\theta}_T)$  only in that the former takes the summation over a larger range, the fourth inequality holds because  $\widehat{Q}_T(\widehat{\theta}_T) \leq \widehat{Q}_T(\theta_0)$  by the definition of  $\widehat{\theta}_T$  and the convergence holds by (57) and (58).

Now we show (58). Consider the derivation

$$\begin{aligned} &\sup_{\theta \in \Theta} \left| \sqrt{\widehat{Q}_{0,T}(\theta)} - \sqrt{\widehat{Q}_T^*(\theta)} \right| \\ &= \sup_{\theta \in \Theta} \left| \sqrt{\sum_{g \in \mathcal{G}_0} \mu(g) \left\{ [\bar{m}_T^u(\theta, g)]_-^2 + [\bar{m}_T^\ell(\theta, g)]_+^2 \right\}} - \sqrt{\sum_{g \in \mathcal{G}_0} \mu(g) \left\{ \bar{m}_T(\theta, g)^2 \right\}} \right| \end{aligned}$$

$$\begin{aligned}
&\leq \sup_{\theta \in \Theta} \left| \sqrt{\sum_{g \in \mathcal{G}_0} \mu(g) \left\{ \left( \sqrt{[\bar{m}_T^u(\theta, g)]_-^2 + [\bar{m}_T^\ell(\theta, g)]_+^2} - |\bar{m}_T(\theta, g)| \right)^2 \right\}} \right| \\
&\leq \sup_{\theta \in \Theta} \left| \sqrt{\sum_{g \in \mathcal{G}_0} \mu(g) \left\{ ([\bar{m}_T^u(\theta, g)]_- - [\bar{m}_T(\theta, g)]_-)^2 + ([\bar{m}_T^\ell(\theta, g)]_+ - [\bar{m}_T(\theta, g)]_+)^2 \right\}} \right| \\
&\leq \sup_{\theta \in \Theta} \left| \sqrt{\sum_{g \in \mathcal{G}_0} \mu(g) \left\{ (\bar{m}_T^u(\theta, g) - \bar{m}_T(\theta, g))^2 + (\bar{m}_T^\ell(\theta, g) - \bar{m}_T(\theta, g))^2 \right\}} \right| \\
&\leq \sqrt{\sup_{\theta \in \Theta} \sup_{g \in \mathcal{G}_0} |\bar{m}_T^u(\theta, g) - \bar{m}_T(\theta, g)|^2 + \sup_{\theta \in \Theta} \sup_{g \in \mathcal{G}_0} |\bar{m}_T^\ell(\theta, g) - \bar{m}_T(\theta, g)|^2} \\
&\xrightarrow{p} 0, \tag{60}
\end{aligned}$$

where the first inequality holds by the triangular inequality for the norm

$$\|a(\cdot)\| := \sqrt{\sum_{g \in \mathcal{G}_0} \mu(g) a(g)^2 / \sum_{g \in \mathcal{G}_0} \mu(g)},$$

the second inequality holds by the triangular inequality for the Euclidean norm, the third inequality holds because  $|[x]_- - [y]_-| \leq |x - y|$  and  $[x]_+ = [-x]_-$ , and the fourth inequality holds because  $\mu : \mathcal{G} \rightarrow [0, 1]$  is a probability measure on  $\mathcal{G}$  and  $\mathcal{G}_0 \subseteq \mathcal{G}$ , and the convergence holds by Lemma 3(a). Therefore (58) is proved.  $\square$

## C.2 Proof of Lemma 3

*Proof of Lemma 3.* First we show part (a). Let  $\sup_{j,t:z_{jt} \in \mathcal{Z}_0}$  abbreviate  $\sup_{t=1,\dots,T} \sup_{j=1,\dots,J_t:z_{jt} \in \mathcal{Z}_0}$ . Consider the derivation:

$$\begin{aligned}
&\sup_{\theta \in \Theta} \sup_{g \in \mathcal{G}_0} |\bar{m}_T^u(\theta, g) - \bar{m}_T(\theta, g)| \\
&= \sup_{\lambda \in \Lambda} \sup_{g \in \mathcal{G}_0} \left| \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\hat{\delta}_{jt}^u(s_t, \lambda) - \delta_{jt}(\pi_t, \lambda)) g(z_{jt}) \right| \\
&\leq \sup_{\lambda \in \Lambda} \sup_{j,t:z_{jt} \in \mathcal{Z}_0} |\hat{\delta}_{jt}^u(s_t, \lambda) - \delta_{jt}(\pi_t, \lambda)| \\
&\leq \sup_{j,t:z_{jt} \in \mathcal{Z}_0} |\log(s_{jt} + \iota_u/n_t) - \log(\tilde{s}_{jt})| + \sup_{\lambda \in \Lambda} \sup_{j,t:z_{jt} \in \mathcal{Z}_0} |\delta_{jt}(\tilde{s}_t, \lambda) - \delta_{jt}(\pi_t, \lambda)|,
\end{aligned}$$

where the first inequality holds by the definition of  $\mathcal{G}_0$ .

Assumptions 4(f) and  $0 < \iota_u < \infty$  (Assumption 1(b) or Assumption 2(b)) together imply

that  $\sup_{j,t:z_{jt} \in \mathcal{Z}_0} |s_{jt} + \iota_u/n_t - \tilde{s}_{jt}| \rightarrow_p 0$ . Also, by equation (47) and Assumption 4(f)

$$\frac{n_t^{1/2}}{T^{1/4}} \sup_t \|\tilde{s}_t - \pi_t\|_f = O_p(1). \quad (61)$$

These, and Assumptions 3, 4(g), and 5(a) together imply that

$$\Pr \left( \inf_{j,t:z_{jt} \in \mathcal{Z}_0} \pi_{jt} \wedge \pi_{0t} > \underline{\varepsilon}_0, \inf_{j,t:z_{jt} \in \mathcal{Z}_0} s_{jt} + \iota_u/n_t > \underline{\varepsilon}_0/2, \inf_{j,t:z_{jt} \in \mathcal{Z}_0} \tilde{s}_{jt} \wedge \tilde{s}_{0t} > \underline{\varepsilon}_0/2 \right) \rightarrow 1.$$

This combined with (61), Assumptions 4(g), and Assumption 6 implies that

$$\sup_{\lambda \in \Lambda} \sup_{j,t:z_{jt} \in \mathcal{Z}_0} |\delta_{jt}(\tilde{s}_t, \lambda) - \delta_{jt}(\pi_t, \lambda)| \rightarrow_p 0.$$

Also, we have

$$\sup_{j,t:z_{jt} \in \mathcal{Z}_0} (|\log(s_{jt} + \iota_u/n_t) - \log(\tilde{s}_{jt})|) \rightarrow_p 0.$$

because the logarithm function is uniformly continuous on the closed interval  $[\underline{\varepsilon}_0/2, 1]$ . Therefore, the first convergence in Lemma 3(a) holds. The second convergence holds by analogous arguments.

Now we show part (b). We separate the two cases, one where Assumption 1 is satisfied and the other where Assumption 2 is satisfied and  $J_t$  is bounded.

Case 1: Assumption 1 is satisfied. In this case, the arguments for the first convergence and the second convergence in part (b) are exactly analogous. Thus, we only discuss the first. Consider the derivation:

$$\begin{aligned} \bar{m}_T^u(\theta_0, g) &= \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\hat{\delta}_{jt}^u(s_t, \lambda_0) - x'_{jt}\beta_0)g(z_{jt}) \\ &\geq \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} \xi_{jt}g(z_{jt}) + \\ &\quad \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\log(s_{jt} + \iota_u/n_t) - \log(\pi_{jt}))g(z_{jt}) + \\ &\quad \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\check{\delta}_{jt}(\tilde{s}_t, \lambda_0) - \check{\delta}_{jt}(\pi_t, \lambda_0))g(z_{jt}), \end{aligned} \quad (62)$$

where the inequality holds because  $\iota_u \geq \underline{\iota}_u$ ,  $\xi_{jt} = \delta_{jt}(\pi_t, \lambda_0) - x'_{jt}\beta_0$  and  $\check{\delta}(\cdot, \lambda) = \delta(\cdot, \lambda) -$

$\log(\cdot_j)$ . We analyze the three summands one by one. For the first summand, observe that  $E[\xi_{jt}^2] \leq M$  by Assumption 4(e). We can then apply Lemma 7 in Appendix C.4 (with  $w_{jt} = \xi_{jt}$ ) and get,

$$E \sup_{g \in \mathcal{G}} \left| \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} \{\xi_{jt} g(z_{jt}) - E[\xi_{jt} g(z_{jt})]\} \right|^2 \leq \frac{CM \sum_{t=1}^T J_t^2}{T^2 \bar{J}_T^2} = O(T^{-1}), \quad (63)$$

where the equality holds because we assume  $T^{-1} J_t^2 / \bar{J}_T^2$  is bounded when Assumption 1 holds. Lemma 7 applies due to Assumptions 4(c)-(e). Also, by Assumption 4(c),  $E[\xi_{jt} g(z_{jt})] = 0$ . This and (63) together imply that

$$\sup_{g \in \mathcal{G}} \left| \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} \xi_{jt} g(z_{jt}) \right| = O_p(T^{-1/2}). \quad (64)$$

Similar arguments apply to the second summand in (62) and yields

$$\begin{aligned} & E \sup_{g \in \mathcal{G}} \left| \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\log(s_{jt} + \underline{\iota}_u/n_t) - \log(\pi_{jt})) g(z_{jt}) - E[(\log(s_{jt} + \underline{\iota}_u/n_t) - \log(\pi_{jt})) g(z_{jt})] \right|^2 \\ & \leq \frac{C \sum_{t=1}^T J_t^2}{T^2 \bar{J}_T^2} \max_{j,t} E[(\log(s_{jt} + \underline{\iota}_u/n_t) - \log(\pi_{jt}))^2] \\ & \leq \frac{2C \sum_{t=1}^T J_t^2}{T^2 \bar{J}_T^2} \max_t [|\log(\underline{\iota}_u/n_t)|^2 + |\log(\underline{\varepsilon}_1/n_t)|^2] \\ & \leq \frac{4C \sum_{t=1}^T J_t^2 (2(\log \bar{n}_T)^2 + (\log \underline{\iota}_u)^2 + (\log \underline{\varepsilon}_1)^2)}{T^2 \bar{J}_T^2} \\ & \rightarrow 0, \end{aligned} \quad (65)$$

where the second inequality holds by  $s_{jt} \in [0, 1]$  and Assumption 5(b) and the convergence holds by Assumptions 4(f) and the boundedness of  $T^{-1} \sum_{t=1}^T J_t^2 / \bar{J}_T^2$ . By the definition of  $\underline{\iota}_u$ , we have  $E[(\log(s_{jt} + \underline{\iota}_u/n_t) - \log(\pi_{jt})) | \pi_{jt}, z_{jt}] \geq 0$ , which then implies that  $E[(\log(s_{jt} + \underline{\iota}_u/n_t) - \log(\pi_{jt})) g(z_{jt})] \geq 0$  for all  $g \in \mathcal{G}$ . Therefore, for any  $c > 0$ ,

$$\lim_{T \rightarrow \infty} \Pr \left( \inf_{g \in \mathcal{G}} \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\log(s_{jt} + \underline{\iota}_u/n_t) - \log(\pi_{jt})) g(z_{jt}) < -c \right) = 0. \quad (66)$$

For the third summand in (62), consider the derivation

$$\sup_{g \in \mathcal{G}} \left| \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\check{\delta}_{jt}(\tilde{s}_t, \lambda_0) - \check{\delta}_{jt}(\pi_t, \lambda_0)) g(z_{jt}) \right| \leq \sup_{j,t} |\check{\delta}_{jt}(\tilde{s}_t, \lambda_0) - \check{\delta}_{jt}(\pi_t, \lambda_0)| \rightarrow_p 0, \quad (67)$$

by (61) and Assumptions 1(a) and 4(g). Finally, (62), (64), (66), and (67) combined imply that for any  $c > 0$ ,

$$\lim_{T \rightarrow \infty} \Pr \left( \inf_{g \in \mathcal{G}} \bar{m}_T^u(\theta_0, g) < -c \right) = 0,$$

which then implies the first convergence in Lemma 3(b) since  $[\bar{m}_T^u(\theta_0, g)]_- = \max\{0, -\bar{m}_T^u(\theta_0, g)\}$ .

Case 2: Assumption 2 is satisfied. We begin with the first convergence in Lemma 3(b). Consider the decomposition:

$$\begin{aligned} \bar{m}_T^u(\theta_0, g) &= \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\hat{\delta}_{jt}^u(s_t, \lambda_0) - x'_{jt} \beta_0) g(z_{jt}) \\ &= \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} \xi_{jt} g(z_{jt}) + \\ &\quad \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\log(s_{jt} + \nu_u/n_t) - \log(\tilde{s}_{jt})) g(z_{jt}) + \\ &\quad \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\delta_{jt}(\tilde{s}_t, \lambda_0) - \delta_{jt}(\pi_t, \lambda_0)) g(z_{jt}). \end{aligned} \quad (68)$$

The first summand is  $O_p(T^{-1/2})$  uniformly over  $g \in \mathcal{G}$  by (64). The second summand is nonnegative almost surely because  $\tilde{s}_{jt} = s_{jt} + 1/n_t$  and  $\nu_u \geq 1$  (Assumption 2(d)). For the third summand, similar to (65), we get for some generic constant  $C$ ,

$$\begin{aligned} &E \sup_{g \in \mathcal{G}} \left| \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\delta_{jt}(\tilde{s}_t, \lambda_0) - \delta_{jt}(\pi_t, \lambda_0)) g(z_{jt}) - E[(\delta_{jt}(\tilde{s}_t, \lambda_0) - \delta_{jt}(\pi_t, \lambda_0)) g(z_{jt})] \right|^2 \\ &\leq \frac{C \sum_{t=1}^T J_t^2}{T^2 \bar{J}_T^2} \max_{j,t} E[(\delta_{jt}(\tilde{s}_t, \lambda_0) - \delta_{jt}(\pi_t, \lambda_0))^2] \\ &\leq \frac{2CC_0 \sum_{t=1}^T J_t^2 \log(\bar{n}_T)^2}{T^2 \bar{J}_T^2} \\ &\rightarrow 0, \end{aligned} \quad (69)$$

where the second inequality holds by Assumption 2(d) also using Assumptions 2(c) and

5(b), and the convergence holds by Assumption 4(g) and the boundedness of  $\sup_{t=1,\dots,T} J_t$ . Moreover, Assumption 2(b) implies that  $E[(\delta_{jt}(\tilde{s}_t, \lambda_0) - \delta_{jt}(\pi_t, \lambda_0))g(z_{jt})] \geq 0$ . This combined with (69) implies that, for any  $c > 0$ ,

$$\lim_{T \rightarrow \infty} \Pr \left( \inf_{g \in \mathcal{G}} \frac{1}{T\bar{J}} \sum_{t=1}^T \sum_{j=1}^{J_t} (\delta_{jt}(\tilde{s}_t, \lambda_0) - \delta_{jt}(\pi_t, \lambda_0))g(z_{jt}) < -c \right) = 0. \quad (70)$$

This combined with the arguments for the first two summands of (62) above yields: for any  $c > 0$ ,

$$\lim_{T \rightarrow \infty} \Pr \left( \inf_{g \in \mathcal{G}} \bar{m}_T^u(\theta_0, g) < -c \right) = 0,$$

which then implies the first convergence in Lemma 3(b) because  $[\bar{m}_T^u(\theta_0, g)]_- = \max\{0, -\bar{m}_T^u(\theta_0, g)\}$ .

Now we show the second convergence in Lemma 3(b) for Case 2. Note that

$$\begin{aligned} \bar{m}_T^\ell(\theta_0, g) &= \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\hat{\delta}_{jt}^\ell(s_t, \lambda_0) - x'_{jt}\beta_0)g(z_{jt}) \\ &\leq \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} \xi_{jt}g(z_{jt}) + \\ &\quad \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\log(s_{jt} + \bar{\nu}_\ell/n_t) - \log(\pi_{jt}))g(z_{jt}) + \\ &\quad \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\check{\delta}_{jt}(\tilde{s}_t, \lambda_0) - \check{\delta}_{jt}(\pi_t, \lambda_0))g(z_{jt}), \end{aligned} \quad (71)$$

where the inequality holds because  $\nu_\ell \leq \bar{\nu}_\ell$  by Assumption 2(d). The first summand is  $O_p(T^{-1/2})$  by (64). Arguments analogous to those for (66) apply to the second summand to yield, for any  $c > 0$ ,

$$\lim_{T \rightarrow \infty} \Pr \left( \inf_{g \in \mathcal{G}} \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\log(s_{jt} + \bar{\nu}_\ell/n_t) - \log(\pi_{jt}))g(z_{jt}) > c \right) = 0. \quad (72)$$

For the third summand in (71), we can apply the same arguments as those for (70) where we use Assumption 2(a) in place of Assumption 2(b). Such arguments yield, for all  $c > 0$ ,

$$\lim_{T \rightarrow \infty} \Pr \left( \inf_{g \in \mathcal{G}} \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\check{\delta}_{jt}(\tilde{s}_t, \lambda_0) - \check{\delta}_{jt}(\pi_t, \lambda_0))g(z_{jt}) > c \right) = 0. \quad (73)$$

Therefore, for any  $c > 0$ ,

$$\lim_{T \rightarrow \infty} \Pr \left( \inf_{g \in \mathcal{G}} \bar{m}_T^\ell(\theta_0, g) > c \right) = 0,$$

which then implies the second convergence in Lemma 3(b) because we have  $[\bar{m}_T^\ell(\theta_0, g)]_+ = \max\{0, \bar{m}_T^\ell(\theta_0, g)\}$ .  $\square$

### C.3 Proof of Asymptotic Normality

To prove Theorem 2, we first give an auxiliary theorem that shows the convergence rate of  $\hat{\theta}_T$ .

**Theorem 3.** *Suppose that either Assumption 1 holds and  $T^{-1} \sum_{t=1}^T J_t^2 / \bar{J}_T^2$  is bounded, or Assumption 2 holds and  $\sup_{t=1, \dots, T} J_t$  is bounded. Also suppose that Assumptions 3-9 hold. Then we have  $\hat{\theta}_T^s - \theta_0^s = O_p(T^{-1/2})$ .*

Theorem 3 is proved using the following three lemmas. Theorem 3 and two of the lemmas together imply Theorem 2 as we explain immediately below. We give the proofs of Theorem 3 and the three lemmas in turn following the proof of Theorem 2.

**Lemma 4.** *Suppose that either Assumption 1 holds and  $T^{-1} \sum_{t=1}^T J_t^2 / \bar{J}_T^2$  is bounded, or Assumption 2 holds and  $\sup_{t=1, \dots, T} J_t$  is bounded. Also suppose that Assumptions 3-9 hold. Then we have for any sequence  $\theta_T$  such that  $\theta_T^s - \theta_0^s = O_p(T^{-1/2})$ ,  $\hat{Q}_T(\theta_T) - \hat{Q}_{0,T}(\theta_T) = o_p(T^{-1})$ .*

**Lemma 5.** *Suppose that either Assumption 1 holds and  $T^{-1} \sum_{t=1}^T J_t^2 / \bar{J}_T^2$  is bounded, or Assumption 2 holds and  $\sup_{t=1, \dots, T} J_t$  is bounded. Also suppose that Assumptions 3-9 hold. Then we have*

(a) *for an open ball  $B_c(\theta_0^s)$  of radius  $c > 0$  around  $\theta_0^s$ , we have that*

$$\sup_{\theta \in \Theta: \theta^s \in B_c(\theta_0^s)} \left| \sqrt{\hat{Q}_{0,T}(\theta)} - \sqrt{\hat{Q}_T^*(\theta)} \right| = o_p(T^{-1/2}), \text{ and}$$

(b)  $\hat{Q}_T^*(\theta_0) = O_p(T^{-1})$ .

**Lemma 6.** *Suppose that either Assumption 1 holds and  $T^{-1} \sum_{t=1}^T J_t^2 / \bar{J}_T^2$  is bounded, or Assumption 2 holds and  $\sup_{t=1, \dots, T} J_t$  is bounded. Also suppose that Assumptions 3-8 hold. For any sequence of random vectors  $\theta_T$  such that  $\|\theta_T^s - \theta_0^s\| \rightarrow_p 0$ , we have*

(a)  $\hat{Q}_T^*(\theta_T) - \hat{Q}_T^*(\theta_0) = (\theta_T^s - \theta_0^s)' \hat{\Upsilon}_T (\theta_T^s - \theta_0^s) + 2W_T'(\theta_T^s - \theta_0^s) + o_p(1) \|\theta_T^s - \theta_0^s\|^2$ , where

$$\hat{\Upsilon}_T = \sum_{g \in \mathcal{G}_0} \mu(g) \hat{\Gamma}_T(g) \hat{\Gamma}_T(g)'$$

$$W_T = \sum_{g \in \mathcal{G}_0} \mu(g) \bar{m}_T(\theta_0, g) \widehat{\Gamma}_T(g)$$

$$\widehat{\Gamma}_T(g) = (T \bar{J}_T)^{-1} \sum_{t=1}^T \sum_{j=1}^{J_t} g(z_{jt}) \partial m_{jt}(\lambda_0), \text{ and}$$

(b)  $\widehat{\Upsilon}_T \rightarrow_p \Upsilon$  and  $T^{1/2}W_T \rightarrow_d N(0, V)$ .

*Proof of Theorem 2.* We use Theorem 2 of Sherman (1993) to prove the theorem. By Theorem 2 of Sherman (1993), the conclusion of our Theorem 2 holds under two conditions:

(i)  $\|\widehat{\theta}_T^s - \theta_0^s\| = O_p(T^{-1/2})$ ,

(ii) uniformly over  $\theta^s$  in a  $O_p(T^{-1/2})$  neighborhood of  $\theta_0^s$ ,  $\widehat{Q}_T(\theta) - \widehat{Q}_T(\theta_0) = (\theta^s - \theta_0^s)' \Upsilon(\theta^s - \theta_0^s) + 2T^{-1/2} B_T'(\theta^s - \theta_0^s) + o_p(T^{-1})$  for a random vector  $B_T$  such that  $B_T \rightarrow_d N(0, V)$ .

Condition (i) is implied by Theorem 3. To establish condition (ii), consider the derivation: for any sequence  $\theta_T$  such that  $\theta_T^s - \theta_0^s = O_p(T^{-1/2})$ ,

$$\begin{aligned} \widehat{Q}_T(\theta_T) - \widehat{Q}_T(\theta_0) &= [\widehat{Q}_T(\theta_T) - \widehat{Q}_{0,T}(\theta_T)] + [\widehat{Q}_{0,T}(\theta_T) - \widehat{Q}_T^*(\theta_T)] + \\ &\quad [\widehat{Q}_T^*(\theta_T) - \widehat{Q}_T^*(\theta_0)] + [\widehat{Q}_T^*(\theta_0) - \widehat{Q}_{0,T}(\theta_0)] + [\widehat{Q}_{0,T}(\theta_0) - \widehat{Q}_T(\theta_0)] \\ &= o_p(T^{-1}) + [\widehat{Q}_{0,T}(\theta_T) - \widehat{Q}_T^*(\theta_T)] + \\ &\quad [\widehat{Q}_T^*(\theta_T) - \widehat{Q}_T^*(\theta_0)] + [\widehat{Q}_T^*(\theta_0) - \widehat{Q}_{0,T}(\theta_0)] + o_p(T^{-1}), \end{aligned} \quad (74)$$

where the second equality holds by Lemma 4. For the summand  $[\widehat{Q}_{0,T}(\theta_T) - \widehat{Q}_T^*(\theta_T)]$ , consider the derivation:

$$\begin{aligned} \widehat{Q}_{0,T}(\theta_T) - \widehat{Q}_T^*(\theta_T) &= \left( \sqrt{\widehat{Q}_{0,T}(\theta_T)} - \sqrt{\widehat{Q}_T^*(\theta_T)} \right)^2 + 2 \left( \sqrt{\widehat{Q}_{0,T}(\theta_T)} - \sqrt{\widehat{Q}_T^*(\theta_T)} \right) \left( \sqrt{\widehat{Q}_T^*(\theta_T)} \right) \\ &= o_p(T^{-1}) + o_p(T^{-1/2}) \sqrt{\widehat{Q}_T^*(\theta_T) - \widehat{Q}_T^*(\theta_0) + \widehat{Q}_T^*(\theta_0)} \\ &= o_p(T^{-1}) + o_p(T^{-1/2}) \sqrt{\widehat{Q}_T^*(\theta_T) - \widehat{Q}_T^*(\theta_0) + O_p(T^{-1})} \\ &= o_p(T^{-1}) + o_p(T^{-1/2}) \sqrt{O_p(T^{-1}) + O_p(T^{-1})} \\ &= o_p(T^{-1}), \end{aligned} \quad (75)$$

where the second equality holds by Lemma 5(a), the third equality holds by Lemma 5(b), and the fourth equality holds by Lemma 6(a)-(b). Similar arguments show that the summand  $[\widehat{Q}_{0,T}(\theta_0) - \widehat{Q}_T^*(\theta_0)] = o_p(T^{-1})$ . Therefore,

$$\widehat{Q}_T(\theta) - \widehat{Q}_T(\theta_0) = o_p(T^{-1}) + \widehat{Q}_T^*(\theta_T) - \widehat{Q}_T^*(\theta_0) \quad (76)$$



This combined with Lemma 6(a)-(b) shows the condition (ii) where  $B_T = T^{1/2}W_T$ . This concludes the proof of Theorem 2.  $\square$

*Proof of Theorem 3.* We prove Theorem 3 using Lemmas 4-6. The three lemmas imply that

$$\begin{aligned}
& (eig_{\min}(\Upsilon) + o_p(1))\|\widehat{\theta}_T^s - \theta_0^s\|^2 + O_p(T^{-1/2})\|\widehat{\theta}_T^s - \theta_0^s\| \\
& \leq \widehat{Q}_T^*(\widehat{\theta}_T) - \widehat{Q}_T^*(\theta_0) \\
& \leq (\sqrt{\widehat{Q}_{0,T}(\widehat{\theta}_T)} + o_p(T^{-1/2}))^2 - \widehat{Q}_T^*(\theta_0) \\
& \leq 2\widehat{Q}_{0,T}(\widehat{\theta}_T) + o_p(T^{-1}) - \widehat{Q}_T^*(\theta_0) \\
& \leq 2\widehat{Q}_T(\widehat{\theta}_T) + o_p(T^{-1}) - \widehat{Q}_T^*(\theta_0) \\
& \leq 2\widehat{Q}_T(\theta_0) + o_p(T^{-1}) - \widehat{Q}_T^*(\theta_0) \\
& \leq 2(\widehat{Q}_T(\theta_0) - \widehat{Q}_{0,T}(\theta_0)) + 2\widehat{Q}_{0,T}(\theta_0) + o_p(T^{-1}) \\
& = O_p(T^{-1}),
\end{aligned} \tag{77}$$

where  $eig_{\min}(\Upsilon)$  is the smallest eigenvalue of  $\Upsilon$ , the first inequality holds by Lemma 6(a)-(b), the second inequality holds by Lemma 5(a), the third inequality holds by the algebraic inequality  $(a + b)^2 \leq 2a^2 + 2b^2$ , the fourth inequality holds because  $\widehat{Q}_{0,T}(\cdot)$  and  $\widehat{Q}_T(\cdot)$  are defined to be exactly the same, both being weighted sums of nonnegative terms, except that the former sums over fewer terms, the fifth inequality holds because  $\widehat{\theta}_T$  is the minimizer of  $\widehat{Q}_T(\cdot)$ , the sixth inequality holds because  $\widehat{Q}_T^*(\theta_0) \geq 0$ , and the equality holds by Lemmas 4 and 5(a)-(b). Let  $\zeta$  be an arbitrary positive number, we next show that we can find a constant  $M_1$  large enough so that

$$\limsup_{T \rightarrow \infty} \Pr \left( T^{1/2}\|\widehat{\theta}_T^s - \theta_0^s\| > M_1 \right) < \zeta. \tag{78}$$

This shows that  $\|\widehat{\theta}_T^s - \theta_0^s\| = O_p(T^{-1/2})$ . To show (78), consider that

$$\begin{aligned}
& \Pr \left( T^{1/2}\|\widehat{\theta}_T^s - \theta_0^s\| > M_1 \right) \\
& \leq \Pr \left( T^{1/2}\|\widehat{\theta}_T^s - \theta_0^s\| > M_1, o_p(1) \geq -eig_{\min}(\Upsilon)/2 \right) + \Pr(o_p(1) < -eig_{\min}(\Upsilon)/2) \\
& \leq \Pr \left( T\|\widehat{\theta}_T^s - \theta_0^s\|^2(eig_{\min}(\Upsilon) + o_p(1)) > \frac{eig_{\min}(\Upsilon)M_1^2}{2}, T^{1/2}\|\widehat{\theta}_T^s - \theta_0^s\| > M_1 \right) + o(1) \\
& \leq \Pr \left( T\|\widehat{\theta}_T^s - \theta_0^s\|^2(eig_{\min}(\Upsilon) + o_p(1)) > \frac{eig_{\min}(\Upsilon)M_1^2}{2}, T^{1/2}\|\widehat{\theta}_T^s - \theta_0^s\| > M_1, O_p(1) \geq -M_2 \right) \\
& \quad + \Pr(O_p(1) < -M_2) + o(1)
\end{aligned}$$

$$\begin{aligned}
&\leq \Pr \left( T \|\widehat{\theta}_T^s - \theta_0^s\|^2 (eig_{\min}(\Upsilon) + o_p(1)) + O_p(T^{1/2}) \|\widehat{\theta}_T^s - \theta_0^s\| > \frac{eig_{\min}(\Upsilon)M_1^2}{2} - M_1M_2 \right) \\
&\quad + \Pr(O_p(1) < -M_2) + o(1) \\
&\leq \Pr \left( O_p(1) > \frac{eig_{\min}(\Upsilon)M_1^2}{2} - M_1M_2 \right) + \Pr(O_p(1) < -M_2) + o(1),
\end{aligned}$$

where the last inequality holds by (77), and the different  $O_p(1)$  terms appearing above are not necessarily the same ones. Fix  $M_2$  at a value such that the limsup of the second term in the last line is less than  $\zeta/2$ . Note that  $\frac{eig_{\min}(\Upsilon)M_1^2}{2} - M_1M_2$  can be made arbitrarily large by increasing  $M_1$  (by Assumption 8(d),  $eig_{\min}(\Upsilon) > 0$ ). Thus, we can choose a  $M_1$  large enough so that the limsup of the first term in the last line is also less than  $\zeta/2$ . Therefore, a large enough  $M_1$  exists such that (78) holds.  $\square$

*Proof of Lemma 4.* Note that

$$\widehat{Q}_T(\theta_T) - \widehat{Q}_{0,T}(\theta_T) = \sum_{g \in \mathcal{G} \setminus \mathcal{G}_0} \mu(g) [\bar{m}_T^u(\theta_T, g)]_-^2 + \sum_{g \in \mathcal{G} \setminus \mathcal{G}_0} \mu(g) [\bar{m}_T^\ell(\theta_T, g)]_+^2.$$

Thus, it suffices to show that

$$\sum_{g \in \mathcal{G} \setminus \mathcal{G}_0} \mu(g) [\bar{m}_T^u(\theta_T, g)]_-^2 = o_p(T^{-1}), \text{ and} \tag{79}$$

$$\sum_{g \in \mathcal{G} \setminus \mathcal{G}_0} \mu(g) [\bar{m}_T^\ell(\theta_T, g)]_+^2 = o_p(T^{-1}). \tag{80}$$

We separate the two cases, one where Assumption 1 is satisfied and the other where Assumption 2 is satisfied.

Case 1: Assumption 1 is satisfied. In this case, arguments for (79) and (80) are analogous. Thus, we give the detailed proof for (79) only. First consider that

$$\sum_{g \in \mathcal{G} \setminus \mathcal{G}_0} \mu(g) [\bar{m}_T^u(\theta_T, g)]_-^2 \leq \sum_{g \in \mathcal{G} \setminus \mathcal{G}_0} \mu(g) [A_T(g) + B_T(g) + C_T(g)]_-^2,$$

where

$$\begin{aligned}
A_T(g) &= \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\delta_{jt}(\pi_t, \lambda_T) - x'_{jt}\beta_T)g(z_{jt}) \\
B_T(g) &= \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\check{\delta}_{jt}(\tilde{s}_t, \lambda_T) - \check{\delta}_{jt}(\pi_t, \lambda_T))g(z_{jt})
\end{aligned}$$

$$C_T(g) = \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\log(s_{jt} + \iota_u/n_t) - \log(\pi_{jt}))g(z_{jt}). \quad (81)$$

The inequality holds because  $\iota_u \geq \underline{\iota}_u$  (Assumption 1(b)). For  $A_T(g)$ , consider that

$$\begin{aligned} A_T(g) &= \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} \xi_{jt}g(z_{jt}) + \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\delta_{jt}(\pi_t, \lambda_T) - \delta_{jt}(\pi_t, \lambda_0))g(z_{jt}) \\ &\quad - \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} x_{jt}(\beta_T - \beta_0)g(z_{jt}). \end{aligned}$$

Equation (64) in the proof of Lemma 3 implies that

$$\sup_{g \in \mathcal{G}} \left| \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} \xi_{jt}g(z_{jt}) \right| = O_p(T^{-1/2}). \quad (82)$$

Also,

$$\begin{aligned} &\sup_{g \in \mathcal{G}} \left| \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\delta_{jt}(\pi_t, \lambda_T) - \delta_{jt}(\pi_t, \lambda_0))g(z_{jt}) \right| \\ &= \sup_{g \in \mathcal{G}} \left| \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} \frac{\partial \delta_{jt}(\pi_t, \tilde{\lambda}_T)}{\partial \lambda'} (\lambda_T - \lambda_0)g(z_{jt}) \right| \\ &\leq \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} \left\| \frac{\partial \delta_{jt}(\pi_t, \tilde{\lambda}_T)}{\partial \lambda'} \right\| \|\lambda_T - \lambda_0\| \\ &= O_p(1)\|\lambda_T - \lambda_0\|, \end{aligned} \quad (83)$$

where the first equality holds by a mean-value expansion for  $\tilde{\lambda}_T$  lying on the line segment connecting  $\lambda_T$  and  $\lambda_0$ , the inequality holds because  $g(z_{jt}) \in (0, 1)$ , the first equality holds by Assumption 8(c) and the condition that  $\lambda_T \rightarrow \lambda_0$  given in the lemma. Moreover,

$$\frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} x_{jt}^s(\beta_T^s - \beta_0^s)g(z_{jt}) \leq \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} \|x_{jt}^s\| \|\beta_T^s - \beta_0^s\| = O_p(1) \|\beta_T^s - \beta_0^s\|, \quad (84)$$

where the equality holds by Assumption 8(c).

Therefore, combining (82), (83), (84) and  $\|\theta_T^s - \theta_0^s\| = O_p(T^{-1/2})$ , we have

$$\sup_{g \in \mathcal{G}} |A_T(g)| = O_p(\log(T)T^{-1/2}). \quad (85)$$

Now consider  $B_T(g)$ . Let  $B_T^0(g) = \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\check{\delta}_{jt}(\tilde{s}_t, \lambda_0) - \check{\delta}_{jt}(\pi_t, \lambda_0))g(z_{jt})$ . Consider that

$$\begin{aligned} \sup_{g \in \mathcal{G}} |B_T(g) - B_T^0(g)| &\leq \sup_{g \in \mathcal{G}} \left| \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\delta_{jt}(\pi_t, \lambda_T) - \delta_{jt}(\pi_t, \lambda_0))g(z_{jt}) \right| \\ &\quad + \sup_{g \in \mathcal{G}} \left| \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\delta_{jt}(\tilde{s}_t, \lambda_T) - \delta_{jt}(\tilde{s}_t, \lambda_0))g(z_{jt}) \right|. \end{aligned}$$

The first summand is less than or equal to  $O_p(1)\|\lambda_T - \lambda_0\|$  by (83). The second summand is also less than or equal to  $O_p(1)\|\lambda_T - \lambda_0\|$  due to the same arguments as those for (83) and the convergence  $\sup_{t=1, \dots, T} \|s_t - \pi_t\|_f \rightarrow_p 0$  implied by (47) and Assumption 4(g). Those combined with  $\|\theta_T^s - \theta_0^s\| = O_p(T^{-1/2})$  shows that:

$$\sup_{g \in \mathcal{G}} |B_T(g) - B_T^0(g)| = O_p(T^{-1/2}). \quad (86)$$

For  $B_T^0(g)$ , consider that

$$\begin{aligned} B_T^0(g) &= \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\check{\delta}_{jt}(\tilde{s}_t, \lambda_0) - \check{\delta}_{jt}(\pi_t, \lambda_0))g(z_{jt}) \\ &= \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} g(z_{jt}) \frac{\partial \check{\delta}_{jt}(\pi_t, \lambda_0)}{\partial \pi'} (\tilde{s}_t - s_t) \\ &\quad + \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} g(z_{jt}) \frac{\partial \check{\delta}_{jt}(\pi_t, \lambda_0)}{\partial \pi'} (s_t - \pi_t) \\ &\quad + \frac{1}{2T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} g(z_{jt}) (\tilde{s}_t - \pi_t)' \frac{\partial^2 \check{\delta}_{jt}(\tilde{\pi}_t, \lambda_0)}{\partial \pi \partial \pi'} (\tilde{s}_t - \pi_t), \end{aligned} \quad (87)$$

where  $\tilde{\pi}_t$  is a point on the line segment connecting  $\tilde{s}_t$  and  $\pi_t$ . For the first summand, note that, by the Cauchy-Schwartz inequality and  $g(z) \in [0, 1]$ , its absolute value is less than or equal to

$$\left( \sup_{t=1, \dots, T} n_t \|\tilde{s}_t - s_t\|_f \right) \left( \frac{1}{T\bar{J}_T \underline{n}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} \left\| \frac{\partial \check{\delta}_{jt}(\pi_t, \lambda_0)}{\partial \pi'} \right\| \right) = O_p(1) \underline{n}_T^{-1} O_p(\sqrt{J_T^{\max}}) = o_p(T^{-1/2}),$$

where the first equality holds by Assumption 4(f),

$$E \left( \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} \left\| \frac{\partial \check{\delta}_{jt}(\pi_t, \lambda_0)}{\partial \pi'} \right\| \right) \leq \sup_{j,t} E \left\| \frac{\partial \check{\delta}_{jt}(\pi_t, \lambda_0)}{\partial \pi'} \right\| = O(\sqrt{J_T^{\max}})$$

(by Assumption 9(b)), and Markov's inequality, and the second equality holds by Assumption 4(g). For the second summand of (87), we can apply Lemma 7 and get

$$\begin{aligned} & E \left[ \sup_{g \in \mathcal{G}} \left( \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} g(z_{jt}) \frac{\partial \check{\delta}_{jt}(\pi_t, \lambda_0)}{\partial \pi'} (s_t - \pi_t) \right)^2 \right] \\ & \leq \frac{C \sum_{t=1}^T J_t^2}{T^2 \bar{J}_T^2} \max_{j,t} E \left( \frac{\partial \check{\delta}_{jt}(\pi_t, \lambda_0)}{\partial \pi'} (s_t - \pi_t) \right)^2 \\ & = \frac{C \sum_{t=1}^T J_t^2}{T^2 \bar{J}_T^2} \max_{j,t} E \left( \frac{\partial \check{\delta}_{jt}(\pi_t, \lambda_0)}{\partial \pi'} \frac{\text{diag}(\pi_t) - \pi_t \pi_t'}{n_t} \frac{\partial \check{\delta}_{jt}(\pi_t, \lambda_0)}{\partial \pi} \right) \\ & \leq \frac{C \sum_{t=1}^T J_t^2}{T^2 \bar{J}_T^2 \underline{n}_T} \max_{j,t} E \left( \left\| \frac{\partial \check{\delta}_{jt}(\pi_t, \lambda_0)}{\partial \pi'} \right\|^2 \right) \\ & = O(\underline{n}_T^{-1} T^{-1} J_T^{\max}) \\ & = o(T^{-1}), \end{aligned}$$

where the first equality holds by  $E[(s_t - \pi_t)(s_t - \pi_t)' | \pi_t, n_t] = \frac{\text{diag}(\pi_t) - \pi_t \pi_t'}{n_t}$  which holds under Assumption 4(b), the second inequality holds because  $\text{diag}(\pi_t) - \pi_t \pi_t'$  is positive semi-definite and its largest eigenvalue does not exceed the highest  $\pi_{jt}$  which does not exceed 1 and because  $n_t \geq \underline{n}_T$  for all  $t = 1, \dots, T$ , the second equality holds by Assumption 9(b) and the boundedness of  $\sum_{t=1}^T J_t^2 / (T \bar{J}_T^2)$ , and the last equality holds by Assumption 4(g). Therefore, the Markov inequality applies and shows that the second summand of (87) is  $o_p(T^{-1/2})$  uniformly over  $g \in \mathcal{G}$ . For the third summand of (87), consider that

$$\begin{aligned} & \sup_{g \in \mathcal{G}} \left| \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} g(z_{jt}) (\tilde{s}_t - \pi_t)' \frac{\partial^2 \check{\delta}_{jt}(\tilde{\pi}_t, \lambda_0)}{\partial \pi \partial \pi'} (\tilde{s}_t - \pi_t) \right| \\ & \leq_{w.p.a.1.} \sup_{j,t} \sup_{\pi: \|\pi - \pi_t\| \leq c} \left\| \frac{\partial^2 \check{\delta}_{jt}(\pi, \lambda_0)}{\partial \pi \partial \pi'} \right\| T^{-1} \sum_{t=1}^T (\tilde{s}_t - \pi_t)' (\tilde{s}_t - \pi_t) \\ & \leq O_p(J_T^{\max}) 2 \left[ T^{-1} \sum_{t=1}^T \|\tilde{s}_t - s_t\|^2 + T^{-1} \sum_{t=1}^T \|s_t - \pi_t\|^2 \right] \\ & = O_p(J_T^{\max}) O_p(\underline{n}_T^{-1}) + O_p(J_T^{\max}) T^{-1} \sum_{t=1}^T \|s_t - \pi_t\|^2 \end{aligned}$$

$$\begin{aligned}
&= O_p(J_T^{\max})O_p(\underline{n}_T^{-1}) + O_p(J_T^{\max})O_p(\underline{n}_T^{-1}) \\
&= o_p(T^{-1/2}),
\end{aligned}$$

where the first inequality holds because  $\sup_t \|\tilde{s}_t - \pi_t\| \leq \sup_t \|\tilde{s}_t - \pi_t\|_f \leq c$  w.p.a.1. by Assumption 4(f,g) and equation (47) and also because  $g(z) \in [0, 1]$ , the second inequality holds by Assumption 9(b), the first equality holds by Assumption 4(f), the second equality holds by Markov's inequality and  $E\|s_t - \pi_t\|^2 = E \sum_{j=1}^{J_t} \pi_{jt}(1 - \pi_{jt})/n_t \leq \underline{n}_T^{-1}$ , and the last equality holds by Assumption 4(g). Combining the arguments for all the three summands in (87), we have

$$\sup_{g \in \mathcal{G}} |B_T^0(g)| = o_p(T^{-1/2}). \quad (88)$$

This and (86) together imply that

$$\sup_{g \in \mathcal{G}} |B_T(g)| = o_p(T^{-1/2}). \quad (89)$$

Next consider  $C_T(g)$ . Using the moment bound derived in (65) in the proof of Theorem 3 and the Markov inequality, we can derive

$$\sup_{g \in \mathcal{G}} |C_T(g) - E[C_T(g)]| = O_p\left(\frac{\log \bar{n}_T}{T^{1/2}}\right) = O_p\left(\frac{\log T}{T^{1/2}}\right), \quad (90)$$

where the second equality holds by  $\bar{n}_T T^{-2} \rightarrow_p 0$  (Assumption 8(e)).

Let  $r_T(g)$  denote  $A_T(g) + B_T(g) + C_T(g) - E[C_T(g)]$ . Then  $\bar{m}_T^u(\theta_T, g) \geq r_T(g) + E[C_T(g)]$ . And by equations (85), (89), and (90), we have

$$\sup_{g \in \mathcal{G}} |r_T(g)| = O_p(T^{-1/2} \log T). \quad (91)$$

For a sequence  $c_T$  such that  $T^{-1/2} \log T = o(c_T)$ , consider:

$$\begin{aligned}
\sum_{g \in \mathcal{G} \setminus \mathcal{G}_0: E[C_T(g)] > c_T} \mu(g) [\bar{m}_T^u(\theta_T, g)]_-^2 &\leq \sum_{g \in \mathcal{G} \setminus \mathcal{G}_0: E[C_T(g)] > c_T} \mu(g) [r_T(g) + c_T]_-^2 \\
&\leq \sup_{g \in \mathcal{G}} [r_T(g) + c_T]_-^2 \\
&= [o_p(c_T) + c_T]_-^2 \\
&=_{w.p.a.1} 0,
\end{aligned} \quad (92)$$

where the first inequality holds because  $[\cdot]_-^2$  is nonincreasing, the second inequality holds

because  $\mu(g)$  is a probability mass function, the first equality holds by (91). Thus, the expression  $\sum_{g \in \mathcal{G} \setminus \mathcal{G}_0: E[C_T(g)] > c_T} \mu(g) [\bar{m}_T^u(\theta_T, g)]_-^2$  converges in probability to zero at arbitrary rate. Further restrict  $c_T$  so that  $c_T = o((\log T)^{-2/\eta})$ . This is possible because for any finite  $\eta > 0$ ,  $\log(T)^{1+2/\eta} = o(T^{1/2})$ . Also consider

$$\begin{aligned}
\sum_{g \in \mathcal{G} \setminus \mathcal{G}_0: E[C_T(g)] \leq c_T} \mu(g) [\bar{m}_T^u(\theta_T, g)]_-^2 &\leq \sum_{g \in \mathcal{G} \setminus \mathcal{G}_0: E[C_T(g)] \leq c_T} \mu(g) [r_T(g)]_-^2 \\
&\leq \sup_{g \in \mathcal{G}} |r_T(g)|^2 \sum_{g \in \mathcal{G} \setminus \mathcal{G}_0: E[C_T(g)] \leq c_T} \mu(g) \\
&= O_p(T^{-1}(\log T)^2) c_T^\eta \\
&= o_p(T^{-1}),
\end{aligned} \tag{93}$$

where the first inequality holds because  $\bar{m}_T^u(\theta_T, g) = r_T(g) + E[C_T(g)]$  and

$$E[C_T(g)] = (T \bar{J}_T)^{-1} \sum_{t=1}^T \sum_{j=1}^{J_t} E[\log(s_{jt} + \underline{l}_u/n_t) - \log(\pi_{jt})g(z_{jt})] \geq 0$$

by the definition of  $\underline{l}_u$ , and the first equality holds by the first part of Assumption 9(a). Therefore, combining (92) and (93), we have

$$\sum_{g \in \mathcal{G} \setminus \mathcal{G}_0} \mu(g) [\bar{m}_T^u(\theta_T, g)]_-^2 = o_p(T^{-1}). \tag{94}$$

Case 2: Assumption 2 is satisfied. We prove (79) first. Observe that

$$\sum_{g \in \mathcal{G} \setminus \mathcal{G}_0} \mu(g) [\bar{m}_T^u(\theta_T, g)]_-^2 = \sum_{g \in \mathcal{G} \setminus \mathcal{G}_0} \mu(g) [A_T(g) + \Delta_T(g) + S_T(g)]_-^2,$$

where

$$\begin{aligned}
A_T(g) &= \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\delta_{jt}(\pi_t, \lambda_T) - x'_{jt} \beta_T) g(z_{jt}) \\
\Delta_T(g) &= \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\delta_{jt}(\tilde{s}_t, \lambda_T) - \delta_{jt}(\pi_t, \lambda_T)) g(z_{jt}) \\
S_T(g) &= \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\log(s_{jt} + \underline{l}_u/n_t) - \log(\tilde{s}_{jt})) g(z_{jt}).
\end{aligned} \tag{95}$$

The same arguments showing (85) in Case 1 still applies in Case 2 since neither Assumption

1 or Assumption 2 is involved. Thus, (85) holds. For  $\Delta_T(g)$ , the same arguments as those for (86) shows that

$$\sup_{g \in \mathcal{G} \setminus \mathcal{G}_0} \left| \Delta_T(g) - \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\delta_{jt}(\tilde{s}_t, \lambda_0) - \delta_{jt}(\pi_t, \lambda_0)) g(z_{jt}) \right| = O_p(1) \|\lambda_T - \lambda_0\|. \quad (96)$$

Equation (69) in Case 2 of the proof of Theorem 3 shows that

$$\begin{aligned} & E \sup_{g \in \mathcal{G}} \left| \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\delta_{jt}(\tilde{s}_t, \lambda_0) - \delta_{jt}(\pi_t, \lambda_0)) g(z_{jt}) - E[(\delta_{jt}(\tilde{s}_t, \lambda_0) - \delta_{jt}(\pi_t, \lambda_0)) g(z_{jt})] \right|^2 \\ &= O\left(\frac{\log(\bar{n}_T)^2}{T}\right). \end{aligned}$$

Thus, by the Markov inequality,

$$\begin{aligned} & \sup_{g \in \mathcal{G}} \left| \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\delta_{jt}(\tilde{s}_t, \lambda_0) - \delta_{jt}(\pi_t, \lambda_0)) g(z_{jt}) - E[(\delta_{jt}(\tilde{s}_t, \lambda_0) - \delta_{jt}(\pi_t, \lambda_0)) g(z_{jt})] \right| \\ &= O_p\left(\frac{\log(\bar{n}_T)}{T^{1/2}}\right) \\ &= O_p(\log(T) T^{-1/2}). \end{aligned} \quad (97)$$

where the second equality holds by  $\bar{n}_T T^{-2} \rightarrow_p 0$  (Assumption 8(e)). By Assumption 2(b),  $E[(\delta_{jt}(\tilde{s}_t, \lambda_0) - \delta_{jt}(\pi_t, \lambda_0)) g(z_{jt})] \geq 0$ . This combined with (96), (97), and  $\|\hat{\theta}_T^s - \theta_0^s\| = O_p(T^{-1/2})$  implies that

$$\inf_{g \in \mathcal{G}} \Delta_T(g) \geq O_p((\log(T) T^{-1/2})). \quad (98)$$

For  $S_T(g)$ , note that

$$\begin{aligned} S_T(g) &\geq \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (s_{jt} + \iota_u/n_t)^{-1} ((s_{jt} + \iota_u/n_t) - (\tilde{s}_{jt})) g(z_{jt}) \\ &= (\iota_u - 1) \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} \frac{g(z_{jt})}{n_t s_{jt} + \iota_u}. \end{aligned}$$

Applying Lemma 7 and using the fact that  $E[(n_t s_{jt} + \iota_u)^{-2}] \leq \iota_u^{-2}$ , we have

$$E \sup_{g \in \mathcal{G}} \left( \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} \frac{g(z_{jt})}{n_t s_{jt} + \iota_u} - E \left[ \frac{g(z_{jt})}{n_t s_{jt} + \iota_u} \right] \right)^2 = O(T^{-1}).$$



Then by the Markov inequality we have

$$\sup_{g \in \mathcal{G}} \left| \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} \frac{g(z_{jt})}{n_t s_{jt} + \iota_u} - E \left[ \frac{g(z_{jt})}{n_t s_{jt} + \iota_u} \right] \right| = O_p(T^{-1/2}).$$

Thus we have

$$S_T(g) \geq O_p(T^{-1/2}) + \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} E \left[ \frac{g(z_{jt})}{n_t s_{jt} + \iota_u} \right]. \quad (99)$$

Using (85), (98), (99), and the third part of Assumption 9(a), we can apply the same arguments as those for (94) (from (92) to (94)) to conclude that (79) holds.

Finally we prove (80) for Case 2. Note that

$$\sum_{g \in \mathcal{G} \setminus \mathcal{G}_0} \mu(g) [\bar{m}_T^\ell(\theta_T, g)]_+^2 \leq \sum_{g \in \mathcal{G} \setminus \mathcal{G}_0} \mu(g) [A_T(g) + B_T(g) + C_T^\ell(g)]_+^2,$$

where

$$\begin{aligned} A_T(g) &= \frac{1}{T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\delta_{jt}(\pi_t, \lambda_T) - x'_{jt} \beta_T) g(z_{jt}) \\ B_T(g) &= \frac{1}{T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\check{\delta}_{jt}(\tilde{s}_t, \lambda_T) - \check{\delta}_{jt}(\pi_t, \lambda_T)) g(z_{jt}) \\ C_T^\ell(g) &= \frac{1}{T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\log(s_{jt} + \bar{\iota}_\ell/n_t) - \log(\pi_{jt})) g(z_{jt}). \end{aligned} \quad (100)$$

The same arguments showing (85) in Case 1 still applies in Case 2 since neither Assumption 1 or Assumption 2 is involved. Thus, (85) holds. For  $B_T(g)$ , the same arguments for (86) in Case 1 still applies here as well. Thus, (86) holds, and we only need to study  $B_T^0(g)$  to understand the behavior of  $B_T(g)$ . Note that

$$\begin{aligned} B_T^0(g) &= \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\delta_{jt}(\tilde{s}_t, \lambda_0) - \delta_{jt}(\pi_t, \lambda_0)) g(z_{jt}) - E[(\delta_{jt}(\tilde{s}_t, \lambda_0) - \delta_{jt}(\pi_t, \lambda_0)) g(z_{jt})] \\ &\quad - \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\log(\tilde{s}_t) - \log(\pi_t)) g(z_{jt}) - E[(\log(\tilde{s}_t) - \log(\pi_t)) g(z_{jt})] \\ &\quad + \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} E[(\check{\delta}_{jt}(\tilde{s}_t, \lambda_0) - \check{\delta}_{jt}(\pi_t, \lambda_0)) g(z_{jt})]. \end{aligned}$$

Equation (97) shows that the first summand is  $O_p(\log(T)T^{-1/2})$  uniformly over  $g \in \mathcal{G}$ , Equ-

tion (65) and Markov inequality combined show that the second summand is  $O_p(\log(T)T^{-1/2})$  uniformly over  $g \in \mathcal{G}$ . The third summand is non-positive by Assumption 2(a). Therefore

$$\sup_{g \in \mathcal{G}} B_T(g) \leq O_p(\log(T)T^{-1/2}). \quad (101)$$

The same arguments as those for the second summand above shows that  $\sup_{g \in \mathcal{G}} |C_T^\ell(g) - E[C_T^\ell(g)]| = O_p(\log(T)T^{-1/2})$ . Using this, (85), (101), and the second part of Assumption 9(a), we can apply similar arguments as those for (94) (from (92) to (94)) to conclude that (80) holds.  $\square$

*Proof of Lemma 5.* (a) By equation (60) in the proof of Theorem 1, we have

$$\begin{aligned} & \sup_{\theta \in \Theta: \theta^s \in B_c(\theta_0^s)} \left| \sqrt{\widehat{Q}_{0,T}(\theta)} - \sqrt{\widehat{Q}_T^*(\theta)} \right| \\ & \leq \sqrt{\sup_{\theta \in \Theta: \theta^s \in B_c(\theta_0^s)} \sup_{g \in \mathcal{G}_0} |\bar{m}_T^u(\theta, g) - \bar{m}_T(\theta, g)|^2 + \sup_{\theta \in \Theta: \theta^s \in B_c(\theta_0^s)} \sup_{g \in \mathcal{G}_0} |\bar{m}_T^\ell(\theta, g) - \bar{m}_T(\theta, g)|^2}. \end{aligned} \quad (102)$$

Now note that

$$\begin{aligned} \bar{m}_T^u(\theta, g) - \bar{m}_T(\theta, g) &= (T\bar{J}_T)^{-1} \sum_{t=1}^T \sum_{j=1}^{J_t} (\hat{\delta}_{jt}^u(s_t, \lambda) - \delta_{jt}(\pi_t, \lambda))g(z_{jt}) \\ &= (T\bar{J}_T)^{-1} \sum_{t=1}^T \sum_{j=1}^{J_t} (\log(s_{jt} + \iota_u/n_t) - \log(s_{jt}))g(z_{jt}) \\ &\quad + (T\bar{J}_T)^{-1} \sum_{t=1}^T \sum_{j=1}^{J_t} \delta_{jt}(\tilde{s}_t, \lambda) - \delta_{jt}(\pi_t, \lambda))g(z_{jt}). \end{aligned} \quad (103)$$

For the first summand, consider that

$$\begin{aligned} \sup_{g \in \mathcal{G}_0} \left| (T\bar{J}_T)^{-1} \sum_{t=1}^T \sum_{j=1}^{J_t} (\log(s_{jt} + \iota_u/n_t) - \log(s_{jt}))g(z_{jt}) \right| &\leq (T\bar{J}_T)^{-1} \iota_u \sum_{t=1}^T \sum_{j=1}^{J_t} (s_{jt} + \tilde{\iota}/n_t)^{-1} n_t^{-1} \\ &\leq \underline{n}_T^{-1} \iota_u \sup_{j,t: z_{jt} \in \mathcal{Z}_0} s_{jt}^{-1} \\ &= O_p(\underline{n}_T^{-1}) = o_p(T^{-1/2}), \end{aligned} \quad (104)$$

where the first inequality holds with  $\tilde{\iota} \in [0, \iota_u]$  by mean-value expansion and  $|g(z_{jt})| \leq 1$ , the second inequality holds by the definition of  $\mathcal{G}_0$ , the first equality holds because  $s_{jt}$  is

bounded away from zero by Assumptions 3 and equation (46), and the last equality holds by Assumption 8(e). For the second summand in (103), we can apply the same arguments as those for (89) to show that this second summand is  $o_p(T^{-1/2})$  with the following adjustment: (1) Replace  $\mathcal{G}$  by  $\mathcal{G}_0$ , (2) replace  $\check{\delta}_{jt}(\cdot, \cdot)$  with  $\delta_{jt}(\cdot, \cdot)$  and (2) replace the Case 1 version of Assumption 9(b) by the Case 2 version. Therefore, we have

$$\sup_{\theta \in \Theta: \theta^s \in B_c(\theta_0^s)} \sup_{g \in \mathcal{G}_0} |\bar{m}_T^u(\theta, g) - \bar{m}_T(\theta, g)| = o_p(T^{-1/2}).$$

Analogous arguments can be used to show that  $\sup_{\theta \in B_c(\theta_0)} \sup_{g \in \mathcal{G}_0} |\bar{m}_T^\ell(\theta, g) - \bar{m}_T(\theta, g)| = o_p(T^{-1/2})$ . That concludes the proof.  $\square$

(b) Recall that  $\hat{Q}_T^*(\theta_0) = \sum_{g \in \mathcal{G}_0} \mu(g)(\bar{m}_T(\theta_0, g))^2$ , and note that

$$\bar{m}_T(\theta_0, g) = \frac{1}{T\bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} (\delta_{jt}(\pi_t, \lambda_0) - x'_{jt}\beta_0)g(z_{jt}) = \frac{1}{T} \sum_{t=1}^T \sum_{j=1}^{J_t} \xi_{jt}g(z_{jt}).$$

Then by equation(64) in the proof of Theorem 3, we have

$$\sup_{g \in \mathcal{G}_0} |\bar{m}_T(\theta_0, g)| = O_p(T^{-1/2}). \quad (105)$$

This implies part (b).

*Proof of Lemma 6.* (a) First consider that, for  $g \in \mathcal{G}_0$ ,

$$\begin{aligned} & \bar{m}_T(\theta_T, g) - \bar{m}_T(\theta_0, g) \\ &= (T\bar{J}_T)^{-1} \sum_{t=1}^T \sum_{j=1}^{J_t} g(z_{jt}) [\delta_{jt}(\pi_t, \lambda_T) - \delta_{jt}(\pi_t, \lambda_0) + x_{jt}^s(\beta_T^s - \beta_0^s)] \\ &= (T\bar{J}_T)^{-1} \sum_{t=1}^T \sum_{j=1}^{J_t} g(z_{jt}) \partial m_{jt}(\lambda_0)'(\theta_T^s - \theta_0^s) + \\ & \quad (T\bar{J}_T)^{-1} \sum_{t=1}^T \sum_{j=1}^{J_t} g(z_{jt}) (\lambda_T - \lambda_0)' \frac{\partial^2 \delta_{jt}(\pi_t, \tilde{\lambda})}{\partial \lambda \partial \lambda'} (\lambda_T - \lambda_0) / 2 \\ &= \hat{\Gamma}_T(g)'(\theta_T^s - \theta_0^s) + (\lambda_T - \lambda_0)' D_T(g) (\lambda_T - \lambda_0), \end{aligned}$$

where  $\tilde{\lambda}$  is a point on the line segment connecting  $\lambda_T$  and  $\lambda_0$ , and

$$D_T(g) = (2T\bar{J}_T)^{-1} \sum_{t=1}^T \sum_{j=1}^{J_t} g(z_{jt}) \frac{\partial^2 \delta_{jt}(\pi_t, \tilde{\lambda})}{\partial \lambda \partial \lambda'}.$$

Thus, we have

$$\begin{aligned}
& \widehat{Q}_T^*(\theta_T) - \widehat{Q}_T^*(\theta_0) \\
&= \sum_{g \in \mathcal{G}_0} \mu(g) (\bar{m}_T(\theta_T, g) - \bar{m}_T(\theta_0, g))^2 + 2 \sum_{g \in \mathcal{G}_0} \mu(g) \bar{m}_T(\theta_0, g) (\bar{m}_T(\theta_T, g) - \bar{m}_T(\theta_0, g)) \quad (106) \\
&= (\theta_T^s - \theta_0^s) \sum_{g \in \mathcal{G}_0} \mu(g) \widehat{\Gamma}_T(g) \widehat{\Gamma}_T(g)' (\theta_T^s - \theta_0^s) \\
&\quad + 2 \sum_{g \in \mathcal{G}_0} \mu(g) (\lambda_T - \lambda_0)' D_T(g) (\lambda_T - \lambda_0) \widehat{\Gamma}_T(g)' (\theta_T^s - \theta_0^s) \\
&\quad + \sum_{g \in \mathcal{G}_0} \mu(g) \{(\lambda_T - \lambda_0)' D_T(g) (\lambda_T - \lambda_0)\}^2 \\
&\quad + 2 \sum_{g \in \mathcal{G}_0} \mu(g) \bar{m}_T(\theta_0, g) \widehat{\Gamma}_T(g)' (\theta_T^s - \theta_0^s) \\
&\quad + 2 \sum_{g \in \mathcal{G}_0} \mu(g) \bar{m}_T(\theta_0, g) (\lambda_T - \lambda_0)' D_T(g) (\lambda_T - \lambda_0). \quad (107)
\end{aligned}$$

Since  $\tilde{\lambda} \in B_c(\lambda_0)$  whenever  $\lambda_T \in B_c(\lambda_0)$  (which holds with probability approaching one because  $\|\lambda_T - \lambda_0\| \rightarrow_p 0$ ), we have for any  $g \in \mathcal{G}_0$ ,

$$\sup_{g \in \mathcal{G}_0} \|D_T(g)\| \leq_{w.p.a.1} \frac{1}{T \bar{J}_T} \sum_{t=1}^T \sum_{j=1}^{J_t} \sup_{\lambda: \|\lambda - \lambda_0\| \leq c} \left\| \frac{\partial^2 \delta_{jt}(\pi_t, \lambda)}{\partial \lambda \partial \lambda'} \right\| = O_p(1), \quad (108)$$

where the first inequality holds because  $0 \leq g(z) \leq 1$ , and the equality holds by Markov's inequality and Assumption 8(c). This combined with  $\|\theta_T^s - \theta_0^s\| = o_p(1)$  implies that

$$\sum_{g \in \mathcal{G}_0} \mu(g) \{(\lambda_T - \lambda_0)' D_T(g) (\lambda_T - \lambda_0)\}^2 \leq \sup_{g \in \mathcal{G}_0} \|D_T(g)\|^2 \|\theta_T^s - \theta_0^s\|^4 = o_p(1) \|\theta_T^s - \theta_0^s\|^2.$$

Also, using Assumption 8(c) and the same arguments as those for (108), we can show that  $\sup_{g \in \mathcal{G}_0} \|\widehat{\Gamma}_T(g)\| = O_p(1)$ . This combined with (108) and  $\|\theta_T^s - \theta_0^s\| = o_p(1)$  implies that

$$\begin{aligned}
\left| \sum_{g \in \mathcal{G}_0} \mu(g) (\lambda_T - \lambda_0)' D_T(g) (\lambda_T - \lambda_0) \widehat{\Gamma}_T(g)' (\theta_T^s - \theta_0^s) \right| &\leq \|\theta_T^s - \theta_0^s\|^3 \sup_{g \in \mathcal{G}_0} \|D_T(g)\| \|\widehat{\Gamma}_T(g)\| \\
&= o_p(1) \|\theta_T^s - \theta_0^s\|^2.
\end{aligned}$$

Next apply Lemma 7 with  $w_{jt} = \xi_{jt}$  and we get

$$E \sup_{g \in \mathcal{G}_0} (\bar{m}_T(\theta_0, g))^2 = E \sup_{g \in \mathcal{G}_0} \left( (T \bar{J}_T)^{-1} \sum_{t=1}^T \sum_{j=1}^{J_t} \xi_{jt} g(z_{jt}) \right)^2$$

$$\begin{aligned}
&\leq \frac{C \sum_{t=1}^T J_t^2}{T^2 \bar{J}_T^2} \sup_{j,t} E[\xi_{jt}^2 1(z_{jt} \in \mathcal{Z}_0)] \\
&= O(T^{-1}),
\end{aligned}$$

where the second equality holds by Assumption 4(e) and the boundedness of  $T^{-1} \sum_{t=1}^T J_t^2 / \bar{J}_T^2$ . Therefore,

$$\sup_{g \in \mathcal{G}_0} |\bar{m}_T(\theta_0, g)| = O_p(T^{-1/2}). \quad (109)$$

This combined with (108) implies that

$$\sum_{g \in \mathcal{G}_0} \mu(g) \bar{m}_T(\theta_0, g) (\lambda_T - \lambda_0)' D_T(g) (\lambda_T - \lambda_0) = O_p(T^{-1/2}) \|\theta_T^s - \theta_0^s\|^2.$$

Therefore, part (a) holds.

(b) Apply Lemma 7 with  $w_{jt}$  being an element of the random vector  $\partial m_{jt}(\lambda_0)$ , do so for every element of  $\partial m_{jt}(\lambda_0)$ , and we get

$$E \sup_{g \in \mathcal{G}_0} \left\| \hat{\Gamma}_T(g) - \Gamma_T(g) \right\|^2 \leq \frac{C \sum_{t=1}^T J_t^2}{T^2 \bar{J}_T^2} \sup_{j,t} E[\|\partial m_{jt}(\lambda_0)\|^2 1(z_{jt} \in \mathcal{Z}_0)] = O(T^{-1}).$$

The equality is implied by Assumptions 8(c) and the boundedness of  $J_t$ . Thus, we have

$$\sup_{g \in \mathcal{G}_0} \left\| \hat{\Gamma}_T(g) - \Gamma_T(g) \right\| = O_p(T^{-1/2}). \quad (110)$$

Assumption 8(c) also implies that

$$\begin{aligned}
\sup_{g \in \mathcal{G}_0} \|\Gamma_T(g)\| &\leq \sup_{g \in \mathcal{G}_0} (T \bar{J}_T)^{-1} \sum_{t=1}^T \sum_{j=1}^{J_t} E[\|\partial m_{jt}(\lambda_0) g(z_{jt})\|] \\
&\leq \sup_{g \in \mathcal{G}_0} \sup_{j,t} E[\|\partial m_{jt}(\lambda_0)\| 1(z_{jt} \in \mathcal{Z}_0)] \\
&= O(1).
\end{aligned} \quad (111)$$

This and (110) together imply that

$$\hat{\Upsilon}_T = \sum_{g \in \mathcal{G}_0} \mu(g) \hat{\Gamma}_T(g) \hat{\Gamma}_T(g)' = o_p(1) + \sum_{g \in \mathcal{G}_0} \mu(g) \Gamma_T(g) \Gamma_T(g)' \rightarrow_p \Upsilon,$$

where the convergence holds by Assumption 8(d).

For  $W_n$ , first consider the derivation

$$\begin{aligned} \left| T^{1/2} \sum_{g \in \mathcal{G}_0} \mu(g) \bar{m}_T(\theta_0, g) (\hat{\Gamma}_T(g) - \Gamma_T(g)) \right| &\leq \sup_{g \in \mathcal{G}_0} |\bar{m}_T(\theta_0, g)| \sup_{g \in \mathcal{G}_0} T^{1/2} \|\hat{\Gamma}_T(g) - \Gamma_T(g)\| \\ &= O_p(T^{-1/2}) = o_p(1), \end{aligned}$$

by equations (109) and (110). Thus,

$$\begin{aligned} T^{1/2} W_n &= o_p(1) + T^{1/2} \sum_{g \in \mathcal{G}_0} \mu(g) \bar{m}_T(\theta_0, g) \Gamma_T(g) \\ &= o_p(1) + T^{-1/2} \sum_{t=1}^T v_t, \end{aligned}$$

where  $v_t = \bar{J}_T^{-1} \sum_{j=1}^{J_t} \left[ \xi_{jt} \left( \sum_{g \in \mathcal{G}_0} \mu(g) g(z_{jt}) \Gamma_T(g) \right) \right]$ . Observe that  $\{v_t\}_{t=1}^T$  is independent across  $t$  by Assumption 4(d). Also consider the derivation:

$$\begin{aligned} E[v_t] &= E \sum_{j=1}^{J_t} \left[ E[\xi_{jt} | z_{jt}] \left( \sum_{g \in \mathcal{G}_0} \mu(g) g(z_{jt}) \Gamma_T(g) \right) \right] = 0 \\ &T^{-1} \sum_{t=1}^T E[v_t v_t'] \\ &= T^{-1} \sum_{t=1}^T \sum_{g, g^* \in \mathcal{G}_0} \text{Cov} \left( \bar{J}_T^{-1} \sum_{j=1}^{J_t} \xi_{jt} g(z_{jt}), \bar{J}_T^{-1} \sum_{j=1}^{J_t} \xi_{jt} g^*(z_{jt}) \right) \Gamma_T(g) \Gamma_T(g)' \mu(g) \mu(g^*) \rightarrow V, \end{aligned}$$

where the second equality in the first lines holds by Assumptions 4(c), and the convergence holds by 8(d). Also for the  $c$  in Assumption 4(e),

$$\begin{aligned} E(\|v_t\|^{2+c}) &\leq \sup_{j,t} E|\xi_{jt}|^{2+c} \sup_{g \in \mathcal{G}_0} \|\Gamma_T(g)\|^{2+c} \\ &= O(1), \end{aligned}$$

by Assumptions 4(d) and equation (111) above. Therefore, we can apply the Lindeberg central limit theorem and conclude  $T^{-1/2} \sum_{t=1}^T v_t \rightarrow_d N(0, V)$ . Therefore,

$$T^{1/2} W_n \rightarrow_d N(0, V).$$

□

## C.4 Auxiliary Lemmas

In this subsection, we present two auxiliary lemmas. Lemma 7 establishes a maximal inequality for certain empirical processes indexed by  $g$  in a subset of  $\mathcal{G}$ . This is used above to establish the convergence rates of several empirical processes. Lemma 8 establishes a concentration inequality for the  $L_2$  distance between a multinomial random vector and its expectation. This is used to derive a tighter tail bound for  $\|s_t - \pi_t\|$  than that implied by Chernoff's inequality when  $J_t$  is large.

**Lemma 7.** *Let  $\{z_{jt} : j = 1, \dots, J_t, t = 1, \dots, T\}_{T \geq 1}$  be an array of random vectors. Let  $\mathcal{G}$  be the set of instrumental functions defined in (16). Let  $\mathcal{Z}^*$  be a subset of  $\text{supp}(z_{jt})$  and let  $\mathcal{G}^*$  be a subset of  $\mathcal{G}$  for which  $g(z) = 0$  for all  $z \notin \mathcal{Z}^*$  for all  $g \in \mathcal{G}^*$ . Let  $\{w_{jt} : j = 1, \dots, J_t, t = 1, \dots, T\}_{T \geq 1}$  be an array of random variables such that  $E[w_{jt}^2 \mathbf{1}(z_{jt} \in \mathcal{Z}^*)] \leq M_T$  for all  $j, t$  for some  $M_T < \infty$ . Let  $w_t = (w_{1t}, \dots, w_{J_t t})'$  and  $z_t = (z_{1t}, \dots, z_{J_t t})'$ . Suppose that  $(w_t, z_t)$  is independent across  $t$ . Then*

$$E \sup_{g \in \mathcal{G}^*} \left( \sum_{t=1}^T \sum_{j=1}^{J_t} (w_{jt} g(z_{jt}) - E[w_{jt} g(z_{jt})]) \right)^2 \leq C M_T \sum_{t=1}^T J_t^2,$$

for some constant  $C > 0$ .

*Proof.* Recall that  $J_T^{\max} = \max_{t=1, \dots, T} J_t$ . First observe that  $\sum_{j=1}^{J_t} w_{jt} g(z_{jt})$  can be written as  $f_t(g) := \sum_{j=1}^{J_t^{\max}} w_{jt} \mathbf{1}(j \leq J_t) g(z_{jt})$ . Observe that the triangular array of random processes  $\{g(z_{jt}) : g \in \mathcal{G}^* : t = 1, \dots, T\}_{T \geq 1}$  is manageable with respect to the envelope  $\mathbf{1}_T$  for all  $j$  in the sense of Pollard (1990) because  $\mathcal{G}$  is the collection of indicator functions for a Vapnik-Cervonenkis class of sets. Then by parts (a) and (c) of Lemma E1 in Andrews and Shi (2013), we have that the triangular array  $\{f_t(g) : g \in \mathcal{G}^* : t = 1, \dots, T; T \geq 1\}$  is manageable with respect to the envelope function  $F_T = (F_{T1}, \dots, F_{TT})$  where  $F_{Tt} = \sum_{j=1}^{J_t^{\max}} \mathbf{1}(j \leq J_t, z_{jt} \in \mathcal{Z}^*) |w_{jt}| \equiv \sum_{j=1}^{J_t} |w_{jt}| \mathbf{1}(z_{jt} \in \mathcal{Z}^*)$ . Therefore, by the maximal inequality (7.10) in Pollard (1990), we have, for some constant  $C > 0$ ,

$$\begin{aligned} E \sup_{g \in \mathcal{G}^*} \left| \sum_{t=1}^T \sum_{j=1}^{J_t} (w_{jt} g(z_{jt}) - E[w_{jt} g(z_{jt})]) \right|^2 &\leq C \sum_{t=1}^T E[(F_{Tt})^2] \\ &\leq C \sum_{t=1}^T J_t \sum_{j=1}^{J_t} E[w_{jt}^2 \mathbf{1}(z_{jt} \in \mathcal{Z}^*)] \\ &\leq C M_T \sum_{t=1}^T J_t^2, \end{aligned} \tag{112}$$

proving the lemma. □

The following lemma presents a concentration inequality for the  $L_2$  distance from the mean for multinomial random vectors. The tail bound presented here does not depend on the length of the multinomial random vector, and thus can be applied for multinomial distributions with an arbitrarily large number of categories. The proof of the lemma uses Poissonization, a technique that Devroye (1983) employs in his Lemma 3 to derive a concentration inequality for the  $L_1$  distance from the mean for multinomial random vectors. Devroye's bound applies when the number of categories is smaller than a scalar multiple of the sample size.

**Lemma 8.** *Let  $(X_1, \dots, X_J)$  be a multinomial  $(n, p_1, \dots, p_J)$  random vector, where  $p_1, \dots, p_J$  are non-negative numbers that sum up to 1 and  $n$  is a positive integer. Then, for all  $\varepsilon > 0$ ,*

$$\Pr \left( \sum_{j=1}^J (X_j - np_j)^2 > n^2 \varepsilon^2 \right) \leq \frac{128(3n^2 + n)}{n^4 \varepsilon^4}. \quad (113)$$

*Proof.* Let  $U_1, U_2, \dots$  be a sequence of independent and identically distributed  $\{1, \dots, J\}$ -valued variables with probability mass given by  $P(U_1 = j) = p_j, 1 \leq j \leq J$ . Let  $N$  be a Poisson( $n$ ) random variable independent of  $\{U_1, U_2, \dots\}$ . Let  $X_j$  be the number of occurrences of the value  $j$  among  $U_1, \dots, U_n$ , and let  $\tilde{X}_j$  be the number of occurrences of the value  $j$  among  $U_1, \dots, U_N$ . It is clear that  $X_1, \dots, X_J$  is a multinomial  $(n, p_1, \dots, p_J)$  random vector, and that  $\tilde{X}_1, \dots, \tilde{X}_J$  are independent Poisson random variables with means  $np_1, \dots, np_J$ . By the inequality  $(a + b)^2 \leq 2a^2 + 2b^2$ , we have

$$\sum_{j=1}^J (X_j - np_j)^2 \leq 2 \sum_{j=1}^J (X_j - \tilde{X}_j)^2 + 2 \sum_{j=1}^J (\tilde{X}_j - np_j)^2. \quad (114)$$

Thus,

$$\begin{aligned} & \Pr \left( \sum_{j=1}^J (X_j - np_j)^2 > n^2 \varepsilon^2 \right) \\ & \leq \Pr \left( \sum_{j=1}^J (X_j - \tilde{X}_j)^2 > n^2 \varepsilon^2 / 4 \right) + \Pr \left( \sum_{j=1}^J (\tilde{X}_j - np_j)^2 > n^2 \varepsilon^2 / 4 \right). \end{aligned} \quad (115)$$



For  $\sum_{j=1}^J (X_j - \tilde{X}_j)^2$  consider the derivation:

$$\begin{aligned}
(X_j - \tilde{X}_j)^2 &= \left( \sum_{i=n+1}^N 1\{U_i = j\} \right)^2 1\{n < N\} + \left( \sum_{i=N+1}^n 1\{U_i = j\} \right)^2 1\{n > N\} \\
&= \left( \sum_{i=n+1}^N 1\{U_i = j\} + 2 \sum_{i \neq i', i, i' = n+1, \dots, N} 1\{U_i = j\} 1\{U_{i'} = j\} \right) 1\{n < N\} \\
&\quad + \left( \sum_{i=N+1}^n 1\{U_i = j\} + 2 \sum_{i \neq i', i, i' = N+1, \dots, n} 1\{U_i = j\} 1\{U_{i'} = j\} \right) 1\{n > N\}. \quad (116)
\end{aligned}$$

Thus

$$\begin{aligned}
\sum_{j=1}^J (X_j - \tilde{X}_j)^2 &= \left( \left( \sum_{i=n+1}^N 1 \right) + \sum_{i \neq i', i, i' = n+1, \dots, N} 1\{U_i = U_{i'}\} \right) 1\{n < N\} \\
&\quad + \left( \left( \sum_{i=N+1}^n 1 \right) + \sum_{i \neq i', i, i' = N+1, \dots, n} 1\{U_i = U_{i'}\} \right) 1\{n > N\} \\
&\leq |N - n + (N - n)(N - n - 1)| \\
&= (N - n)^2. \quad (117)
\end{aligned}$$

Therefore, using Markov's inequality, we have

$$\begin{aligned}
\Pr \left( \sum_{j=1}^J (X_j - \tilde{X}_j)^2 > n^2 \varepsilon^2 / 4 \right) &\leq \Pr(|N - n|^2 > n^2 \varepsilon^2 / 4) \\
&\leq \frac{16E[(N - n)^4]}{n^4 \varepsilon^4} = \frac{16(3n^2 + n)}{n^4 \varepsilon^4}. \quad (118)
\end{aligned}$$

where the equality holds by  $N \sim \text{Poisson}(n)$ . For  $\sum_{j=1}^J (\tilde{X}_j - np_j)^2$  consider the derivation:

$$\begin{aligned}
&E \left[ \left( \sum_{j=1}^J (\tilde{X}_j - np_j)^2 \right)^2 \right] \\
&= \sum_{j=1}^J E[(\tilde{X}_j - np_j)^4] + \sum_{j \neq j', j, j' = 1, \dots, J} E[(\tilde{X}_j - np_j)^2] E[(\tilde{X}_{j'} - np_{j'})^2] \\
&= \sum_{j=1}^J (3n^2 p_j^2 + np_j) + \sum_{j \neq j', j, j' = 1, \dots, J} n^2 p_j p_{j'} \\
&= 2n^2 \sum_{j=1}^J p_j^2 + n + n^2 \leq 3n^2 + n. \quad (119)
\end{aligned}$$

Therefore,

$$\Pr \left( \sum_{j=1}^J (\tilde{X}_j - np_j)^2 > n^2 \varepsilon^2 / 4 \right) \leq \frac{16(3n^2 + n)}{n^4 \varepsilon^4} \quad (120)$$

Equations (115), (118), and (120) together concludes the proof.  $\square$

## D Random Coefficient Logit

In this section, we prove a lemma that establishes Assumption 1 for the random coefficient logit model.

**Lemma 9.** *Consider the random coefficient logit model in Example 4.2. Also assume that (i)  $w_{jt}$  is bounded, i.e.  $\|w_{jt}\| \leq \bar{w}$ ; (ii)  $\sup_{\lambda \in \Lambda} \sup_{\|w\| \leq \bar{w}} \int \exp(2w'v) dF(v; \lambda) < \infty$ , (iii)  $\inf_{t=1, \dots, T} \inf_{\pi_t \in \Delta_{J_t}^0} \pi_{0t} \geq \underline{\varepsilon}_0 > 0$  for all  $T$ , and (iv) there exists  $e_1 > 0$  and  $0 < e_2 < \underline{\varepsilon}_0/2$  such*

*that, the maximum eigenvalue of  $\int \tilde{\pi}_t(v) \tilde{\pi}_t(v)' dF(v; \lambda)$   $\begin{pmatrix} \tilde{\pi}_{1t}^{-1} & 0 & \dots & 0 \\ 0 & \tilde{\pi}_{2t}^{-1} & \dots & 0 \\ \dots & \dots & \ddots & 0 \\ 0 & 0 & 0 & \tilde{\pi}_{J_t}^{-1} \end{pmatrix}$  is less than  $1 - e_1$  for all  $\lambda \in \Lambda$ , and all  $\tilde{\pi}_t \in \Delta_{J_t}^{e_2}$  for all  $t = 1, \dots, T$  and  $T = 1, 2, 3, \dots$ , where*

$$\tilde{\pi}_{jt}(v) = \frac{\exp(w'_{jt}v + \delta_{jt}(\tilde{\pi}_t; \lambda))}{1 + \sum_{k=1}^{J_t} \exp(w'_{kt}v + \delta_{kt}(\tilde{\pi}_t; \lambda))}.$$

*Then Assumption 1(a) is satisfied.*

*Proof.* Without loss of generality, consider the derivative with respect to  $\pi_{\ell t}$ . For  $j = 1, \dots, J_t$ , take partial derivative with respect to  $\pi_{\ell t}$  on both sides of (24), and we get:

$$\begin{aligned} & \frac{\partial \check{\delta}_{jt}(\tilde{\pi}_t; \lambda)}{\partial \pi_{\ell t}} \\ &= \int \frac{\exp(w'_{jt}v) \exp(\check{\delta}_{jt}(\tilde{\pi}_t; \lambda))}{\left(1 + \sum_{k=1}^{J_t} \exp(\check{\delta}_{kt}(\tilde{\pi}_t; \lambda) + w'_{kt}v) \tilde{\pi}_{kt}\right)^2} \\ & \cdot \left( \exp(\check{\delta}_{\ell t}(\tilde{\pi}_t; \lambda) + w'_{\ell t}v) + \sum_{k=1}^{J_t} \tilde{\pi}_{kt} \exp(w'_{kt}v) \exp(\check{\delta}_{kt}(\tilde{\pi}_t; \lambda)) \frac{\partial \check{\delta}_{kt}(\tilde{\pi}_t; \lambda)}{\partial \pi_{\ell t}} \right) dF(v; \lambda) \\ &= \tilde{\pi}_{\ell t}^{-1} \tilde{\pi}_{jt}^{-1} \int \tilde{\pi}_{jt}(v) \tilde{\pi}_{\ell t}(v) dF(v; \lambda) + \sum_{k=1}^{J_t} \left\{ \left[ \tilde{\pi}_{jt}^{-1} \int \tilde{\pi}_{jt}(v) \tilde{\pi}_{kt}(v) dF(v; \lambda) \right] \frac{\partial \check{\delta}_{kt}(\tilde{\pi}_t; \lambda)}{\partial \pi_{\ell t}} \right\}. \end{aligned}$$

Stacking the  $J_t$  equations in matrix form, we find that

$$H_t(\tilde{\pi}_t, \lambda) \frac{\partial \check{\delta}_t(\tilde{\pi}_t; \lambda)}{\partial \pi_{1t}} = b_{\ell t}(\tilde{\pi}_t; \lambda),$$

where

$$H_t(\tilde{\pi}_t, \lambda) = I - \int \tilde{\pi}_t(v) \tilde{\pi}_t(v)' dF(v; \lambda) \begin{pmatrix} \tilde{\pi}_{1t}^{-1} & 0 & \dots & 0 \\ 0 & \tilde{\pi}_{2t}^{-1} & \dots & 0 \\ \dots & \dots & \ddots & 0 \\ 0 & 0 & 0 & \tilde{\pi}_{J_t t}^{-1} \end{pmatrix},$$

and

$$b_{\ell t}(\tilde{\pi}_t; \lambda) = \begin{pmatrix} \tilde{\pi}_{\ell t}^{-1} \tilde{\pi}_{1t}^{-1} \int \tilde{\pi}_{\ell t}(v) \tilde{\pi}_{1t}(v) dF(v; \lambda) \\ \tilde{\pi}_{\ell t}^{-1} \tilde{\pi}_{2t}^{-1} \int \tilde{\pi}_{\ell t}(v) \tilde{\pi}_{2t}(v) dF(v; \lambda) \\ \vdots \\ \tilde{\pi}_{\ell t}^{-1} \tilde{\pi}_{J_t t}^{-1} \int \tilde{\pi}_{\ell t}(v) \tilde{\pi}_{J_t t}(v) dF(v; \lambda) \end{pmatrix}.$$

By condition (iv), we have that the eigenvalues of  $H_t(\tilde{\pi}_t, \lambda)$  are positive and bounded away from zero for all  $t$ , all  $\lambda$  and all  $\tilde{\pi}_t \in \Delta_{J_t}^{e_2}$ . Next we show that the elements  $b_{\ell t}(\tilde{\pi}_t; \lambda)$  is bounded uniformly over  $\ell$  and  $t$ , which will then imply that

$$\sup_{t=1, \dots, T; T=1, 2, \dots} \sup_{j, \ell=1, \dots, J_t} \sup_{\tilde{\pi}_t \in \Delta_{J_t}^{e_2}} \sup_{\lambda \in \Lambda} \left| \frac{\partial \check{\delta}_{jt}(\tilde{\pi}_t; \lambda)}{\partial \pi_{\ell t}} \right| \leq M < \infty.$$

for some  $M$ . Consider the derivation

$$\begin{aligned} \check{\delta}_{jt}(\hat{\pi}_t; \lambda) - \check{\delta}_{jt}(\pi_t; \lambda) &= \frac{\partial \check{\delta}_{jt}(\tilde{\pi}_t; \lambda)}{\partial \pi'_{\ell t}} (\hat{\pi}_t - \pi_t) \\ &\leq \|\hat{\pi}_t - \pi_t\| \left\| \frac{\partial \check{\delta}_{jt}(\tilde{\pi}_t; \lambda)}{\partial \pi'_{\ell t}} \right\| \\ &\leq \sqrt{J_t} M \|\hat{\pi}_t - \pi_t\| \\ &\leq \sqrt{J_t} M \|\hat{\pi}_t - \pi_t\|_f. \end{aligned} \tag{121}$$

Thus Assumption 1(a) holds.

To show that  $b_{\ell t}(\tilde{\pi}_t; \lambda)$  is uniformly bounded, we first show that  $\check{\delta}_{jt}(\tilde{\pi}_t; \lambda)$  is uniformly bounded. Without loss of generality, consider  $\check{\delta}_{1t}(\tilde{\pi}_t; \lambda)$ :

$$\begin{aligned} \check{\delta}_{1t}(\tilde{\pi}_t; \lambda) &= -\log \int \frac{\exp(w'_{jt} v)}{1 + \sum_{k=1}^{J_t} \exp(\check{\delta}_{kt}(\tilde{\pi}_t; \lambda) + w'_{kt} v) \tilde{\pi}_{kt}} dF(v; \lambda) \\ &\geq -\log \int \exp(w'_{jt} v) dF(v; \lambda) \end{aligned}$$

$$\geq -\log \sup_{\lambda \in \Lambda} \sup_{\|w\| \leq \bar{w}} \int \exp(w'v) dF(v; \lambda),$$

where the second inequality holds by condition (i). Then by condition (ii), we have  $\inf_{t, \lambda, \tilde{\pi}_t} \check{\delta}_{1t}(\tilde{\pi}_t; \lambda) > -\infty$ . To show that  $\sup_{t, \lambda, \tilde{\pi}_t} \check{\delta}_{1t}(\tilde{\pi}_t; \lambda) < \infty$ , consider the outside share:

$$\tilde{\pi}_{0t} = \int \frac{1}{1 + \sum_{k=1}^{J_t} \exp(\check{\delta}_{kt}(\tilde{\pi}_t; \lambda) + w'_{kt}v) \tilde{\pi}_{kt}} dF(v; \lambda).$$

By  $|\tilde{\pi}_{0t} - \pi_{0t}| \leq \|\tilde{\pi}_t - \pi_t\| < e_2 < \varepsilon_0/2$  and  $\pi_{0t} \geq \varepsilon_0$ , we have  $\tilde{\pi}_{0t} \geq \varepsilon_0/2$ . Then there must exists  $\bar{v}$  large enough such that  $\int_{\|v\| \leq \bar{v}} \frac{1}{1 + \sum_{k=1}^{J_t} \exp(\check{\delta}_{kt}(\tilde{\pi}_t; \lambda) + w'_{kt}v) \tilde{\pi}_{kt}} dF(v; \lambda) \geq \varepsilon_0/4$ . Then

$$\begin{aligned} \check{\delta}_{1t}(\tilde{\pi}; \lambda) &\leq -\log \int_{\|v\| \leq \bar{v}} \frac{\exp(w'_{1t}v)}{1 + \sum_{k=1}^{J_t} \exp(\check{\delta}_{kt}(\tilde{\pi}_t; \lambda) + w'_{kt}v) \tilde{\pi}_{kt}} dF(v; \lambda) \\ &\leq -\log \left\{ \left[ \min_{\|w\| \leq \bar{w}, \|v\| \leq \bar{v}} \exp(w'v) \right] \int_{\|v\| \leq \bar{v}} \frac{1}{1 + \sum_{k=1}^{J_t} \exp(\check{\delta}_{kt}(\tilde{\pi}_t; \lambda) + w'_{kt}v) \tilde{\pi}_{kt}} dF(v; \lambda) \right\} \\ &\leq - \left[ \min_{\|w\| \leq \bar{w}, \|v\| \leq \bar{v}} (w'v) \right] - \log(\varepsilon_0/4). \end{aligned}$$

Thus,  $\sup_{t, \lambda, \tilde{\pi}_t} \check{\delta}_{1t}(\tilde{\pi}_t; \lambda) < \infty$ .

Now we show that  $b_{\ell t}(\tilde{\pi}_t; \lambda)$  is uniformly bounded. By Cauchy-Schwarz inequality, it suffices to consider the  $\ell$ th element of  $b_{\ell t}(\tilde{\pi}_t; \lambda)$ :

$$\begin{aligned} \tilde{\pi}_{\ell t}^{-2} \int \tilde{\pi}_{\ell t}(v)^2 dF(v; \lambda) &= \int \left( \frac{\exp(w'_{\ell t}v + \check{\delta}_{\ell t}(\tilde{\pi}_t; \lambda))}{1 + \sum_{k=1}^{J_t} \exp(w'_{kt}v + \check{\delta}_{kt}(\tilde{\pi}_t; \lambda)) \pi_{kt}} \right)^2 dF(v; \lambda) \\ &\leq \exp(2\check{\delta}_{\ell t}(\tilde{\pi}_t; \lambda)) \int \exp(2w'_{\ell t}v) dF(v; \lambda). \end{aligned}$$

Then by condition (ii) and  $\sup_{t, \lambda, \tilde{\pi}_t} \check{\delta}_{\ell t}(\tilde{\pi}_t; \lambda) < \infty$ , we have

$$\sup_t \sup_{\lambda} \sup_{\|\tilde{\pi}_t - \pi_t\| \leq e_2} \|\tilde{\pi}_{\ell t}^{-2} \int \tilde{\pi}_{\ell t}(v)^2 dF(v; \lambda)\| < \infty.$$

This shows the uniform boundedness of the elements of  $b_{\ell t}(\tilde{\pi}_t; \lambda)$ .  $\square$

## E Approximate Log Share

In this section, we show some theoretical derivation that provides further support for the finiteness of  $\bar{v}_\ell$  and the approximate value of  $\underline{v}_u$ . Lemma 10 shows that  $\iota^*(n, n\pi)$  approaches

$(1 - \pi)/2$  when  $n\pi$  is large. Lemma 11 shows that  $\iota^*(n, n\pi)$  approaches  $n\pi$  when  $n\pi$  is small. Both align well with the numerical results shown in Figure 2 and Table 2. Thus, we are confident that the conclusions regarding the approximate values of  $\bar{\iota}_\ell$  and  $\underline{\iota}_u$  drawn in Section 3.3 are correct, even though a fully rigorous theoretical proof is out of reach due to the lack of an analytical solution for the expectation of the logarithm of mean-shifted binomial or Poisson random variables. For two sequences of positive numbers  $a_n$  and  $b_n$ , we denote  $a_n \propto b_n$  if  $a_n/b_n = O(1)$  and  $b_n/a_n = O(1)$ .

**Lemma 10.** *Let  $q$  follow a binomial distribution with parameters  $(n, \pi)$ . Consider a sequence of binomial distributions such that  $\pi \propto n^{-\nu}$  with  $\nu \in [0, 1)$ . Then along this sequence we have:*

(a) *For any fixed constant  $\iota > 0$ ,*

$$E[\log(q + \iota) - \log(n\pi)] = \frac{\iota}{n\pi} - \frac{1 - \pi}{2n\pi} + o((n\pi)^{-1}).$$

where the  $o((n\pi)^{-1})$  is uniform over  $\iota$  in any bounded closed subinterval of  $(0, \infty)$ .

(b)  $\iota^*(n, n\pi) - (1 - \pi)/2 \rightarrow 0$ .

**Lemma 11.** *Let  $q$  follow a binomial distribution with parameters  $(n, \pi)$ . Consider a sequence of binomial distributions such that  $\pi \propto n^{-\nu}$  with  $\nu > 1$ . Then along this sequence we have*

$$\frac{\iota^*(n, n\pi)}{n\pi} \rightarrow 1.$$

*Proof of Lemma 10.* (a) First note that by the Chernoff's inequality, we have for any  $c > 0$

$$\Pr(|q - n\pi| > (n\pi)^{0.5+c}) \leq 2 \exp\left(-\frac{(n\pi)^{2c}}{3}\right). \quad (122)$$

Decompose  $E[\log(q + \iota) - \log(n\pi)]$  as

$$\begin{aligned} E[\log(q + \iota) - \log(n\pi)] &= E[(\log(q + \iota) - \log(n\pi))1\{|q - n\pi| \leq (n\pi)^{0.5+c}\}] \\ &\quad + E[(\log(q + \iota) - \log(n\pi))1\{|q - n\pi| > (n\pi)^{0.5+c}\}] \Pr(|q - n\pi| > (n\pi)^{0.5+c}) \end{aligned} \quad (123)$$

To bound each of the summands of the right-hand side of (123), first consider the derivation

$$\begin{aligned} \exp(-(n\pi)^{2c}/3) &= \exp(-(n\pi)^{2c}/3)n^2n^{-2} \\ &= \exp(-(n\pi)^{2c}/3 + 2 \log n)n^{-2} = o(n^{-2}) \end{aligned} \quad (124)$$

where the last equality holds because  $(n\pi)^{2c}/(3 \log n) \propto n^{2c(1-\nu)}/(3 \log n) \rightarrow \infty$ .

Now consider the second summand of the right-hand side of (123). By (122), it is bounded by

$$\begin{aligned} & 2 \max\{|\log(\iota) - \log(n\pi)|, |\log(n + \iota) - \log(n\pi)|\} \exp(-(n\pi)^{2c}/3) \\ & \leq C \log(n) \exp(-(n\pi)^{2c}/3) \\ & = o(n^{-2} \log n) = o((n\pi)^{-1}), \end{aligned} \tag{125}$$

where  $C$  is a universal constant, the inequality holds because  $\iota$  is a fixed positive constant and  $\log(n\pi) \leq \log n$ , the first equality holds due to (124) and the second equality holds because  $n\pi \propto n^{(1-\nu)} = o(n^2/\log(n))$ . It is easy to see that the  $o(\cdot)$ 's are uniform in  $\iota$  on any compact interval on  $(0, \infty)$ .

Write  $(q + \iota)/(n\pi) = 1 + (q - n\pi + \iota)/(n\pi)$ , and use a Taylor series expansion of the logarithm around 1, and we can write the first summand of the right-hand side of (123) as

$$\begin{aligned} & E[(\log(q + \iota) - \log(n\pi))1\{|q - n\pi| \leq (n\pi)^{0.5+c}\}] \\ & = (n\pi)^{-1} E[(q - n\pi + \iota)1\{|q - n\pi| \leq (n\pi)^{0.5+c}\}] \\ & - 2^{-1}(n\pi)^{-2} E[(q - n\pi + \iota)^2 1\{|q - n\pi| \leq (n\pi)^{0.5+c}\}] \\ & + 2(3!)^{-1}(n\pi)^{-3} E[(1 + \tilde{x})^{-3}(q - n\pi + \iota)^3 1\{|q - n\pi| \leq (n\pi)^{0.5+c}\}], \end{aligned} \tag{126}$$

where  $\tilde{x}$  is a value on the interval  $[0, (q - n\pi + \iota)/(n\pi)]$ . Consider the derivation:

$$\begin{aligned} (n\pi)^{-3} E[(1 + \tilde{x})^{-3}(q - n\pi + \iota)^3 1\{|q - n\pi| \leq (n\pi)^{0.5+c}\}] & \leq C(n\pi)^{-3}(\iota^3 + (n\pi)^{1.5+3c}) \\ & = o((n\pi)^{-1}), \end{aligned} \tag{127}$$

where  $C$  is a universal constant, and the equality holds when we pick a  $c \in (0, 1/6)$ . Also consider the derivation

$$\begin{aligned} & E[(q - n\pi + \iota)^2 1\{|q - n\pi| \leq (n\pi)^{0.5+c}\}] \\ & = E[(q - n\pi + \iota)^2] - E[(q - n\pi + \iota)^2 | |q - n\pi| > (n\pi)^{0.5+c}] \Pr(|q - n\pi| > (n\pi)^{0.5+c}) \\ & = E[(q - n\pi + \iota)^2] - O(n^2 \exp(-(n\pi)^{2c}/3)) \\ & = E[(q - n\pi + \iota)^2] + o(1) \\ & = n\pi(1 - \pi) + \iota^2 + o(1), \end{aligned} \tag{128}$$

where the second equality holds by (122) and  $q \leq n$ , the third equality holds by (124) and the last equality holds because  $q$  follows the binomial distribution with parameters  $(n, \pi)$ .

By similar derivation, we have

$$E[(q - n\pi + \iota)1\{|q - n\pi| \leq (n\pi)^{0.5+c}\}] = \iota + o(1). \quad (129)$$

Combining (126)-(129), we have

$$E[(\log(q + \iota) - \log(n\pi))1\{|q - n\pi| \leq (n\pi)^{0.5+c}\}] = (n\pi)^{-1}(\iota - 2^{-1}(1 - \pi) + o(1)). \quad (130)$$

Equations (123), (125), and (130) together prove part (a) of the lemma.

(b) It is without loss of generality to assume that  $\pi \rightarrow \pi_\infty$  for some  $\pi_\infty \in [0, 1]$  as  $n \rightarrow \infty$ . (If not, we can consider subsubsequences of arbitrary subsequences of  $\{n\}$  along with  $\pi$  converges. Such subsubsequences always exist because  $[0, 1]$  is a compact set.) We consider two cases:

- (i)  $\pi_\infty = 1$ . Suppose there exists a  $\underline{c} > 0$  such that  $\iota^*(n, n\pi) > \underline{c}$  infinitely often. Then by the monotonicity of the logarithm function, we have, infinitely often,

$$(n\pi)E[\log(q + \underline{c}) - \log(n\pi)] \leq (n\pi)E[\log(q + \iota^*(n, n\pi)) - \log(n\pi)] = 0, \quad (131)$$

where the equality holds by the definition of  $\iota^*(n, n\pi)$ . On the other hand, part(a) of the lemma implies that

$$(n\pi)[\log(q + \underline{c}) - \log(n\pi)] \rightarrow \underline{c} > 0. \quad (132)$$

This and (131) form a contradiction. Thus, there does not exist a  $\underline{c} > 0$  such that  $\iota^*(n, n\pi) > \underline{c}$  infinitely often. This implies that  $\iota^*(n, n\pi) \rightarrow 0$ , and in turn implies part (b) of the lemma.

- (ii)  $\pi_\infty \in (0, 1)$ . Let  $\underline{c} = (1 - \pi_\infty)/4$  and  $\bar{c} = (1 - \pi_\infty)$ . Suppose that  $\iota^*(n, n\pi) < \underline{c}$  ( $\iota^*(n, n\pi) > \bar{c}$ ) infinitely often. Then by the monotonicity of the logarithm function, we have, infinitely often,  $(n\pi)E[\log(q + \underline{c}) - \log(n\pi)] \geq 0$  ( $(n\pi)E[\log(q + \bar{c}) - \log(n\pi)] \leq 0$ ). But part(a) of the lemma implies that  $(n\pi)E[\log(q + \underline{c}) - \log(n\pi)] \rightarrow \underline{c} - (1 - \pi_\infty)/2 < 0$  ( $(n\pi)E[\log(q + \bar{c}) - \log(n\pi)] \rightarrow \bar{c} - (1 - \pi_\infty)/2 > 0$ ). These form a contradiction. Thus,  $\iota^*(n, n\pi) \in [\underline{c}, \bar{c}]$  eventually. This, the compactness of the interval  $[\underline{c}, \bar{c}]$ , and part (a) of the lemma together imply that

$$(n\pi)[\log(q + \iota^*(n, n\pi)) - \log(n\pi)] - (\iota^*(n, n\pi) - (1 - \pi)/2) \rightarrow 0. \quad (133)$$

But  $\log(q + \iota^*(n, n\pi)) - \log(n\pi) = 0$  by the definition of  $\iota^*(n, n\pi)$ . Thus,  $\iota^*(n, n\pi) -$

$(1 - \pi)/2 \rightarrow 0$ , which shows part (b) of the lemma. □

*Proof of Lemma 11.* By the definition of  $\iota^*(n, n\pi)$ , we have

$$\begin{aligned} 0 &= E[\log(q + \iota^*(n, n\pi)) - \log(n\pi)] \\ &= E[\log((q + \iota^*(n, n\pi))/(n\pi))] \\ &= (1 - \pi)^n \log(\iota^*(n, n\pi)/(n\pi)) + (1 - (1 - \pi)^n)E[\log((q + \iota^*(n, n\pi))/(n\pi)) | q > 0]. \end{aligned} \quad (134)$$

Thus,

$$\log\left(\frac{\iota^*(n, n\pi)}{n\pi}\right) = -\frac{1 - (1 - \pi)^n}{(1 - \pi)^n} E[\log((q + \iota^*(n, n\pi))/(n\pi)) | q > 0] \quad (135)$$

For  $1 \leq q \leq n$ , since  $\iota^*(n, n\pi) \leq n$ , we have  $0 < \log((q + \iota^*(n, n\pi))/(n\pi)) \leq \log(2) - \log(\pi)$ .

Thus,

$$|E[\log((q + \iota^*(n, n\pi))/(n\pi)) | q > 0]| \leq (\log(2) - \log(\pi)) \leq 2|\log(\pi)|. \quad (136)$$

where the second inequality holds for large enough  $n$  since  $n\pi \rightarrow 0$ . Also consider

$$\begin{aligned} 1 - (1 - \pi)^n &= \binom{n}{1}\pi - \binom{n}{2}\pi^2 + \cdots + (-1)^{n+1}\binom{n}{n}\pi^n \\ &\leq n\pi + (n\pi)^2 + \cdots + (n\pi)^n = \frac{n\pi(1 - (n\pi)^n)}{1 - n\pi} \end{aligned} \quad (137)$$

Thus, for large enough  $n$ , we have

$$\begin{aligned} \left| \log\left(\frac{\iota^*(n, n\pi)}{n\pi}\right) \right| &\leq 2|\log \pi| \frac{n\pi(1 - (n\pi)^n)}{1 - n\pi - n\pi(1 - (n\pi)^n)} \\ &= 2n\pi|\log(\pi)|(1 + o(1)) \\ &= o(1), \end{aligned} \quad (138)$$

where the inequality holds by (136) and (137), and the equalities hold by  $\pi \propto n^{-v}$  with  $v > 1$ . This proves the lemma. □



## References

- Andrews, D. W. K. and Shi, X. (2013). Inference based on conditional moment inequality models. *Econometrica*, 81.
- Berry, S., Levinsohn, J., and Pakes, A. (1995). Automobile prices in market equilibrium. *Econometrica: Journal of the Econometric Society*, pages 841–890.
- Devroye, L. (1983). The equivalence of weak, strong and complete convergence in  $l_1$  for kernel density estimates. *The Annals of Statistics*, 11:896–904.
- Dezhbakhsh, H., Rubin, P. H., and Shepherd, J. M. (2003). Does capital punishment have a deterrent effect? new evidence from postmoratorium panel data. *American Law and Economics Review*, 5(2):344–376.
- Head, K., Mayer, T., et al. (2013). Gravity equations: Workhorse, toolkit, and cookbook. *Handbook of international economics*, 4.
- Nurski, L. and Verboven, F. (2016). Exclusive dealing as a barrier to entry? evidence from automobiles. *The Review of Economic Studies*, 83(3):1156.
- Pollard, D. (1990). Empirical process theory and application, nsf-cbms regional conference series in probability and statistics. II:Institute of Mathematical Statistics.
- Quan, T. W. and Williams, K. R. (2015). Product variety, across-market demand heterogeneity, and the value of online retail. *Working Paper*.
- Sherman, R. P. (1993). The limiting distribution of the maximum rank correlation estimator. *Econometrica*, 61:123–137.