

Limit Theorems for Estimating the Parameters of Differentiated Product Demand Systems*

Steve Berry[†], Oliver B. Linton[‡] and Ariel Pakes[§]

Yale University, London School of Economics, and Harvard University

May 20, 2003

Abstract

We provide an asymptotic distribution theory for a class of Generalized Method of Moments estimators that arise in the study of differentiated product markets when the number of observations is associated with the number of products within a given market. We allow for three sources of error: sampling error in estimating market shares, simulation error in approximating the shares predicted by the model, and the underlying model error. It is shown that the estimators are CAN provided the size of the consumer sample and the number of simulation draws grow at a large enough rate relative to the number of products. We consider the implications of the results for Berry, Levinsohn, and Pakes' (1995) random coefficient logit model and the pure characteristic model analyzed in Berry and Pakes (2002). The required rates differ for these two frequently used demand models. A small Monte Carlo study shows that the difference in asymptotic properties of the two models are reflected, in quite a striking way, in the models' small sample properties. Moreover the limit distributions provide a good approximation to the actual monte carlo distribution of the parameter estimates. The results have important implications for the computational burden of the two models.

Journal of Economic Literature Classification: C13, C15, C35, L13

Key Words: Choice models; Method of Moments; Multinomial; Random Coefficients; Vertical Model.

1 Introduction

We are often interested in estimating parameters of demand (or production) functions from data on the quantity, price, characteristics (and perhaps the production inputs) of a set of

*We would like to thank the National Science Foundation for financial support.

[†]Department of Economics, Yale University, 27 Hillhouse Avenue, New Haven, CT 06520-8281, USA. e-mail: steven.berry@econ.yale.edu

[‡]Department of Economics, London School of Economics, Houghton Street, London WC2A 2AE, United Kingdom. e-mail: lintono@lse.ac.uk

[§]Department of Economics, Harvard University, Cambridge, MA. e-mail: ariel@ariel.fas.harvard.edu

products that interact in an imperfectly competitive market. In the simplest case, there is a national market with one observation per product, and the approximations used for the distribution of the estimators are obtained by taking the limit as the number of those products, say J , grows large.

This paper is concerned with the limiting properties of such estimators. The fact that we allow for interactions among firms implies that to obtain consistency and asymptotic normality as J grows large we need to modify standard limiting arguments. The modifications are quite general and apply to a broad class of models used in empirical Industrial Organization. We are particularly concerned with applications to the demand for differentiated products; applications which are complicated by the presence of sampling and simulation errors which enter the estimating equations in a non-linear fashion. Indeed our main focus is on conditions which guarantee the consistency and asymptotic distributions of estimators of the parameters of familiar differentiated product discrete-choice demand systems.

Before proceeding we explain why we made certain choices in the way we conduct the analysis. Of particular importance is our choice of taking limits in the number of products. The argument here is entirely practical. First Industrial Organization often has to deal with markets in which both: J is quite large (large enough to think limiting approximations in dimension J are likely to be relevant), and the theory of imperfect competition is clearly relevant (partly because of spatial competition and multi-product firms). That is using approximations obtained from limits as J grows large is often appropriate. Second, though dynamic models often make J endogenous, the relationship between J and its observable determinants (market size, ownership of products, etc.) varies with the relevant model for the problem at hand (see, for e.g., Sutton (1991)). Hence in order to take limits in the determinants of J (rather than in J itself) we would need to discuss a host of issues which are not central to the goals of this paper. Instead we suffice with conditions on how shares behave as J increases (condition S below); conditions which are likely to be satisfied by a range of dynamic models.

We focus the analysis on a process which generates a single cross section of products even though estimates of parameters of differentiated product demand systems are often obtained from richer data generating processes. For example, micro data which matches individuals to the products they choose, or regional and/or time series variance in the product level data are often also available. However, as discussed below, in most (though not all) of these cases J will still be one of the relevant limiting dimensions, and as a result arguments similar to those given here will still have to be used to rationalize the limiting properties of the parameter estimates.¹

One final point. We have chosen not to condition our major results on the nature of competition in the product market (e.g., Nash in prices or quantities). This adds to the generality of our results but, as we will explain, leaves open an interesting and important set of questions on the efficiency of alternative estimators. At appropriate points we will comment on different aspects of this efficiency issue. However, for reasons explained below, we do not have a *practical* solution to the problem of obtaining efficient estimators for the models typically used in empirical I.O.

¹Indeed we do not know of any empirical work on differentiated product demand systems which does not generate their objective function by forming averages over the products in a given market. As a result they all have to worry about the interactions between products that lie at the heart of our analysis.

Background on the Model and Results

Discrete choice differentiated product demand systems posit that the utility of the consuming unit is a function of: parameters, θ , observed product characteristics, x , random consumer tastes, and unobserved (by the econometrician) product characteristics, ξ . Some of the observed characteristics (e.g., price) may be correlated with ξ . The consuming unit either chooses one of the J products marketed or it chooses not to spend any money on the goods in this market (in which case we say the consumer chooses the “outside” alternative). Each unit makes the choice that maximizes its utility. The choices of different consumers differ because of their tastes, and the distribution of those tastes is denoted by P^0 .

Our estimate of the model’s market shares, say $\sigma(\theta, x, \xi, P)$, are generated by simply adding up over the choices of consuming units with taste distribution P , where P is typically the empirical distribution of tastes from a random sample drawn from P^0 . We observe the *actual* market shares, s . Up to sampling error, these are assumed to be the market shares generated by the model at the true (θ^0, P^0) .

To see how we find estimators for our model note that the true value of the unobservables are implicitly defined by the system

$$\sigma(\theta^0, x, \xi, P^0) = s^0, \quad (1)$$

where θ^0 is the true value of the parameter vector and s^0 is the population value of the market shares (it does not contain sampling error).

The equation $\sigma(\theta, x, \xi, P) = s$ can be solved for ξ as a function of (θ, x, s, P) . An identifying assumption on the conditional distribution of $\xi(\theta^0, x, s^0, P^0)$ is made and the θ vector is estimated by method of moments. For example, if we assume a zero covariance restriction between some exogenous vector of instruments, z , and the unobserved characteristics, our moment restriction would be

$$E[G_J(\theta)] \equiv E\left[\frac{1}{J} \sum_{j=1}^J z_j \xi_j(\theta, x, s, P)\right] = 0 \quad (2)$$

at $\theta = \theta_0$, and our estimate of θ would minimize a norm in $\frac{1}{J} \sum_{j=1}^J z_j \xi_j(\theta, x, s, P)$.

Several econometric issues arise in this context. First, unlike a traditional microeconomic cross-section, when we add new observations (products) to the market, we expect the shares and prices of the existing products to change, indicating dependence in the data generating process. Similar problems arise in other contexts involving interacting firms (e.g., production function estimation). To our knowledge no analysis of the limiting properties of parameter estimates as the number of products grow large in an imperfectly competitive market is available, although those properties seem fundamental to empirical work in industrial organization.

In our context the interdependence of firms’ decisions implies that away from the true value of θ the observations on $\xi_j(\theta, x, s^0, P^0)$ are not independent from one another. That is since both s_j and p_j are endogenously determined as a function of the characteristics of other products (as well as of own-product characteristics) there is conditional dependence in the estimate of ξ when $\theta \neq \theta^0$. As a result, consistency proofs that require uniform convergence of objective functions, uniform over all possible values of θ , cannot be used (at least not without a specification for how prices and shares behave as the number of products

grow large – two aspects of the problem we do not want to hinge our proofs on). Relatedly, efficient instruments are likely to be a function of the characteristics of *all* of the products, and this generates instruments that are not independent over j .

We show how to obtain a consistency proof based on a property of the limiting value of the objective function that can frequently be evaluated *a priori*. Given consistency, we then require only local properties of the objective function to characterize the limit distribution of the parameter estimates. As a result we are able to use a “triangular array” argument for the limit distribution of the objective function at $\theta = \theta^0$, together with simpler local convergence results (smoothness assumptions will do), to prove asymptotic normality of the parameter estimates. Our approach to these problems should be broadly applicable to a wide range of models of equilibrium markets.

A further problem turns out to be quite important in estimating demand parameters when either; [i] the function $\sigma(\cdot)$ is an integral which cannot be calculated analytically and as a result is estimated by Monte Carlo simulation, or [ii] the observed market shares, s , are based on a random sample of consumers of size n and hence are subject to sampling error. In these cases, the disturbances generated by the simulation and sampling processes also impact on the distribution of the estimators. Importantly, it works out that the impact of those disturbances differ markedly depending on which of the available differentiated product demand models are used. That is the nature of competition in demand space feeds back into the asymptotic limit theory and causes the rates of convergence of the estimators for different demand models to differ. This generates different limit theorems for the different demand models. Moreover as clearly illustrated by our Monte Carlo results these differences imply that the computational burdens of the two models are different, and the importance of the differences depend on the characteristics of the data.

In particular we show that under fairly general conditions the estimators of the parameters of the random coefficients logit (or probit) demand systems discussed in Berry, Levinsohn and Pakes (1995; henceforth BLP) will be consistent if $J \log J/n$ and $J \log J/R$ converge to zero as J increases. For asymptotic normality at rate \sqrt{J} in these cases we require J^2/n and J^2/R to be bounded. That is, to obtain a consistent and asymptotically normal estimator for the parameters of these models we require the number of simulation draws and the size of the consumer sample to grow as the *square* of the growth in the number of products. So to obtain precise parameter estimates from these models we expect to need to use a relatively large number of simulation draws, especially when the number of products is large.

The second class of demand models we consider in detail are the “pure characteristic” models. Their theoretical lineage dates back at least to Hotelling’s 1929 horizontal model and have been used extensively in the context of the vertical model introduced by Shaked and Sutton (1982). Berry and Pakes (2002) endow the pure characteristics model with an estimation algorithm analogous to the estimation algorithm provided in BLP and discuss the advantages of the pure characteristics framework (focusing on the analysis of the demand for, and the welfare implications of, new goods).

We show that to estimate the parameters of the uni-dimensional (vertical) pure characteristic model consistently we require only that n and R increase at rate $\log J$, while for asymptotic normality we require only that J/n and J/R stay bounded. We also explain why the multidimensional pure characteristic model is likely to obey the same rate restrictions, but do not have a formal proof to that effect. Since the rate at which n and R must grow

for asymptotically normal parameter estimates given the pure characteristics model is the square root of the rate at which they must grow to obtain asymptotically normal estimates for BLP’s model, we expect to need much smaller numbers of simulation and sampling draws to obtain precise parameter estimates in the pure characteristics case.

The difference in results arises because differences in the nature of competition between the two models imply differences in the properties of the share functions; i.e. of $\sigma(\cdot)$ in (1). The system in equation (1) must be solved for ξ in order to implement our method of moments estimation algorithm. In the models with “diffuse” substitution patterns, such as the random coefficient logit model of BLP, all goods are substitutes for all other goods. Thus when we decrease ξ_j consumers who leave good j distribute among all other goods and as the number of goods grows large each of the (growing number of) partials, $\partial\sigma_i(\cdot)/\partial\xi_j$, goes to zero. The sampling error is in the shares, and to analyze that error’s impact on the objective function (equation 2) we have to find how it effects our estimates of ξ . The first order impact of the sampling error on ξ is $\partial\xi/\partial\sigma$, and as the elements $\partial\sigma/\partial\xi$ go to zero, the elements of $\partial\xi/\partial\sigma$ grow large (when J is large a little bit of simulation or sampling error in s causes large changes in the computed value of ξ). To obtain the asymptotic properties of our estimator we have to control the impact of sampling error on the objective function. Since the impact of a unit of sampling error grows rapidly in J , we need the variance of the sampling error to fall rapidly in J . This is accomplished by allowing the size of the consumer (and simulation) samples to grow at a rapid rate as J grows large.²

In contrast the pure characteristic model has “local” competition (products are only substitutes with a finite number of other products). As the number of products grows large each product will tend to have the same number of substitutes, and, as a result, the elements of $\partial\sigma/\partial\xi$ will stay bounded away from zero. Again the first order impact of sampling error on the objective function is given by $\partial\xi/\partial\sigma$, but if the elements of $\partial\sigma/\partial\xi$ are bounded away from zero, the elements of $\partial\xi/\partial\sigma$ remain bounded. In the pure characteristic model then the impact of a given simulation or sampling error on the computed value of ξ does not increase with J . This suggests that for large J we should be able to obtain “well behaved” parameter estimates from the pure characteristic model with fewer simulation draws than we need to use in estimating BLP’s model. We provide a Monte Carlo study which indicates that the difference is rather dramatic.

Since the number of simulation draws needed to obtain precise estimates of the objective function is likely to be larger in BLP’s model than in the pure characteristic model, the computational burden of simulation in BLP’s model is expected to be larger than in the pure characteristics model. Berry and Pakes (2002) show, however, that the computational burden of obtaining the $\xi(\cdot)$ from the system in (1) is typically larger for the pure characteristics model than it is for BLP’s model. So there is a trade off to be considered when comparing the computational burden of the two models between the ease of simulation in one model

²There is an analogy here to the impact of simulation error on the maximum likelihood estimators of discrete choice models when the choice probabilities are simulated. In that case the probabilities that determine the likelihood acts like our $\sigma(\cdot)$ function, and the impact of simulation error on the log-likelihood is larger when the underlying probabilities of a choice are small. If we let the number of choices (our J) grow large all but possibly a small number of probabilities would have to go to zero, so for consistency we would need the simulation error to converge to zero at a faster rate than the rate at which J grows. Unfortunately the analogy to maximum likelihood does not carry over to the pure characteristics model where $\sigma(\cdot)$ has notably different properties, see below.

and the ease of computing ξ in the other.

Generalizations and Limitations

There are two common ways in which actually estimated differentiated products demand systems differ from our setup. First, the same demand model can be applied to richer types of data. Second, one can add a supply side to the model. For example, one might assume a Nash pricing equilibrium and use the pricing equation together with the demand equation to estimate the demand parameters. While these extensions can greatly aid in obtaining precise parameter estimates, in most cases there is still an interest in how the estimates behave as J becomes large.

On the demand side, richer data could be either [i] observations on multiple markets across time and/or regions or else [ii] direct observations on consumers, matching observed attributes of the consumers to their choices. First consider adding more detailed consumer data within a single market. As explained in Berry, Levinsohn and Pakes (2001), the consumer data can allow one to obtain more precise estimates of parameters governing the interaction between consumer attributes and product characteristics. However, by itself the consumer data does not reveal the mean effects of the product characteristics on demand. That paper shows that in the single-market case with observed consumer choices and unobserved product characteristics, ξ , it is still necessary to take limits in J .

If one has data on multiple independent markets then the situation is more complicated and the relevant assumptions depend on the nature of the data. If the same products, or a subset of the same products, appear in every market, as in Nevo's (2001) analysis of breakfast cereals, then the observations on the unobserved quality of the product are not independent across markets, and we again require limits in dimension J (although different assumptions might be used for identification). A similar situation occurs when we have data on a given market over time, and the same, or related products, appear in different time periods (as in BLP's study of auto demand which had data on twenty years with about a hundred products per year).

If, on the other hand, there were a large number of markets with products whose unobservable characteristics are independent across markets, then one may be able to obtain CAN estimators by taking limits solely in the number of markets and not in J .³ However even then our paper's implication for how simulation and sampling error affect the estimation error in the different models should still be useful (especially if J is large). In other cases the number of markets and the number of products may each be moderately large so we will want limits as both grow but their ratio remains bounded.

Turning to the supply side, many studies have found that adding a pricing equation and then jointly estimating all parameters from the combined pricing and demand equations can markedly increase the precision of demand-parameter estimates. While the strategy has a cost in additional assumptions, the presence of the demand parameters in the pricing equation adds efficiency to the demand estimates. In this case, though, the need for asymptotics in J does not change and limit arguments similar to those given here are required (indeed it is straightforward to add the pricing equation to the analysis below).

³Retail or service sectors in different cities might be a candidate here, at least sectors that are not dominated by chains.

Adding the pricing equation does add some clarity to questions about the optimal choice of instruments for our problem (as in Chamberlain (1987)). It makes clear that optimal instruments for price will depend on the characteristics of rival products, rendering semi-parametric analysis of optimal instruments (as in Newey (1990) and (1993)) difficult if not impossible. We shall illustrate these problems in the context of our examples and provide some heuristic guidance for the choice of instruments; but we do not currently have a practical answer to the questions of optimal instruments.

Organization and Notation

The paper is organized as follows. In section 2 we present the underlying model, an overview of the main results, and the intuition underlying them. This includes two subsections which introduce our leading examples and explain the differences between them. Section 3 provides the main mathematical details of the arguments (formal proofs are relegated to an appendix) and explains how to determine rates of convergence for our models. Section 4 returns to our examples and verifies that they satisfy the conditions set out in section 4. A small Monte Carlo study is presented in section 5. This validates the theoretical arguments, and shows that the limit distributions provide an adequate approximation to the empirical distribution of the Monte Carlo estimates.

We use $\|A\| = \{\text{tr}(A'A)\}^{1/2}$ to denote the Euclidean norm of any $m \times n$ matrix A , \xrightarrow{P} to denote convergence in probability, and \implies to mean convergence in distribution. For a matrix $A_{J \times J}$, we say $A = O(g(J))$ if the absolute value of the maximum element of the matrix is of order $g(J)$.

2 The Model and An Overview of the Results

We consider a market with J competing products and an outside good. The vectors of product characteristics will be denoted by (ξ_j, x_{1j}) . The $\xi_j \in \mathbb{R}$ are characteristics which are not observed by the econometrician whereas the $x_{1j} \in \mathcal{X}_1 \subset \mathbb{R}^{d_1}$ are observed. As noted in BLP (1995) they are analogous to the disturbance in the specification of traditional demand systems and are included to account for the fact that the list of product characteristics used in estimation does not contain all the product characteristics that consumers care about. Note also that without these disturbances the model could not rationalize the data. In large markets, where sampling error in the shares is essentially absent, the model predicts that the estimated shares should fit the observed shares exactly. This would typically be impossible if there were no disturbances.

We assume that the sequence $\{\xi_j\}_{j=1}^J$ are independent draws, and, for the most part maintain the assumption that

$$E[\xi_j|x_1] = 0 \quad \text{and} \quad \sup_{1 \leq j} E[\xi_j^2|x_{1j}] < \infty \quad (3)$$

with probability one, where $x_1 = (x_{11}, \dots, x_{1J})$. The role and content of the assumption that the ξ_j 's have a conditional mean of zero is discussed in BLP. It can be replaced by other identifying assumptions without changing the logic of the underlying limit theorem. Note that (3) allows for conditional heteroskedasticity of quite general form.

In addition to the “exogenous” characteristics [those that satisfy $E(\xi_j|x_{1j}) = 0$], we allow products to have additional characteristics, say $x_{2j} \in \mathcal{X}_2 \subset \mathbb{R}^{d_2}$, which are “endogenous” (like price) in the sense of being related to the $\{\xi_j\}$. This produces a problem analogous to the traditional simultaneity problem in demand and supply estimation. We let $x_2 = (x_{21}, \dots, x_{2J})$, $x = (x_1, x_2)$, and $\xi = (\xi_1, \dots, \xi_J)$. At times we will also need explicit assumptions on the process generating x .

The model determines the purchasing decision of a household as a function of its attributes and the characteristics of the products marketed. In all our examples we will be able to explicitly aggregate over certain dimensions of household heterogeneity to obtain average purchase probabilities conditional on a finite dimensional vector of remaining household attributes, say conditional on a $\lambda \in \mathbb{R}^v$. So the model produces a map from λ , a parameter vector, $\theta \in \Theta$, where Θ is a compact subset of \mathbb{R}^k , and the vectors of product characteristics, (x, ξ) , into purchase probabilities. Let that map be $\omega(x, \xi, \lambda, \theta)$. If P is the distribution of λ , then the vector of aggregate market shares predicted by our model and given values for θ , ξ , and P are

$$\sigma(\xi, \theta, P) = \int \omega(x, \xi, \lambda, \theta) dP(\lambda), \quad (4)$$

where we have suppressed the dependence of σ on x for convenience. The map $\sigma : D \rightarrow \mathcal{S}_J$, where D is the appropriate product space, and \mathcal{S}_J is the $J + 1$ dimensional unit simplex, i.e.,

$$\mathcal{S}_J = \{(s_0, \dots, s_J)' \mid 0 \leq s_j \leq 1 \text{ for } j = 0, \dots, J, \text{ and } \sum_{j=0}^J s_j = 1\}.$$

The actual market shares in the population are given by evaluating this function at (θ^0, P^0) the true value of θ and P . We designate this vector by $s^0 = \sigma(\xi, \theta^0, P^0)$. The vector s^0 is a random quantity uniquely determined by the realization of ξ .

Although P^0 is assumed to be known, we typically will not be able to calculate $\sigma(\xi, \theta, P^0)$ analytically and will have to make do with a simulator of it, say $\sigma(\xi, \theta, P^R)$, where P^R is the empirical measure of some i.i.d. sample $\lambda_1, \dots, \lambda_R$. For example,

$$\sigma(\xi, \theta, P^R) = \int \omega(x, \xi, \lambda, \theta) dP^R(\lambda) = \frac{1}{R} \sum_{r=1}^R \omega(x, \xi, \lambda_r, \theta).$$

Relatedly the observed vector of market shares, say $s^n \in \mathcal{S}_J$ will typically be constructed from n i.i.d. draws from the population of consumers, and hence not be precisely equal to the population market shares (our s^0). Assuming then that for any fixed (ξ, θ) , say (ξ_1, θ_1) , the function $\sigma(\xi_1, \theta_1, P^R)$ is constructed from R independent, unbiased, simulation draws makes it natural to assume the disturbances emanating from the simulation and sampling processes abide by A1.

ASSUMPTION A1. *The market shares $s_\ell^n = \frac{1}{n} \sum_{i=1}^n 1(C_i = \ell)$, where C_i is the choice of the i^{th} consumer, and C_i are i.i.d. across i . For any fixed (ξ, θ) ,*

$$\sigma_\ell(\xi, \theta, P^R) - \sigma_\ell(\xi, \theta, P^0) = \frac{1}{R} \sum_{r=1}^R \varepsilon_{\ell,r}(\xi, \theta),$$

where $\varepsilon_{\ell,r}(\xi, \theta)$ are independent across r and have mean zero, while the function $\varepsilon_{\ell,r}(\xi, \theta)$ is bounded, continuous, and differentiable in ξ, θ . Define the $J \times J$ matrices $V_2 = nE_*[(s^n - s^0)(s^n - s^0)'] = \text{diag}[s^0] - s^0 s^0'$ and $V_3 = RE_*[(\sigma(\xi, \theta^0, P^R) - \sigma(\xi, \theta^0, P^0))(\sigma(\xi, \theta^0, P^R) - \sigma(\xi, \theta^0, P^0))']$, where ξ here are the true values.

Here, $\text{diag}[x]$ is notation for a diagonal matrix with x on the principal diagonal and E_* denotes expectations conditional on product characteristics (integrating out over the simulation and/or sampling disturbances). We can allow for more general simulators like those based on importance sampling advocated by BLP by simply replacing the V_3 given in A1 with the appropriate importance sampling variance covariance matrix in the results that follow.

We will make the following regularity assumptions on $\sigma(\xi, \theta, P)$.

ASSUMPTION A2. (regularity conditions for share function) *For every finite J , for all $\theta \in \Theta$, and for all P in a neighborhood of P^0 , $\partial\sigma_j(\xi, \theta, P)/\partial\xi_k$ exists, and is continuously differentiable in both ξ and θ , with $\partial\sigma_j(\xi, \theta, P)/\partial\xi_j > 0$, and for $k \neq j$, $\partial\sigma_j(\xi, \theta, P)/\partial\xi_k \leq 0$ (for $k, j = 1, \dots, J$). The matrix $\partial\sigma(\xi, \theta, P)/\partial\xi'$ is invertible for all J . Moreover, $s_j^0 > 0$ for all j .*

Note that although these properties must hold for each finite J , they need not hold in the limit. Thus although we assume that $s_\ell^0 > 0$ for all ℓ , we have $s_\ell^0 \rightarrow 0$ as $J \rightarrow \infty$ for all but possibly a finite subset of the products. Although we do not explicitly model the process which generates the products with positive market shares, below we require the process that generates the (ξ, x) tuples to satisfy certain regularity conditions.

We now outline the logic of the estimation procedure. Elsewhere, [BLP, 1995, and Berry and Pakes (2002)] we provide quite general conditions which insure that for every $(\theta, s, P) \in \Theta \times \mathcal{S}_J^o \times \mathbf{P}$, where $\mathcal{S}_J^o = \{s : 0 < s_\ell < 1 \text{ for all } \ell\}$ and \mathbf{P} is a family of probability measures, there is a unique solution for the $\xi(\theta, s, P)$ that satisfies

$$s - \sigma(\xi, \theta, P) = 0. \quad (5)$$

Thus s, ξ are in one-to one relation for any θ, P . By the implicit function theorem, Dieudonné (1969, Theorem 10.2.1), and A1, the mapping $\xi(\theta, s, P)$ is continuously differentiable in θ, s, P , in some neighborhood.⁴ The true value of ξ , $\xi(\theta^0, s^0, P^0)$, is obtained as the solution to

$$s^0 - \sigma(\xi, \theta^0, P^0) = 0. \quad (6)$$

Define the instrument matrix $z = (z_1, \dots, z_J)$ whose components $z_q = z(x_{11}, \dots, x_{1J})_q \in \mathbb{R}^\ell$, where $z(\cdot)_q : (\mathbb{R}^{d1})^J \rightarrow \mathbb{R}^\ell$, and $\ell \geq k$ (k is the dimension of θ), for $q = 1, \dots, J$. Note that we allow the value of the instruments for the j^{th} observation to be a function of the values of the characteristics of all the observations. This is because most notions of equilibrium in use [e.g., Nash in prices or quantities] imply that the endogenous variables we are instrumenting [i.e., price] are functions of both the product's own and its competitors' characteristics. We will require only weak regularity conditions on the z_q and will introduce them where needed below.

⁴With regard to P , this would be some notion of functional differentiability, see Fernholz (1983) for discussion.

Now let

$$G_J(\theta, s, P) \equiv \frac{1}{J} \sum_{j=1}^J z_j \xi_j(\theta, s, P). \quad (7)$$

The assumption that $E(\xi_j|x_1) = 0$ ensures that $E[G_J(\theta^0, s^0, P^0)] = 0$. If we were able to calculate $\xi_j(\theta, s^0, P^0)$, then (2) would suggest using as our estimate of θ the method of moments estimator, Hansen (1982), obtained by minimizing a norm of $G_J(\theta, s^0, P^0)$.

Unfortunately we observe only s^n and not s^0 , and we cannot calculate $\sigma(\xi, \theta, P^0)$ but only $\sigma(\xi, \theta, P^R)$. Consequently, what we do is substitute an estimate of ξ , obtained as that value of ξ that sets $s^n - \sigma(\xi, \theta, P^R)$ to zero and denoted by $\xi(\theta, s^n, P^R)$, into (2) and minimize the resulting objective function. Thus our estimator of θ , say $\hat{\theta}$, is defined as any random variable that satisfies

$$\|G_J(\hat{\theta}, s^n, P^R)\| = \inf_{\theta \in \Theta} \|G_J(\theta, s^n, P^R)\| + o_p(1/\sqrt{J}). \quad (8)$$

2.1 The Main Results

$\|G_J(\theta, s^n, P^R)\|$ has a distribution determined by three independent sources of randomness: randomness generated from the draws on the vectors $\{\xi_j, x_{1j}\}$, randomness generated from the sampling distribution of s^n , and that generated from the simulated distribution P^R . Analogously there are three dimensions in which our sample can grow: as n , as R , and as J grow large. Our limit theorems will require rates of growth for each dimension. Throughout we let $J \rightarrow \infty$ and make n and R deterministic functions of J , i.e., we write $n(J)$ and $R(J)$ and let $n(J), R(J) \rightarrow \infty$ at some specified rate. If $n(J), R(J) \rightarrow \infty$ at a fast enough rate, then the contribution from simulation and sampling error will be of smaller order, and the asymptotics will be dominated by the randomness of ξ . To enable us to evaluate the contribution of simulation and sampling error to the asymptotic distribution of the estimator we make assumptions on the rates of growth on n and R which insure that all three sources of randomness contribute to the asymptotics. Finally, keep in mind that both s^n and $\sigma(\xi, \theta, P^R)$ take values in \mathbb{R}^J , where J is one of the dimensions that we let grow in our limiting arguments (for expositional ease we have not indexed these functions by J , but our assumptions should be interpreted as holding for each finite J).

The fact that the dimension of the share function grows with J makes the proofs of our major results quite detailed. So we begin with an overview of the major steps in the proofs and their implications for our two leading examples. The proofs themselves are provided in the next section.

The consistency argument is established by showing that

- (i) $\sup_{\theta \in \Theta} \|G_J(\theta, s^n, P^R) - G_J(\theta, s^0, P^0)\|$ converges to zero in probability, and
- (ii) an estimator that minimized $\|G_J(\theta, s^0, P^0)\|$ over $\theta \in \Theta$ would be consistent for θ^0 .

(i) insures that neither simulation nor sampling error impacts on the consistency of our estimator. To establish it we assume that the instruments satisfy a boundedness condition

and then provide conditions which insure that $\|\xi(\theta, s^n, P^R) - \xi(\theta, s^0, P^0)\|^2/J$ converges to zero in probability uniformly in $\theta \in \Theta$.⁵

To establish (ii), we apply a version of Pakes and Pollard (1989, Theorem 3.1). This requires that: (a) $G_J(\theta^0, s^0, P^0)$ converges to zero, and (b) for all θ outside of a neighborhood θ^0 , $G_J(\theta, s^0, P^0)$ stays bounded away from zero. Since at $\theta = \theta_0$, the $\xi_j(\theta^0, s^0, P^0)$ are indeed conditionally independent of one another (conditional on all the z_j), standard laws of large numbers can be used to insure (a). The problem in using standard uniform convergence arguments to guarantee (b) is that to verify them we would require a model for how the distribution of product characteristics (including price) evolves as the number of products grows large. What we do instead is provide an asymptotic identification condition which bounds the function $\|E[G_J(\theta, s^0, P^0)]\|$ uniformly away from zero when θ lies far enough away from θ^0 . This condition, which suffices for (b), does not require that $G_J(\theta, s^0, P^0)$ converges at all, and puts only weak restrictions on how the characteristic distribution changes as J grows large. We provide the intuition underlying why we expect the identification condition to hold in the context of our examples presently.

We turn next to the asymptotic normality result. Write

$$\xi(\theta, s^n, P^R) = \xi(\theta, s^0, P^0) + \{\xi(\theta, s^n, P^R) - \xi(\theta, s^0, P^R)\} + \{\xi(\theta, s^0, P^R) - \xi(\theta, s^0, P^0)\}. \quad (9)$$

Next we express the last two terms in this expression in terms of the simulation and sampling errors and the parameters of the model. The simulation and sampling errors are defined by the $J \times 1$ vectors

$$\varepsilon^n = s^n - s^0 \quad \text{and} \quad \varepsilon^R(\theta) = \sigma^R(\theta) - \sigma(\theta),$$

where $\sigma^R(\theta) = \sigma[\xi(\theta, s^0, P^0), \theta, P^R]$ and $\sigma(\theta) = \sigma[\xi(\theta, s^0, P^0), \theta, P^0]$. By A1 both ε^n and $\varepsilon^R(\theta)$ are sums of i.i.d. mean zero random vectors with known covariance matrix.

From equation (5) and the definition of ε^n and $\varepsilon^R(\theta)$,

$$s^0 + \varepsilon^n - \varepsilon^R(\theta) = \sigma[\xi(\theta, s^n, P^R), \theta, P^0].$$

We can therefore expand the inverse map from (θ, s^n, P) to $\xi(\cdot)$ around s^0 . More formally by assumption A2, for each J , almost every P , almost all ξ , and every $\theta \in \Theta$, the function $\sigma(\xi, \theta, P)$ is differentiable in ξ , and its derivative has an inverse, say

$$H^{-1}(\xi, \theta, P) = \left\{ \frac{\partial \sigma(\xi, \theta, P)}{\partial \xi'} \right\}^{-1}. \quad (10)$$

Abbreviate $\sigma(\theta, s, P) = \sigma(\xi(\theta, s, P), \theta, P)$, $H(\theta, s, P) = H(\xi(\theta, s, P), \theta, P)$, and $H_0 = H(\theta^0, s^0, P^0)$.

Now two Taylor expansions give us the last two terms in equation (9) in terms of H_0^{-1} , ε^n and $\varepsilon^R(\theta)$. That is, $\xi(\theta, s^n, P^R) \simeq \xi(\theta, s^0, P^0) + H_0^{-1}\{\varepsilon^n - \varepsilon^R(\theta)\}$, where the approximation sign indicates that we have omitted the second order terms from the Taylor's expansion.

Substituting our approximation for $\xi(\theta, s^n, P^R)$ into the objective function, we obtain our linear approximation to $G_J(\theta, s^n, P^R)$ as

⁵Note that (s^n, P^R) is a "function valued" nuisance parameter, similar to the nuisance parameters used in semiparametric estimation; see, for e.g., Newey (1994). However, unlike the usual semiparametric case the entire vector $s^n = (s_1^n, \dots, s_J^n)'$ affects each ξ_j , and $E(s^n) = s^0$ for all J (i.e., the function of interest depends on all the observations, but the 'nonparametric estimator' has zero bias at finite J).

$$\mathcal{G}_J(\theta) = G_J(\theta, s^0, P^0) + \frac{1}{J} z' H_0^{-1} \{ \varepsilon^n - \varepsilon^R(\theta^0) \}. \quad (11)$$

Next we provide conditions under which insure that

- (a) $\sup_{\|\theta - \theta^0\| \leq \delta_J} \sqrt{J} [\mathcal{G}_J(\theta) - G_J(\theta, s^n, P^R)] \xrightarrow{P} 0$ when $\delta_J \rightarrow 0$, and
- (b) An estimator that minimized $\|\mathcal{G}_J(\theta)\|$ over $\theta \in \Theta$ would be; (i) asymptotically normal at rate \sqrt{J} , and (ii) have a variance-covariance matrix which is the sum three mutually independent terms (one resulting from randomness in the draws on product characteristics, one from sampling error, and one from simulation error).

Given consistency, a consequence of (a) is that the estimator obtained from minimization of the criterion function $\|\mathcal{G}_J(\theta)\|$, has the same limit distribution as our estimator (i.e., as $\hat{\theta}$ as defined in equation 8). Since the former is easier to analyze, we work with it. We prove (a) using a stochastic equicontinuity argument and pointwise convergence results.

To establish (b) we provide a slight generalization to Theorem 3.3 in Pakes and Pollard (1989). The generalization allows for the fact that the underlying distributions of the random variables we are taking averages of may depend on J . The proof of (b) also requires a smoothness condition on the non-random function $E[G_J(\theta, s^0, P^0)]$ at $\theta = \theta^0$, and a further stochastic equicontinuity condition on the stochastic process $G_J(\theta, s^0, P^0)$ similar to condition (iii) of Theorem 3.3 of Pakes and Pollard (1989).

As in that theorem we obtain the asymptotic distribution of $\sqrt{J}(\hat{\theta} - \theta^0)$ in terms of $\partial E[G_J(\theta^0, s^0, P^0)]/\partial \theta$ and $\text{var}[\sqrt{J}\mathcal{G}_J(\theta^0)]$. However in our case the random vector $\sqrt{J}\mathcal{G}_J(\theta^0)$ is the sum of the three terms

$$T_{J1} = \frac{1}{\sqrt{J}} \sum_{j=1}^J z_j \xi_j \quad , \quad T_{J2} = \frac{1}{\sqrt{J}} z' H_0^{-1} \varepsilon^n \quad , \quad \text{and}, \quad T_{J3} = \frac{1}{\sqrt{J}} z' H_0^{-1} \varepsilon^R(\theta^0). \quad (12)$$

These random variables are mutually independent and asymptotically normal [at rates determined by the growth of $n(J)$ and $R(J)$]. Thus $\text{var}[\sqrt{J}\mathcal{G}_J(\theta^0)] = \text{var}[T_{J1}] + \text{var}[T_{J2}] + \text{var}[T_{J3}]$. We choose $n(J)$ and $R(J)$ so that all three terms are of the same magnitude, i.e., so that the effects of share estimation and simulation are captured by our approximations.

The next section provides the details that insure that the arguments explained in this subsection are in fact correct (these guarantee the required notions of uniform convergence and that the second order terms in the Taylor expansion which produces \mathcal{G}_J from G_J do not affect the limiting properties of our estimator). We conclude this subsection by taking the heuristic argument one step further and using it to outline our results for the estimators of the two models of demand that are the focus of our attention; i) the logit model and its extension to the random coefficients logit as discussed in BLP(1995), and ii) the ‘‘pure characteristics model’’ which first appeared as the horizontal model of Hotelling (1929) [see also Shaked and Sutton’s (1982) vertical model], and was endowed with an estimation algorithm by Berry and Pakes (2002).

The Simple Logit

The utility the i^{th} individual derives from consuming product j is

$$u_{ij} = x_j' \theta + \xi_j + \epsilon_{ij} \equiv \delta_j + \epsilon_{ij}, \quad (13)$$

where x_j is a vector of observed product characteristics which typically includes price, ξ_j is an unobserved characteristic, and ϵ_{ij} is an i.i.d. (over both products and individuals) extreme value error term. Since we can add an individual specific constant to all utilities without changing the distribution of choices, there is a free normalization in this model. This is customarily resolved by setting the utility of the outside good $u_{i0} = \epsilon_{i0}$.

Individual i chooses the product which maximizes its utility. The market share function is obtained by solving for that maximum and then integrating out over the distribution of ϵ to obtain

$$\sigma_j(x, \xi, \theta) = \frac{e^{x_j' \theta + \xi_j}}{1 + \sum_{k=1}^J e^{x_k' \theta + \xi_k}}, \quad j = 1, \dots, J, \quad (14)$$

while $\sigma_0(x, \xi, \theta) = (1 + \sum_{k=1}^J e^{x_k' \theta + \xi_k})^{-1}$. Note that this is one of the few models which has an analytic form for the market share function; consequently we need not simulate that function and there is no simulation error in this model (i.e., $\epsilon^R(\theta) \equiv 0$).

The model predicts that market shares are determined by the random variables $x_j' \theta + \xi_j$. For now assume this family of random variables has bounded support [because say x_j, ξ_j , and θ have bounded support] and density bounded away from zero on this support. Note that this implies that (with probability one); (a) market shares are all of magnitude $O(1/J)$, and (b) that for all finite J all products have market shares which are strictly positive.

From (14) the model also has an analytic expression for the unobserved product characteristic

$$\xi_j(\theta, s, P^0) = \ln(s_j) - \ln(s_0) - x_j' \theta. \quad (15)$$

So our estimator is found by minimizing a norm of

$$G_J(\theta, s^n, P^0) = \frac{1}{J} \sum_{j=1}^J z_j \xi_j(\theta, s^n, P^0) = \frac{1}{J} \sum_{j=1}^J z_j [\ln(s_j^n) - \ln(s_0^n) - x_j' \theta],$$

and can be interpreted as a linear instrumental variable estimator.

Assume temporarily that $\sup_{\theta \in \Theta} \|G_J(\theta, s^n, P^0) - G_J(\theta, s^0, P^0)\|$ converges to zero in probability. Then all we require for consistency is that for all θ outside of a neighborhood θ^0 , $G_J(\theta, s^0, P^0)$ stays bounded away from zero. But

$$\|G_J(\theta, s^0, P^0) - G_J(\theta^0, s^0, P^0)\| = \left\| \frac{1}{J} \sum_{j=1}^J z_j x_j' (\theta - \theta^0) \right\|,$$

where z_j is a vector of instruments of dimension at least as large as that of x_j . Thus a sufficient condition for identification is that for J sufficiently large $J^{-1} \sum_{j=1}^J z_j' x_j$ is of full column rank with probability arbitrarily close to one.

Typically z_j will consist of the $x_{1,j}$, or the exogenous product characteristics, and instruments for price (which will frequently be treated as endogenous characteristic in the sense of

being correlated with the ξ_j). So our identification condition requires the price of the product to be a function of observables which are not collinear with that product's exogenous characteristics. To formally verify whether this condition holds we would have to specify the nature of the pricing equilibrium. However all assumptions used to approximate equilibria in differentiated product markets imply that a product's price is a function of; the structure of ownership, the characteristics of competing products, and its own and its competing product's factor prices. As a result it is typically straightforward to construct instruments that are not collinear with price; indeed the more substantive question is *which* among the possible instruments to choose. We return to this question presently.

The expansion that underlies our asymptotic normality result requires the matrix $H(\cdot)$ whose elements are $\partial\sigma_k(\cdot)/\partial\xi_j$. From (14)

$$\frac{\partial\sigma_j(x, \xi, \theta)}{\partial\xi_k} = \begin{cases} \sigma_j(x, \xi, \theta)(1 - \sigma_j(x, \xi, \theta)) & k = j \\ -\sigma_k(x, \xi, \theta)\sigma_j(x, \xi, \theta) & \text{if } k \neq j, \end{cases} \quad (16)$$

Let $S = \text{diag}(s)$ and $i = (1, \dots, 1)'$. Then $H(\theta, s, P) = S - ss'$. This is a diagonally dominant matrix with inverse

$$H(\theta, s, P)^{-1} = S^{-1} + ii'/s_0.$$

From (11) the contribution of sampling error to $\sqrt{J}\mathcal{G}_J(\theta^0)$, or T_{J2} in (12), is $J^{-1/2}z'H_0^{-1}\varepsilon^n$. Recall that to obtain a limiting distribution of the estimator we require a rate of growth for $n(J)$ that produces a finite variance for T_{J2} . To get that rate of growth for the simplest case let z contain a single variable, and suppose that $\underline{c}/J \leq s_j^0 \leq \bar{c}/J$, for $j = 0, 1, \dots$. Then since A1 insures that $\text{Var}_*(\varepsilon^n) = H_0$, if we let \bar{z}_J be the sample average of z , we have

$$\text{Var}_*\left(\frac{1}{\sqrt{J}}z'H_0^{-1}\varepsilon^n\right) = \frac{1}{nJ}z'H_0^{-1}H_0H_0^{-1}z = \frac{1}{nJ}\sum_{j=1}^J\left[\frac{z_j^2}{s_j^0} + \frac{(\sum_{j=1}^J z_j)^2}{s_0^0}\right] \leq \frac{1}{\underline{c}n}\left[\sum_{j=1}^J z_j^2 + J^2\bar{z}_J^2\right].$$

That is asymptotic normality requires $n(J)$ to grow like J^2 (this assumes \bar{z}_J is bounded). To see why this must be the case we need to go back to the relationship between the estimates of ξ and the observed shares. In the logit model an increase in any particular ξ has a small impact on the shares of all products, and since

$$\left|\sum_{k \neq j} \frac{\partial\sigma_k}{\partial\xi_j}\right| = \frac{\partial\sigma_j}{\partial\xi_j} < \infty,$$

as J grows large the impact of a ξ on the share of any given product goes down like $1/J$. It is the inverse map from changes in s to the implied $\xi(\cdot)$ that determines the influence of sampling and simulation error on our estimates of ξ . Since the elements of $H(\cdot)$ go to zero as J grows large, the elements of this inverse grow without bound in J . To counteract this effect we need to increase the number of sampling and simulation draws, i.e., reduce the variance in those errors, at a rate faster than J ; in particular we need $n \propto J^2$. Below we provide the formalities that prove this result and show that the same rate conditions hold for the random coefficient logit analyzed in BLP.

We now briefly return to the question of efficient instruments. Temporarily assume $\text{Var}(\xi_j|x_{1,j}) = \text{Var}(\xi_j)$. Then Chamberlain (1987) leads us to expect the efficient instruments to be

$$E \left[\frac{\partial \xi(\theta_0, \cdot)}{\partial \theta} | x_1 \right], \quad (17)$$

where $x_1 = [x_{1,1}, \dots, x_{1,J}]$. Given (15) the instruments for the coefficients of the $x_{1,j}$ should be the $x_{1,j}$, while the instrument for the endogenous characteristic, say p or price, should be $E[p|x_1, \theta_0]$.

To proceed further then we need a model for how prices are set. For simplicity assume that all competitors are single product firms and that equilibrium is Nash in prices. Then

$$p_j = mc_j + \sigma_j(x, \xi, \theta_0) / \left[\frac{\partial \sigma_j(x, \xi, \theta_0)}{\partial p} \right].$$

Note first that $E[p|x_1, \theta_0]$ depends on the whole distribution of the ξ ; a distribution not needed for the rest of the analysis. Moreover even assuming we knew (or were willing to estimate) that distribution, $E[p|x_1, \theta_0]$ could not be computed unless there were a unique solution to the vector of pricing equations for each ξ . In many models we use this uniqueness condition is not satisfied in which case $E[p|x_1, \theta_0]$ is not well defined (though Nash equilibrium prices are unique for the simple logit example when each product is owned by a different firm; see Caplin and Nalebuff (1991)).

Assuming uniqueness, the integral defining the instrument could be approximated by simulation; take random draws on ξ , solve the system for the equilibrium prices for each of the draws, and then average. This would add significantly to the computational burden of the estimator. An alternative would be to use a semiparametric approximation of those instruments [as in Newey (1990) and (1993)]. However the prices depend on the characteristics and factor prices of all competing products, so the number of dimensions needed for a semiparametric approximation to those instruments will typically grow in J . As shown in Pakes (1994), since most of our models imply that the pricing function is exchangeable in the order of the characteristic vectors of the competing products, the dimensionality problem can be reduced by using an exchangeable basis (a basis whose dimension does not grow in J) in forming the semiparametric estimators (see, for e.g., BLP). However in many practical situations a low order exchangeable basis does not seem to do much better than instruments obtained from less formal intuitive arguments.

Indeed though it might be difficult to construct (consistent estimates of) efficient instruments for many of our problems, once we combine Chamberlain's results with the form of the pricing function, we can often provide quite a bit of guidance for which functions of the exogenous variables are likely to produce more effective instruments. This is perhaps easiest to see in the vertical model, a model we turn to now.

The Vertical Model

Perhaps the simplest among the models with a finite set of product characteristics is the "vertical" model of Shaked and Sutton (1982). In this model the utility function is

$$u_{ij} = \delta_j - \lambda_i p_j,$$

where $\delta_j = x_j\beta + \xi_j$ and we normalize the outside alternative so that $\delta_0 = p_0 = 0$.

Order the products so $0 = \delta_0 < \delta_1 < \delta_2 < \delta_3 < \dots$. Let $F(\cdot; \theta_p)$, where θ_p is a parameter to be estimated, denote the distribution of the price coefficient (of λ), assume it is increasing over its domain and let

$$\Delta_j = (\delta_j - \delta_{j-1}) / (p_j - p_{j-1}), \quad \text{for } j = 1, \dots, J.$$

Then necessary and sufficient conditions for all goods to have positive market share in this model are that $0 = p_0 < p_1 < p_2 < \dots$, and $\Delta_j = (\delta_j - \delta_{j-1}) / (p_j - p_{j-1})$ are ordered as $\Delta_1 > \Delta_2 > \dots$. In this case the market shares are given by

$$s_0 = 1 - F(\Delta_1; \theta_p), \quad s_j = F(\Delta_j; \theta_p) - F(\Delta_{j+1}; \theta_p), \quad \text{for } j = 1, \dots, J-1, \quad \text{and } s_J = F(\Delta_J; \theta_p). \quad (18)$$

We analyze this model in detail in section 6.1.2. Here we simply want to point out two properties of its share function. First, in contrast to the logit model, in this model

$$\frac{\partial \sigma_k}{\partial \xi_j} = 0 \quad \text{for } j \notin \{j-1, j, j+1\}.$$

That is competition is “local”, i.e. only a small number of cross partials are nonzero. Consequently as J grows large none of the nonzero elements of $H(\cdot) \equiv \partial \sigma / \partial \xi$ go to zero, and the elements of $H^{-1}(\cdot)$ remain bounded. This implies that simulation and sampling error are likely to have less impact on estimators of the vertical model than on the horizontal model. Indeed it will allow us to prove an asymptotic normality result when both the number of simulation and the number of sampling draws grows at rate J (rather than J^2 as required for the logit model).

The local nature of competition in the vertical model makes it relatively easy to consider questions related to the choice of instruments for this model. Recall that we want our instruments to approximate the formula in (17) for $\theta = (\theta_p, \beta)$. Not surprisingly, then, the efficient instrument for θ_p is the conditional expectation of a function of price.

Assuming again that there is a Nash pricing equilibrium, and that each product is owned by a distinct firm

$$p_j = mc_j + \frac{F(\Delta_j; \theta_p) - F(\Delta_{j+1}; \theta_p)}{f(\Delta_j; \theta_p) \frac{\delta_j - \delta_{j-1}}{(p_j - p_{j-1})^2} + f(\Delta_{j+1}; \theta_p) \frac{\delta_{j+1} - \delta_j}{(p_{j+1} - p_j)^2}} \quad (19)$$

for $j = 1, \dots, J$, where $f(\cdot)$ is the density for $F(\cdot)$. So the price of product “ j ” depends directly on the characteristics and factor prices of the products adjacent to j and only indirectly on the factor prices and characteristics of the other products (through their effects on the prices of the adjacent products). Thus “good” instruments are likely to depend more on the characteristics and factor prices of adjacent than non-adjacent products.

In more general pure characteristics models there is more than one characteristic whose coefficient distributes among consumers and so different consumers will have different “adjacent” products (products they would substitute to in response to a sufficiently large price increase of the given product). However the characteristics and factor prices of products which “border” (in some dimension) the good whose price we are instrumenting will have a relatively large influence on its price (indeed a related result holds for the random coefficients

logit specification of BLP (1995)). When there are multi-product firms then the impact of the characteristics of adjacent products on price would differ depending on whether those products were owned by the firm which owns the product whose price is being instrumented or by a competing firm (see BLP).

The next section provides our consistency and asymptotic normality theorems, while section 4 returns to our examples and details the implications of our theorems on the properties of the estimators of their parameters.

3 A Formal Treatment of Asymptotic Properties

To formalize the arguments given in the last section we need to specify the way in which the large vector s^n approaches s^0 [as well as the corresponding model-predicted market shares]. Since these are expanding vectors in which almost all of the individual elements of s^0 are decreasing to zero, it will not suffice to specify how each component s_j^n approaches s_j^0 ; we will require stronger, uniform, notions of convergence, as is common in semiparametric estimation problems. It is also helpful to have some restrictions on the rate at which s_j^0 approaches zero. In most of the discussion we focus on the special case defined by the following assumption:

CONDITION S. *There exists positive finite constants \underline{c} and \bar{c} such that with probability one*

$$\frac{\underline{c}}{J} \leq s_\ell^0 \leq \frac{\bar{c}}{J}, \quad \ell = 0, 1, \dots, J. \quad (20)$$

In this case, each market share, including the outside option, is declining to zero at the same rate.

We will work with the product space $\Theta \times \mathcal{S}_J \times \mathbf{P}$, where \mathbf{P} is the set of probability measures, and endow the marginal spaces with (pseudo) metrics: the L_∞ metric on \mathbf{P} , $\rho_P(P, Q) = \sup_{B \in \mathcal{B}} |P(B) - Q(B)|$, where \mathcal{B} is the class of all Borel sets on \mathbb{R}^k , the Euclidean metric on Θ , $\rho_E(\theta, \theta') = \|\theta - \theta'\|$, and a metric ρ_{α, s^0} on \mathcal{S}_J defined below. We suppose that

$$\rho_{\alpha, s^0}(s, s^*) = \begin{cases} \max_{0 \leq j \leq J} \left| \frac{(s_j)^\alpha - (s_j^*)^\alpha}{(s_j^0)^\alpha} \right| & \text{if } 0 < \alpha \leq 1 \\ \max_{0 \leq j \leq J} |s_j - s_j^*| & \text{if } \alpha = 0. \end{cases} \quad (21)$$

The metric $\rho_{\alpha, s}$ depends on the parameter α ; the higher α is, the stronger is the metric.⁶ We state the theory for some α , and the choice of α will depend on the application. In the logit-like case, we use $\alpha = 1$, while in the vertical case we take $\alpha = 0$. We also put a metric on the space where ξ lives and for this we shall just take the averaged Euclidean metric $\rho_\xi(\xi, \xi^*) = J^{-1} \|\xi - \xi^*\|^2 = J^{-1} \sum_{j=1}^J (\xi_j - \xi_j^*)^2$. Finally, define for each ϵ , the following neighborhoods of θ^0 , P^0 , and s^0 : $\mathcal{N}_{P^0}(\epsilon) = \{P : \rho_P(P, P^0) \leq \epsilon\}$ and $\mathcal{N}_{s^0}(\epsilon) = \{s : \rho_{\alpha, s}(s, s^0) \leq \epsilon\}$, $\mathcal{N}_{\theta^0}(\epsilon) = \{\theta : \rho_E(\theta, \theta^0) \leq \epsilon\}$, and for each θ and any $\epsilon > 0$, define $\mathcal{N}_{\xi^0}(\theta; \epsilon) = \{\xi : \rho_\xi(\xi, \xi(\theta, s^0, P^0)) \leq \epsilon\}$.

⁶We include the $j = 0$ term because the uniform convergence of all other terms does not imply the convergence of this term.

Note that the space \mathcal{S}_J and metric ρ_s both change with J ; nevertheless, the space can be embedded in the limiting space consisting of all infinite sequences.

3.1 Consistency

The consistency result uses several assumptions in addition to A1-A2; assumptions we introduce now. A3 controls the way in which s^n approaches s^0 (and likewise for the corresponding model quantities), A4 puts fairly mild restrictions on the instruments, and A5 and A6 are our identification conditions.

ASSUMPTION A3. *The random sequences s^n and $\sigma^R(\theta)$ are consistent with respect to the corresponding metrics, i.e.,*

$$(a) \rho_{\alpha, s^0}(s^n, s^0) \xrightarrow{P} 0 \quad ; \quad (b) \sup_{\theta \in \Theta} \rho_{\alpha, \sigma(\theta)}(\sigma^R(\theta), \sigma(\theta)) \xrightarrow{P} 0, \quad (22)$$

where $\sigma^R(\theta) = \sigma[\xi(\theta, s^0, P^0), \theta, P^R]$ and $\sigma(\theta) = \sigma[\xi(\theta, s^0, P^0), \theta, P^0]$. Furthermore, we suppose that the true market shares satisfy [for the α defined in (21)]

$$(c) \frac{1}{nJ^\alpha} \sum_{j=0}^J \frac{s_j^0(1-s_j^0)}{(s_j^0)^{2\alpha}} \xrightarrow{P} 0 \quad ; \quad (d) \sup_{\theta \in \Theta} \left| \frac{1}{R \cdot J^\alpha} \sum_{j=0}^J \frac{\sigma_j(\theta)(1-\sigma_j(\theta))}{(\sigma_j(\theta))^{2\alpha}} \right| \xrightarrow{P} 0.$$

Assumption A3(a,b) is complicated by the fact that the dimensions of the vectors increase with J . Note that each $(s_\ell^n, \sigma_\ell^R)$ is a sum of independent bounded random variables with expectation s_ℓ^0 , conditional on the realization of ξ . Therefore to verify assumption A3 we require restrictions on the growth rates of $n(J)$ and $R(J)$, and on the limiting behavior of the vector s^0 .

Suppose that condition S (equation 20) holds. Then we have that $\text{var}(s_\ell^n) = O(1/nJ)$ by assumption A1. Therefore, $(s_\ell^n - s_\ell^0)/s_\ell^0 = O_p(\sqrt{J/n})$ for each $\ell = 0, 1, \dots, J$. This gives the pointwise rate of convergence. To obtain the sup-norm convergence rate [with respect to the pseudo-metric $\rho_s(s^1, s^2) = \max_{0 \leq \ell \leq J} |s_\ell^1 - s_\ell^2|/s_\ell^0$], we apply the Bonferroni and Bernstein inequalities [see Pollard (1984)] to obtain

$$\begin{aligned} \Pr \left[\max_{0 \leq \ell \leq J} \left| \frac{s_\ell^n - s_\ell^0}{s_\ell^0} \right| > \epsilon \right] &\leq \sum_{\ell=0}^J \Pr \left[\left| \frac{s_\ell^n - s_\ell^0}{s_\ell^0} \right| > \epsilon \right] \\ &\leq \sum_{\ell=0}^J \exp \left(-\frac{\epsilon^2}{2\text{var}(s_\ell^n/s_\ell^0) + 2\epsilon/n s_\ell^0} \right) \\ &\leq \sum_{\ell=0}^J \exp(-\epsilon^2 O(n/J)). \end{aligned} \quad (23)$$

A sufficient condition for (23) to decrease to zero is that $J^{1+\epsilon}/n \rightarrow 0$ for any $\epsilon > 0$. This guarantees A3(a). Assumption A3(b) is similar but requires uniformity over θ . Assumption A3(c) is implied by $J^\alpha/n \rightarrow 0$ under condition S, likewise A3(d).

Assumption A4 is a fairly mild restriction on the instruments (it will be satisfied if they are bounded). Note that there is no presumption that a law of large numbers holds since to show that we would need to be more specific about the details of how the instruments are constructed and the nature of the equilibrium.

ASSUMPTION A4. *The instruments are such that the matrix $z'z/J$ is stochastically bounded, i.e., for all $\epsilon > 0$ there exists an M_ϵ such that $\Pr[\|z'z/J\| > M_\epsilon] < \epsilon$.*

Next we provide an assumption that ensures the uniform mean square convergence for the vector $\xi(\theta, s^n, P^R)$. We reinterpret solving the equations $s = \sigma(\xi, \theta, P)$ as a minimization problem, i.e. $\xi(\theta, s, P)$ is the unique minimum of $\|s - \sigma(\xi, \theta, P)\|$. In fact it is convenient to take a certain monotonic transform of both sides of the equation $s = \sigma(\xi, \theta, P)$, where the choice of transform will depend on the application. Specifically, we introduce the component-wise transformation $\tau_J : \mathbb{R}^J \rightarrow \mathbb{R}^J$ [i.e., $\tau_J(s) = (\tau_J(s_1), \dots, \tau_J(s_J))'$] and the $J \times 1$ vector $\psi_J(\xi, \theta, s, P) = \tau_J(s) - \tau_J(\sigma(\xi, \theta, P))$. We then define

$$\xi(\theta, s, P) = \arg \min_{\xi \in \mathbb{R}^J} \|\psi_J(\xi, \theta, s, P)\| \quad (24)$$

for any θ, s, P . For any bijective transform $\tau_J(\cdot)$, (24) has the same solution. We already know that there exists a unique solution $\xi(\theta, s, P)$ to $s = \sigma(\xi, \theta, P)$ for all (θ, s, P) ; this is equivalent to saying that $\psi_J(\xi, \theta, s, P) = 0$ if and only if $\xi = \xi(\theta, s, P)$. We use the new definition of $\xi(\theta, s, P)$ as an optimization estimator to guarantee its statistical properties; in view of the increasing dimensions of (ψ_J, ξ) however, we must refine the concept of uniqueness of $\xi(\theta, s, P)$. Let $\tau_J(x) = J^{-\alpha/2} \tau_\alpha(x)$ for

$$\tau_\alpha(x) = \begin{cases} \frac{x^{1-\alpha}-1}{1-\alpha} & \text{if } 0 \leq \alpha < 1 \\ \log x & \text{if } \alpha = 1. \end{cases} \quad (25)$$

For each α the function $\tau_\alpha(\cdot)$ is monotonic. In the logit-like case, we use $\alpha = 1$ and $\tau_\alpha(x) = \log x$, while in the pure characteristics case we take $\alpha = 0$ and $\tau_\alpha(x) = x$. The next condition insures that we can, at least asymptotically, distinguish the ξ that sets the models predictions for shares equal to the actual share from other values of ξ .

ASSUMPTION A5. *For all $\delta > 0$, there exists $C(\delta)$ such that*

$$\lim_{J \rightarrow \infty} \Pr \left[\inf_{\theta \in \Theta} \inf_{\xi \notin \mathcal{N}_{\xi^0}(\theta; \delta)} \|\tau_J(\sigma(\xi, \theta, P^0)) - \tau_J(\sigma(\xi(\theta, s^0, P^0), \theta, P^0))\| > C(\delta) \right] = 1.$$

The assumptions made thus far insure that consistency depends only on the properties of $G_J(\theta, s^0, P^0)$ for $\theta \in \Theta$ as defined in equation (7). Our assumptions imply that $G_J(\theta^0, s^0, P^0) = o_p(1)$. Thus what we require is an assumption on the limiting behavior of $G_J(\cdot, s^0, P^0)$ for θ outside of a neighborhood of θ^0 (c.f. Theorem 3.1 of Pakes and Pollard, 1989). This is the role of A6 below. Note that it does not require convergence of the objective function $G_J(\theta, s^0, P^0)$ at $\theta \neq \theta^0$ (since that would require conditions on both the process generating the x 's and an equilibrium assumption).

ASSUMPTION A6. *For all $\delta > 0$, there exists $C(\delta)$ such that*

$$\lim_{J \rightarrow \infty} \Pr \left[\inf_{\theta \notin \mathcal{N}_{\theta^0}(\delta)} \|G_J(\theta, s^0, P^0) - G_J(\theta^0, s^0, P^0)\| \geq C(\delta) \right] = 1.$$

We can now state our consistency result.

THEOREM 1 [Consistency] *Suppose that A1-A6 hold for some $\alpha \in [0, 1]$ and some $n(J), R(J) \rightarrow \infty$. Then, $\widehat{\theta} \xrightarrow{P} \theta^0$.*

All our proofs are in the appendix. Note that this result applies to a wide range of models and growth rates on $n(J), R(J)$ (a fact we use later).

3.2 Asymptotic Normality

We next establish the asymptotic distribution of $\widehat{\theta}$. We do this by providing conditions under which

$$\mathcal{G}_J(\theta^0) = G_J(\theta^0, s^0, P^0) + \frac{1}{J} z' H_0^{-1} \{ \varepsilon^n - \varepsilon^R(\theta^0) \}$$

is asymptotically normal with bounded variance, while $\sqrt{J}[G_J(\theta, s^n, P^R) - \mathcal{G}_J(\theta)] = o_p(1)$ uniformly over a shrinking neighborhood of θ^0 . Additional standard arguments deliver the asymptotic distribution of $\sqrt{J}(\widehat{\theta} - \theta^0)$ in terms of the variance of $\sqrt{J}\mathcal{G}_J(\theta^0)$ and the derivative of $EG_J(\theta, s^0, P^0)$ with respect to θ . The precise magnitude of the variance of $\sqrt{J}\mathcal{G}_J(\theta^0)$ is determined by the behavior of the matrix H_0^{-1} , an issue we will come back to below.

ASSUMPTION B1. θ^0 is an interior point of Θ .

ASSUMPTION B2. For all θ in some $\delta > 0$ neighborhood of θ^0

$$E [G_J(\theta, s^0, P^0)] = \Gamma^J(\theta - \theta^0) + o(\|\theta - \theta^0\|)$$

uniformly in J . The matrix $\Gamma^J \rightarrow \Gamma$ as $J \rightarrow \infty$, where Γ has full (column) rank.

Note that B2 requires only that the expectation of $G_J(\theta, s^0, P^0)$ be differentiable (not the function itself). This condition is similar to condition (ii) of Theorem 3.3 in Pakes and Pollard (1989). What is different here is that the expectation of $G_J(\theta, s^0, P^0)$ varies with J . This is because the derivative of $\xi(\cdot)$ with respect to θ and the form of the instruments both depend on the number and characteristics of the products marketed.

ASSUMPTION B3. For all sequences of positive numbers δ_J such that $\delta_J \rightarrow 0$,

$$\sup_{\|\theta - \theta^0\| \leq \delta_J} \|\sqrt{J}[G_J(\theta, s^0, P^0) - EG_J(\theta, s^0, P^0)] - \sqrt{J}[G_J(\theta^0, s^0, P^0) - EG_J(\theta^0, s^0, P^0)]\| = o_p(1).$$

This assumption is essentially condition (iii) of Theorem 3.3 in Pakes and Pollard (1989). Given consistency, B1, and B2, it insures that an estimator that minimized $\|G_J(\theta, s^0, P^0)\|$ (an estimator we generally cannot solve for analytically) has the same distribution as the estimator which minimizes the quadratic form $\|\Gamma^J(\theta - \theta_0) + G_J(\theta^0, s^0, P^0)\|$.

To go further we need to work with the disturbances generated by the expansion in (9) and (10). Define the stochastic process in (ξ, θ, P)

$$\nu_J(\xi, \theta, P) = \frac{1}{\sqrt{J}} z' H^{-1}(\xi, \theta, P) \{ \varepsilon^n - \varepsilon^R(\theta) \}, \quad (26)$$

where $\varepsilon^n = (\varepsilon_1^n, \dots, \varepsilon_J^n)'$ and $\varepsilon^R(\theta) = (\varepsilon_1^R(\theta), \dots, \varepsilon_J^R(\theta))'$. We now show that this process has the structure of a sum of independent random variables from a triangular array. Interchanging the order of summation, and letting $z' H^{-1}(\xi, \theta, P) \equiv (a_1(\xi, \theta, P), \dots, a_J(\xi, \theta, P))$, we have

$$\nu_J(\xi, \theta, P) = \sum_{i=1}^n Y_{Ji}(\xi, \theta, P) - \sum_{r=1}^R Y_{J,r}^*(\xi, \theta, P),$$

where

$$Y_{Ji}(\xi, \theta, P) = \frac{1}{n\sqrt{J}} \sum_{j=1}^J a_j(\xi, \theta, P) \varepsilon_{ji} \quad ; \quad Y_{J,r}^*(\xi, \theta, P) = \frac{1}{R\sqrt{J}} \sum_{j=1}^J a_j(\xi, \theta, P) \varepsilon_{j,r}(\theta). \quad (27)$$

The random variables Y_{Ji} (and $Y_{J,r}^*$) are i.i.d. across i (r) with mean zero but a distribution that changes with J . It is this structure can be used to apply laws of large numbers and central limit theorems for triangular arrays to our problem. To be precise we need the following condition.

ASSUMPTION B4. *Let $Y_{Ji} = Y_{Ji}(\xi(\theta^0, s^0, P^0), \theta^0, P^0)$ and $Y_{J,r}^* = Y_{J,r}^*(\xi(\theta^0, s^0, P^0), \theta^0, P^0)$. Suppose that $\lim_{J \rightarrow \infty} E(z' \xi \xi' z / J) = \Phi_1$ and that with probability one*

$$(a) \lim_{J \rightarrow \infty} nE_* [Y_{Ji} Y_{Ji}'] = \Phi_2 \quad ; \quad (b) \lim_{J \rightarrow \infty} RE_* [Y_{J,r}^* Y_{J,r}^{*'}] = \Phi_3 \quad (28)$$

for finite positive definite non-random matrices Φ_q , $q = 1, 2, 3$. Suppose that for some $\delta > 0$, $E(\|z' \xi / \sqrt{J}\|^{2+\delta}) = o(1)$ and with probability one

$$(c) nE_* [\|Y_{Ji}\|^{2+\delta}] = o(1) \quad ; \quad (d) RE_* [\|Y_{J,r}^*\|^{2+\delta}] = o(1). \quad (29)$$

Condition B4 guarantees that $\sqrt{J} \mathcal{G}_J(\theta^0)$ is asymptotically normal with variance $\Phi = \sum_{q=1}^3 \Phi_q$. The reason for condition (29) is that use of the Lyapunov Central Limit Theorem for triangular arrays requires moment conditions holding to power $2 + \delta$. The next section will translate these conditions into restrictions on $n(J)$ and $R(J)$ for our leading cases (this will require more detailed assumptions on z and H_0).

Finally, we use a stochastic equicontinuity condition on the stochastic process (26) to handle remainder terms. This approach to asymptotics is now well established in econometrics, see the recent survey of Andrews (1994).

ASSUMPTION B5. *The process $\nu_J(\xi, \theta, P)$ is stochastically equicontinuous in (ξ, θ, P) at $(\xi(\theta^0, s^0, P^0), P^0, \theta^0)$, that is, for all sequences of positive numbers ϵ_J with $\epsilon_J \rightarrow 0$, we have*

$$\sup_{\|\theta - \theta^0\| \leq \epsilon_J} \sup_{(\xi, P) \in \mathcal{N}_{\xi^0}(\theta^0; \epsilon_J) \times \mathcal{N}_{P^0}(\epsilon_J)} \|\nu_J(\xi, \theta, P) - \nu_J(\xi(\theta^0, s^0, P^0), \theta^0, P^0)\| = o_p(1).$$

In B5 we need to insure that $\sqrt{J}[G_J(\theta, s, P) - EG_J(\theta, s^0, P^0)]$ can be made arbitrarily close to $\sqrt{J}[G_J(\theta^0, s, P) - EG_J(\theta^0, s^0, P^0)]$ (with arbitrarily large probability) by making θ close to θ^0 . This is stronger than the condition needed to make $\sqrt{J}[G_J(\theta, s^0, P^0) - EG_J(\theta, s^0, P^0)]$ close to $\sqrt{J}[G_J(\theta^0, s^0, P^0) - EG_J(\theta^0, s^0, P^0)]$ (we have also to insure that the consumer sampling and the simulation processes do not cause jumps in the disturbance process at values of θ close to θ^0). The stochastic equicontinuity assumption is sufficient to

ensure that the remainder term from the expansions we use is of smaller order in probability than $\sqrt{J}\mathcal{G}_J(\theta^0)$. We verify it for the logit case. With these conditions we can state our asymptotic normality result (again its proof is in the appendix).

THEOREM 2. [Asymptotic Normality] *Suppose that A1-A6 and B1-B5 hold for some α . Then,*

$$\sqrt{J}(\widehat{\theta} - \theta^0) \implies N[0, (\Gamma'\Gamma)^{-1}\Gamma'\Phi\Gamma(\Gamma'\Gamma)^{-1}], \text{ with } \Phi = \Phi_1 + \Phi_2 + \Phi_3.$$

Note that the $\{\Phi_q\}_{q=1}^3$ in the statement of the theorem (and as defined in B4) satisfy

$$\Phi_2 = p \lim_{n \rightarrow \infty} \frac{1}{nJ} z' H_0^{-1} V_2 H_0^{-1'} z \quad ; \quad \Phi_3 = p \lim_{n \rightarrow \infty} \frac{1}{RJ} z' H_0^{-1} V_3 H_0^{-1'} z.$$

We can obtain consistent estimates of the standard errors by substituting consistent estimates of $\{\Phi_q\}_{q=1}^3$ and Γ into the theorem's formula.⁷ We now consider alternative ways of obtaining those estimates.

When G_J is differentiable in θ let $\widehat{\Gamma} = \partial G_J(\widehat{\theta}, s^n, P^R)/\partial\theta$; this will consistently estimate Γ under quite general conditions. When G_J is not differentiable in θ we must use numerical derivatives. Let $\delta = \delta(J)$ be a small positive number and for each (k, l) let

$$\widehat{\Gamma}_{lk} = \frac{G_{Jl}(\widehat{\theta} + \delta e_k, s^n, P^R) - G_{Jl}(\widehat{\theta}, s^n, P^R)}{\delta}$$

be an estimate of Γ_{lk} , where e_k is the k^{th} unit vector. Under the conditions on the rate at which the bandwidth parameter $\delta \rightarrow 0$ as $J \rightarrow \infty$ given in Pakes and Pollard (1989), $\widehat{\Gamma}_{lk} \xrightarrow{P} \Gamma_{lk}$.

To define estimates of Φ_j we just substitute estimates of the unknown quantities in the asymptotic variance formulae and replace expectations by sample averages. Specifically, define the residual vector $\widehat{\xi} = (\widehat{\xi}_1, \dots, \widehat{\xi}_J)'$, where $\widehat{\xi}_j = \xi_j(\widehat{\theta}, s^n, P^R)$, and let

$$\widehat{\Phi}_1 = \frac{1}{J} \sum_{j=1}^J z_j z_j' \widehat{\xi}_j^2.$$

To show that $\widehat{\Phi}_1 \xrightarrow{P} \Phi_1$ we can make use of our results to establish that $J^{-1} \sum_{j=1}^J z_j z_j' (\widehat{\xi}_j^2 - \xi_j^2) \xrightarrow{P} 0$ from our existing conditions, but we also need that $J^{-1} \sum_{j=1}^J z_j z_j' \xi_j^2 \xrightarrow{P} \Phi_1$,

⁷Note that the traditional bootstrap estimator of standard errors is not well defined in our context (at least not without additional assumptions). That is, though we could bootstrap a sample of exogenous characteristics, we would need to make an equilibrium assumption before we could move from that sample to a sample that could be used in estimation. Moreover that equilibrium assumption would have to generate unique price and demand vectors, and uniqueness is not a property generated by the standard equilibrium assumptions (Nash in prices or quantities) for the demand models we typically take to data (either BLP or the pure characteristics model). Moreover for a uniqueness result to apply we would typically have to allow for the multiproduct firms we often see in our data.

which does not follow from our primitive conditions. We shall assume that this holds and so $\widehat{\Phi}_1 \xrightarrow{P} \Phi_1$.⁸

The plug-in estimates of Φ_2 and Φ_3 are consistent under our conditions because we can exploit the independence of the simulation and sampling errors. That is let $\widehat{\varepsilon}_{\ell,r} = \varepsilon_{\ell,r}(\widehat{\xi}, \widehat{\theta})$. Then if

$$\widehat{\Phi}_2 = \frac{1}{nJ} z' \widehat{H}^{-1} \widehat{V}_2 \widehat{H}^{-1'} z \quad ; \quad \widehat{\Phi}_3 = \frac{1}{RJ} z' \widehat{H}^{-1} \widehat{V}_3 \widehat{H}^{-1'} z,$$

where $\widehat{H} = H(\widehat{\theta}, s^n, P^R)$, $\widehat{V}_2 = S^n - s^n s^{n'}$, and $\widehat{V}_3 = (\widehat{V}_3)_{\ell,\ell'}$ with

$$(\widehat{V}_3)_{\ell,\ell'} = \frac{1}{R} \sum_{r=1}^R \widehat{\varepsilon}_{\ell,r} \widehat{\varepsilon}_{\ell',r}.$$

Then $\widehat{\Phi}_q \xrightarrow{p} \Phi_q$, $q = 2, 3$ under our conditions. An alternative estimator of Φ_2 can be obtained by recalculating the objective function at $\widehat{\theta}$ and s^n for independent sets or R simulation draws and computing the variance in these estimates (and a similar procedure could be used for Φ_3).⁹ Consistency requires the number of sets of simulation draws to grow large but in applied work we have found that we obtain a relatively precise estimator quite easily.

Note that the standard errors allow for conditional heteroskedasticity in the ξ_j . However if the market of interest contains multiproduct firms, and the observed characteristics did not include firm specific dummy variables, then one might want to generalize to allow the ξ to be correlated across products (say have a firm specific component).¹⁰

Also, as in Hansen (1982), we can improve the efficiency of $\widehat{\theta}$ by taking the weighted norm criterion, i.e.,

$$\|G_J(\theta, s, P)\|_{W_J}^2 = G_J(\theta, s, P)' W_J G_J(\theta, s, P)$$

for some weighting matrix W_J . The resulting class of estimators can be treated similarly to above: it suffices for asymptotic normality to make the additional assumption that $W_J \rightarrow_p W$ for some symmetric positive definite matrix W , in which case the asymptotic variance is $(\Gamma' W \Gamma)^{-1} \Gamma' W \Phi W \Gamma (\Gamma' W \Gamma)^{-1}$ (see Pakes and Pollard (1989)). The optimal weighting matrix is proportional to Φ^{-1} , and the resulting estimator has asymptotic variance $(\Gamma' \Phi^{-1} \Gamma)^{-1}$, and is efficient within this class.

A few final points on efficiency. First if we make a comparison with the estimator that is optimal when s^0, P^0 are known [and the corresponding moment $G_J(\theta, s^0, P^0)$ can be computed], we find that the variance of our estimator is strictly larger, so an estimator of the

⁸To show this we would have to limit the amount of dependence across the z_j so as to apply a law of large numbers based on a weak dependence concept like mixing. To establish that such a property holds under more primitive conditions on the economic model is a quite challenging and is left for future work. Our simulations show that the standard errors seem to be consistent in the models we examined.

⁹Frequently s^n is constructed from total sales data, rather than from a sample of consumer purchasing patterns, and in these cases n is typically so large that sampling variance has very little impact on the estimator.

¹⁰Similarly, if there were unobserved product characteristics that were determined by the date the product were introduced, and launch dates were not included in the observed x 's, then we might want to allow for a launch date specific component of the ξ . We have not generalized in this way because it would require adding a dimension to all of our indices, which would complicate our notation considerably. Moreover the added generality does not create any econometric issues that cannot be handled in standard ways.

variance of the parameter estimates which ignores sampling and simulation error will be biased downwards. Also, since we are only dealing with the demand subsystem here, our estimator can only be efficient in a limited information sense. That is in virtually all currently used pricing models the pricing equation also depends on the parameters of the demand system. So if we were willing to make an assumption on how prices are set, we could also use the pricing equation to help estimate the demand parameters. Relatedly once we made such an assumption we could look for an efficient estimator under the conditional moment restriction $E[\xi_j|x_{1j}] = 0$ as in Chamberlain (1987); see the discussion the last section.

3.3 A Comment on Rates of Convergence

Establishing conditions under which B4 is true, i.e. conditions under which the random variables $T_{J_2} \equiv J^{-1/2}z'H_0^{-1}\varepsilon^n$ and $T_{J_3} \equiv J^{-1/2}z'H_0^{-1}\varepsilon^R(\theta)$ are asymptotically normal with zero mean and finite non-zero variances, is central to the analysis. The relevant variances are obtained as $p \lim_{J \rightarrow \infty} \Phi_2(J)$ and $p \lim_{J \rightarrow \infty} \Phi_3(J)$, where

$$\Phi_2(J) = \frac{1}{nJ}z'H_0^{-1}V_2H_0^{-1}z \quad ; \quad \Phi_3(J) = \frac{1}{RJ}z'H_0^{-1}V_3H_0^{-1}z. \quad (30)$$

Keep in mind that the matrix H_0 is dimension $J \times J$ and J grows large in our limiting argument.

Since for fixed J both T_{J_2} and T_{J_3} are a sum of i.i.d. random variables central limit theorems for triangular arrays imply that it will be sufficient to find conditions on $n(J)$ and $R(J)$ that guarantee that the Φ matrices are bounded. We consider the term $\Phi_2(J)$ [similar comments apply to $\Phi_3(J)$]. The behavior of the elements of $H^{-1}(\theta, s^0, P^0)$ has a key role here, and, consequently, we will consider different scenarios regarding these quantities as is appropriate for the different demand models.

The differences arise because the different models have different implications for the components of $\partial\sigma(\cdot)/\partial\xi$. In particular in the models with “diffuse” substitution patterns, such as the random coefficient logit model of BLP in which all goods are substitutes for all other goods, that partial goes to zero as the number of products increase, and its inverse grows large. Consequently, when J is large a little bit of sampling error causes large changes in the computed value of ξ . In contrast, in the pure characteristic model, competition is “local”, the more the number of products the “closer” will your nearest competitor tend to be and the larger will be the response to small changes in the quality of the product. In these cases a little bit of simulation or sampling error will have almost no effect on the computed value of ξ .

Formally, if we let $a' = (a_1, \dots, a_J) = z'H_0^{-1}$ and suppose, without loss of generality, that z is a $J \times 1$ vector, we have [conditional on s^0]

$$\Phi_2(J) = \frac{1}{nJ} \left[\sum_{j=1}^J a_j^2 s_j^0 - \left(\sum_{j=1}^J a_j s_j^0 \right)^2 \right], \quad (31)$$

since $V_2 = \text{diag}[s^0] - s^0 s^{0'}$. The magnitude of the matrix Φ_2 depends on the vectors a and s^0 . Note that the term in square brackets in (31) can be considered to be the ‘variance’ of the vector (a_1, \dots, a_J) with respect to the multinomial like measure induced by the sequence

of weights (s_1^0, \dots, s_J^0) [note that depending on the behavior of s_0^0 , these weights do not necessarily sum to one even asymptotically].

There are three factors that influence the magnitude of $\Phi_2(J)$. First, the rate at which s_j^0 , $j = 0, 1, \dots, J$ decline with J . Here we assume Condition S (20) for all of our models (roughly, all shares go down like $1/J$). Second, the rate at which the a_j 's grow or decline with J . Finally, the variability of the sequence $\{a_j\}$ also has a role to play in some cases. The examples establish rates of growth on $n(J), R(J)$, which insure, through their effects on the rate of growth and variability of the sequence $\{a_1, \dots, a_J\}$, finite limits for $\Phi_q(J); q = 2, 3$.

In general, if for some function $g(\cdot)$, we have $|a_j| \leq g(J)$ for $j = 1, \dots, J$, then for all J ,

$$\sum_{j=1}^J a_j^2 s_j^0 - \left(\sum_{j=1}^J a_j s_j^0 \right)^2 \leq \sum_{j=1}^J a_j^2 s_j^0 \leq \left(\max_{1 \leq j \leq J} |a_j| \right)^2 \sum_{j=1}^J s_j^0 \leq g(J)^2. \quad (32)$$

This gives a global bound on the variance matrix $\Phi_2(J)$; it is essentially this bound that was used in BLP to provide sufficient conditions for asymptotic normality.

However, it turns out that for two of our leading cases (the logit and random coefficient logit), there is further structure that can sometimes be exploited to give tighter bounds on $\Phi_2(J)$. Specifically, when (20) hold in these cases we have

$$(a_1, \dots, a_J) = g(J)\{(1, \dots, 1) + O(1/J)\}$$

for some non-decreasing function g [i.e., the normalized a 's have zero sample variability]. Then, we have

$$\begin{aligned} \sum_{j=1}^J a_j^2 s_j^0 - \left(\sum_{j=1}^J a_j s_j^0 \right)^2 &\simeq g(J)^2 \left[\sum_{j=1}^J s_j^0 - \left(\sum_{j=1}^J s_j^0 \right)^2 \right] \\ &= g(J)^2 \left[1 - s_0^0 - (1 - s_0^0)^2 \right] \\ &= g(J)^2 s_0^0 (1 - s_0^0). \end{aligned} \quad (33)$$

When (20) holds, the share of the outside alternative s_0^0 is $O(1/J)$, and so (33) is $O(g(J)^2/J)$, and we get a reduction in the magnitude of the variance from the crude bound (32).¹¹

4 A Detailed Analysis of Our Examples

Section 2 introduced two examples and we now provide a detailed analysis of both of them. The first was the logit model. As noted the logit has “diffuse” substitution patterns which, in turn, make estimators of the parameters of the model quite sensitive to sampling and simulation error. The fact that the simple logit only accommodates very restrictive substitution patterns has caused interest in the random coefficients logit model of BLP (1995), so

¹¹Note that when the share of the outside alternative is $O(1)$, (33) is the larger magnitude $O(g(J)^2)$. In this case, there is no gain and (32) is not improved.

we extend our results to that model also (analogous results hold for the nested logit, the multinomial probit, and the random coefficients probit).

The second example is the vertical model of Shaked and Sutton (1982). This and the horizontal model of Hotelling (1929) are uni-dimensional examples of a class of models Berry and Pakes (2002) call the pure characteristics model, and we consider in more detail below. In these models individual's preferences are defined on a finite dimensional space of product characteristics, and substitution patterns are "local" in the sense that cross price and characteristic elasticities are only non-zero for a finite number of products.

The asymptotic behavior of the estimator of the two model's parameters differ. In the first we require that the variance in both the simulation and the sampling error must decline at a rate faster than J increases for consistency and at the rate J^2 for asymptotic normality. For the second, the variance in the sampling and the simulation error can decline at any rate for consistency and to decline at rate J for asymptotic normality.

4.1 The logit model

Recall from equations (13) and (14) that the market shares predicted by the logit model are

$$\sigma_j(x, \xi, \theta) = \frac{e^{x'_j \theta + \xi_j}}{1 + \sum_{k=1}^J e^{x'_k \theta + \xi_k}}, \quad j = 1, \dots, J \quad \text{while} \quad \sigma_0(x, \xi, \theta) = \frac{1}{(1 + \sum_{k=1}^J e^{x'_k \theta + \xi_k})},$$

and from equation (16)

$$\frac{\partial \sigma}{\partial \xi} \equiv H(\theta, s, P) = S - ss', \quad \text{while} \quad H(\theta, s, P)^{-1} = S^{-1} + ii'/s_0,$$

where $S = \text{diag}[s]$ and $i = (1, \dots, 1)'$. $H(\theta, s, P)$ is the $J \times J$ share matrix derivative evaluated at $\xi = \xi(\theta, s, P)$, and does not depend on the parameter vector θ .

We assume that the random variables $x'_j \theta + \xi_j$ have bounded support and density bounded away from zero on this support.¹² This implies market shares are all of magnitude $O(1/J)$ with probability one (our Condition S, equation 20). Note that this is sufficient to ensure invertibility of the matrix H for every finite J .

From section 3.1 condition S implies A3 provided $J^{1+\epsilon}/n \rightarrow 0$ for some $\epsilon > 0$, and we simply assume that the instruments are stochastically bounded thus satisfying A4. Using $\tau_\alpha(x) = \log(x/\sigma_0)$, it is easy to see that A5 is also satisfied. Thus to prove consistency we need only verify the identification condition in A6

A sufficient condition for A6 is that for each $\epsilon > 0$ there is a $J(\epsilon)$ such that for any $J > J(\epsilon)$, $J^{-1} \sum_j z_j x'_j$ has full column rank with probability $1 - \epsilon$, since then

$$\inf_{\theta \notin \mathcal{N}_{\theta^0}(\delta)} \|G_J(\theta, s^0) - G_J(\theta^0, s^0)\| = \inf_{\theta \notin \mathcal{N}_{\theta^0}(\delta)} \left\| \left(\frac{1}{J} \sum_{j=1}^J z_j x'_j \right) (\theta - \theta^0) \right\| \geq \inf_{\theta \notin \mathcal{N}_{\theta^0}(\delta)} C \|\theta - \theta^0\| \geq C\delta,$$

with probability $1 - \epsilon$. In terms of the pricing problem this requires that the price of a product not be a linear function of that product's demand side attributes. However, we

¹²If instead we assumed $x_j \theta + \xi_j$ has finite variance, it would only affect the argument for normality through the remainder term magnitudes. That is a different argument would be required to insure that the rates of growth of $n(J)$ and $R(J)$ are large enough to cause the remainder terms to converge to zero.

know that the solution to the pricing problem generates a pricing function which depends on competitor characteristics and factor prices, as well as on its own characteristics.

We now move on to the conditions needed for asymptotic normality; in particular B4 (or equations 28 and 29) when condition S is satisfied. Without loss of generality assume z is a vector, and recall that to prove B4 it suffices to find a rate of growth for n that makes the limit, as J grows large, of $\Phi_2(J)$ finite (element by element), where

$$\Phi_2(J) = \frac{1}{nJ} \left[\sum_{j=1}^J a_j^2 s_j^0 - \left(\sum_{j=1}^J a_j s_j^0 \right)^2 \right] \quad \text{and} \quad a_k = z' H_0^{-1} e_k.$$

The formula for H_0^{-1} , and condition S (i.e., all $s_j > \underline{c}/J$) implies

$$a_k = \frac{z_k}{s_k} + \frac{\sum_{j=1}^J z_j}{s_0} = \frac{J^2 \bar{z}_J}{\underline{c}} [1 + O(1/J)], \quad (34)$$

where \bar{z}_J is the sample mean of z , which is bounded by assumption. From the discussion at the end of last section if

$$(a_1, \dots, a_J) = g(J)[(1, \dots, 1) + O(1/J)]$$

then the components of $\Phi_2(J)$ are $O_p[g(J)^2/J^2n]$. Equation (34) implies we satisfy this condition with $g(J) = J^2 O_p(1)$. Thus the components of $\Phi_2(J)$ are $O_p(J^2/n)$; i.e., n must grow like J^2 for asymptotic normality.

We now verify (29). Note that $|\sum_{j=1}^J a_j \varepsilon_{ji}| \leq \max_{1 \leq j \leq J} |a_j| \sum_{j=1}^J |\varepsilon_{ji}| \leq cJ^2$ for some constant c , because $\sum_{j=1}^J |\varepsilon_{ji}| \leq \sum_{j=1}^J [1(C_i = j) + E1(C_i = j)] \leq 2$. and (34) is true. Thus

$$E \left[\left| \frac{1}{n\sqrt{J}} \sum_{j=1}^J a_j \varepsilon_{ji} \right|^{2+\delta} \right] \leq \left(\frac{\bar{c}J^2}{n\sqrt{J}} \right)^{2+\delta}$$

for any δ . Thus $nE[|Y_{Ji}|^{2+\delta}] = O(J^{3+3\delta/2}n^{-(1+\delta)}) = o(1)$, which, after substituting $n(J) = J^2$, satisfies our condition provided $3 + 3\delta/2 - 2(1 + \delta) < 0$. That is condition (29) is satisfied for any $\delta > 2$.

Finally, we turn to the stochastic equicontinuity condition B5. In the logit case, there is no simulation, i.e., P is known exactly, and there is only the sampling error to consider. Furthermore, since $[\partial\sigma/\partial\xi]^{-1} = \partial\xi/\partial\sigma$, where $\xi_j(\theta, s, P^0) = \ln(s_j/(1 - \sum_{k=1}^J s_k)) - x_j\theta$, we can equivalently restrict our attention to the process for s , i.e.

$$\nu_J(s) = \frac{1}{J} z' H^{-1}(\theta, s, P) \varepsilon^n,$$

where $H^{-1}(\theta, s, P) = S^{-1} - ii'/(1 - i's)$. We must show that this process is stochastically equicontinuous. In the appendix we show the equivalent condition that

$$\nu_J(s^n) - \nu_J(s^0) = O_p(J^{3/2}/n). \quad (35)$$

It follows that the remainder terms in the expansion (11) are of smaller order than the leading terms.

In conclusion, the asymptotic variance of $\sqrt{J}(\hat{\theta} - \theta_0)$ is

$$\left(\lim_{J \rightarrow \infty} \frac{1}{J} \sum_{j=1}^J E(z_j x'_j) \right)^{-1} (\Phi_1 + \Phi_2) \left(\lim_{J \rightarrow \infty} \frac{1}{J} \sum_{j=1}^J E(x_j z'_j) \right)^{-1}$$

where

$$\Phi_2 = \lim_{J \rightarrow \infty} \frac{1}{nJ} z' H_0^{-1} z = \lim_{J \rightarrow \infty} \left[\frac{1}{nJ} \sum_{j=1}^J z_j z'_j s_j^{-1} + \frac{J \frac{1}{J} \sum_{j=1}^J E(z_j) \frac{1}{J} \sum_{j=1}^J E(z'_j)}{n s_0} \right].$$

Under our assumptions the first term is $O_p(J/n)$. Since the second term is $O_p(J^2/n)$ it is dominant.¹³ That is

$$\Phi_2 = \lim_{J \rightarrow \infty} \frac{J^2}{n} \times \frac{\lim_{J \rightarrow \infty} \frac{1}{J} \sum_{j=1}^J E(z_j) \frac{1}{J} \sum_{j=1}^J E(z'_j)}{\lim_{J \rightarrow \infty} (J s_0)}.$$

4.1.1 The Random Coefficients Logit

The logit model limits substitution patterns in unrealistic ways. However, the random coefficients logit, given by

$$u_{ij} = \delta_j + x_j \lambda_i + \epsilon_{ij}, \text{ where } \delta_j = x_j \beta + \xi_j,$$

is much more flexible (see BLP (1995)). Usually the δ_j contain the mean of the coefficients of the x and the λ contain the deviations from that mean. It is the variance in preferences for characteristics (the variance in the λ) that is essential for more flexible substitution patterns (especially the variance in the price coefficient).

The market share for this model is given by

$$\sigma_j(x, \xi, \theta) = \int \frac{e^{\delta_j + x_j \lambda}}{1 + \sum_k e^{\delta_k + x_k \lambda}} dP(\lambda) \equiv \int s_j(\lambda) dP(\lambda) \equiv E[s_j(\lambda)], \quad (36)$$

where P is a given probability measure. Note that the integrand, $s_j(\lambda)$, is just the logit market share function evaluated at a particular value of the random coefficients (we have suppressed its other arguments x, θ, ξ). The derivatives of the market share function are

$$\frac{\partial \sigma_j}{\partial \xi_k} = \begin{cases} \int s_j(\lambda) \{1 - s_j(\lambda)\} dP(\lambda) & j = k \\ - \int s_j(\lambda) s_k(\lambda) dP(\lambda) & \text{if } k \neq j. \end{cases}$$

In matrix terms we can write the share matrix

$$H = E[\mathcal{H}(\lambda)], \text{ where } \mathcal{H}(\lambda) = S(\lambda) - s(\lambda)s(\lambda)'$$

in which $S(\lambda) = \text{diag}(s_1(\lambda), \dots, s_J(\lambda))'$ and $s(\lambda) = (s_1(\lambda), \dots, s_J(\lambda))'$. Unfortunately there is no easy expression (that we know of) for the inverse matrix H_0^{-1} for this case. However, we can still characterize its properties sufficiently well to ensure that property (33) holds.

¹³In the case where the outside alternative is $O(1)$, the two terms in this equation are of equal magnitude and we must include both.

By the convexity of the matrix inverse [Groves and Rothenberg (1969)] we have

$$H^{-1} = [E\mathcal{H}(\lambda)]^{-1} \leq E[\mathcal{H}(\lambda)^{-1}]$$

in the positive definite sense. The inverse of any given logit matrix is $S(\lambda)^{-1} + ii'/s_0(\lambda)$. If we assume that $s_j(\lambda) \geq \underline{s}_j$ for all $j = 0, 1, \dots, J$ for some nonrandom sequence of constants \underline{s}_j that obey condition S, then

$$[E\mathcal{H}(\lambda)]^{-1} \leq \underline{S}^{-1} + \frac{ii'}{\underline{S}_0} \equiv \overline{H^{-1}}, \quad (37)$$

where $\underline{S} = \text{diag}(\underline{s}_1, \dots, \underline{s}_J)'$. Furthermore, $H^{-1}V_2H^{-1} \leq \overline{H^{-1}}V_2\overline{H^{-1}}$ by the properties of positive definite symmetric matrices [Anderson (1984, Theorem A1.1)]. We can now apply the results from the previous subsection. Under condition S, the variance term (31) is of order J^2/n as in the fixed coefficient logit case. The remaining arguments of the previous subsection hold here too so that the condition for the central limit theorem is satisfied in the random coefficient case. In fact, we are able to prove in this case that

$$\Phi_2 = \lim_{J \rightarrow \infty} \frac{J^2}{n} \times \frac{\mu_z \mu_z'}{\lim_{J \rightarrow \infty} (J \int s_0(\lambda) dP(\lambda))} \quad (38)$$

$$\Phi_3 \leq \lim_{J \rightarrow \infty} \frac{J^2}{R} \times \frac{\mu_z \mu_z'}{\lim_{J \rightarrow \infty} (J \int s_0(\lambda) dP(\lambda))}. \quad (39)$$

where $\mu_z = \lim_{J \rightarrow \infty} \frac{1}{J} \sum_{j=1}^J E(z_j)$. We could provide more detailed formalizations of both the identification and stochastic equicontinuity conditions, but we really have nothing substantive to say that we have not already said in the context of the fixed coefficient logit model.

4.2 The Vertical Model

Recall from equation (18) that the market shares for the vertical model are given by

$$s_0 = 1 - F(\Delta_1; \theta), \quad s_j = F(\Delta_j; \theta) - F(\Delta_{j+1}; \theta), \quad \text{for } j = 1, \dots, J-1, \quad \text{and } s_J = F(\Delta_J; \theta),$$

where $\Delta_j = (\delta_j - \delta_{j-1})/(p_j - p_{j-1})$, and for all market shares to be positive we require $\Delta_1 > \Delta_2 > \dots$ and $\delta_1 > \delta_2 > \dots$, where $\delta_j = x_j \beta + \xi_j$ for $j = 1, \dots, J$ ($\delta_0 = p_0 = 0$).

Since the simple vertical model only requires integration over one dimension of heterogeneity, we assume there is no simulation error. Further for this model the inversion from shares to ξ is obtained from the recursive system

$$\delta_j - \delta_{j-1} = (p_j - p_{j-1}) F^{-1} \left(1 - \sum_{r=1}^{j-1} s_r \right)$$

which, together with our normalization ($\xi_0 = 0$), implies

$$\xi_j(s^n) - \xi_j(s^0) = \sum_{l=1}^j (p_l - p_{l-1}) \left[F^{-1} \left(1 - \sum_{r=1}^{l-1} s_r^n \right) - F^{-1} \left(1 - \sum_{r=1}^{l-1} s_r^0 \right) \right].$$

Recall that one requirement for consistency is that $J^{-1}\|\xi(s^n) - \xi(s^0)\|^2 \rightarrow_p 0$.

For simplicity we assume that the distribution of λ (i.e., $F(\cdot)$) has bounded support and is strictly increasing (so its inverse satisfies a Lipschitz condition), and that whatever equilibrium is established $\max_{j \leq J} (p_j - p_{j-1}) = c < \infty$. Then for any $\epsilon > 0$

$$\begin{aligned} \Pr \left[\frac{1}{J} \sum_{j=1}^J \{ \xi_j(s^n) - \xi_j(s^0) \}^2 > \epsilon \right] &\leq \max_{j \leq J} \Pr \left[\{ \xi_j(s^n) - \xi_j(s^0) \}^2 > \epsilon \right] \\ &\leq J \max_{j \leq J} \Pr \left[\left\{ \sum_{l=0}^{j-1} s_l^n - \sum_{l=0}^{j-1} s_l^0 \right\}^2 > \epsilon/c \right] \\ &\leq J \exp(-\epsilon n/c), \end{aligned}$$

by Bernstein's inequality (since $\sum_{l=0}^{j-1} s_l^n$ is a sum of n independent random variables each bounded by one). Thus assumption A3 will be satisfied provided $n \rightarrow \infty$ faster than $\log J$.

For the asymptotic normality result we need the elements of the matrix H^{-1} , where $H = \partial\sigma/\partial\xi$. Letting $\alpha_1 = f(\Delta_1)/p_1$, $\alpha_2 = f(\Delta_2)/(p_2 - p_1), \dots, \alpha_J = f(\Delta_J)/(p_J - p_{J-1})$

$$H = \begin{pmatrix} \alpha_1 + \alpha_2 & -\alpha_2 & 0 & \dots & 0 \\ -\alpha_2 & \alpha_2 + \alpha_3 & \ddots & 0 & 0 \\ 0 & \ddots & \ddots & -\alpha_{J-1} & 0 \\ \vdots & 0 & -\alpha_{J-1} & \alpha_J + \alpha_{J-1} & -\alpha_J \\ 0 & 0 & 0 & -\alpha_J & \alpha_J \end{pmatrix}. \quad (40)$$

The matrix H is a band matrix with all elements more than one place from the diagonal being zero. Note also that all row and column sums are zero apart from the first row and column, and so the matrix is not diagonal dominant. Furthermore, it can be verified that

$$H^{-1} = \left(\sum_{r=1}^{\min(i,j)} \frac{1}{\alpha_r} \right)_{i,j} = \begin{pmatrix} \frac{1}{\alpha_1} & \frac{1}{\alpha_1} & \frac{1}{\alpha_1} & \dots & \frac{1}{\alpha_1} \\ \frac{\alpha_1}{\alpha_1} & \frac{\alpha_1}{\alpha_1} + \frac{1}{\alpha_2} & \frac{\alpha_1}{\alpha_1} + \frac{1}{\alpha_2} & \dots & \frac{\alpha_1}{\alpha_1} + \frac{1}{\alpha_2} \\ \frac{1}{\alpha_1} & \frac{1}{\alpha_1} + \frac{1}{\alpha_2} & \frac{1}{\alpha_1} + \frac{1}{\alpha_2} + \frac{1}{\alpha_3} & \dots & \frac{1}{\alpha_1} + \frac{1}{\alpha_2} + \frac{1}{\alpha_3} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{1}{\alpha_1} & \frac{1}{\alpha_1} + \frac{1}{\alpha_2} & \dots & \dots & \frac{1}{\alpha_1} + \dots + \frac{1}{\alpha_J} \end{pmatrix}.$$

Notice that any fixed element of the inverse matrix is of order one as $J \rightarrow \infty$ (this is in contrast to the logit models where the individual elements of the inverse were all of order J).

Assume that the z are bounded. Then, for $k = 1, \dots, J$,

$$\begin{aligned} a_k \equiv z' H^{-1} e_k &\leq \max_{1 \leq l \leq J} |z_l| \times \left[J \sum_{\ell=1}^k \frac{1}{\alpha_\ell} + \sum_{j=1}^{k-1} \left(\sum_{\ell=1}^j \frac{1}{\alpha_\ell} - \sum_{\ell=1}^k \frac{1}{\alpha_\ell} \right) \right] \\ &= \max_{1 \leq l \leq J} |z_l| \times \left[J \sum_{\ell=1}^k \frac{p_\ell - p_{\ell-1}}{f(\Delta_\ell)} - \sum_{j=1}^{k-1} \sum_{\ell=j+1}^k \frac{p_\ell - p_{\ell-1}}{f(\Delta_\ell)} \right], \end{aligned}$$

which gives the individual elements in the vector $z'H^{-1}$. Since prices increase in the order of the products

$$a_k \leq \max_{1 \leq l \leq J} |z_l| \times J \sum_{\ell=1}^k \frac{p_\ell - p_{\ell-1}}{f(\Delta_\ell)}, \quad (41)$$

which is of order J for any fixed k .

For Theorems 1 and 2 we must determine the magnitude of the sample variance of the sequence (a_1, \dots, a_J) with respect to the multinomial (s_1, \dots, s_J) or equivalently given our assumptions the multinomial $(1/J, \dots, 1/J)$. We have

$$\begin{aligned} \frac{1}{J} \sum_{k=1}^J a_k^2 - \left(\frac{1}{J} \sum_{k=1}^J a_k \right)^2 &\leq \frac{1}{J} \sum_{k=1}^J a_k^2 \\ &\leq J^2 \left(\max_{1 \leq l \leq J} |z_l| \right)^2 \times \frac{1}{J} \sum_{k=1}^J \left(\sum_{\ell=1}^k \frac{p_\ell - p_{\ell-1}}{f(\Delta_\ell)} \right)^2 \\ &\leq J^2 \left(\max_{1 \leq l \leq J} |z_l| \right)^2 \times \frac{\frac{1}{J} \sum_{k=1}^J \left(\sum_{\ell=1}^k (p_\ell - p_{\ell-1}) \right)^2}{\{\min_{1 \leq \ell \leq J} f(\Delta_\ell)\}^2} \\ &\leq J^2 \left(\max_{1 \leq l \leq J} |z_l| \right)^2 \times \underline{m}^{-2} \frac{1}{J} \sum_{k=1}^J p_k^2, \end{aligned}$$

with $\underline{m}^{-2} \frac{1}{J} \sum_{k=1}^J p_k^2$ finite (since we assume $\frac{1}{J} \sum_{k=1}^J p_k^2$ is finite and $\min_{1 \leq \ell \leq J} f(\Delta_\ell) \geq \underline{m} > 0$).¹⁴

Given the lower bound \underline{m} and that the price sequence has a finite second moment, (32) holds with $g(J) = J$. Therefore, the covariance matrix $\Phi_2(J)$ is of order J/n . That is, in this case, we obtain consistency if n increases at any rate faster than $\log J$, while the asymptotic normality result holds with all three terms contributing provided n grows like J .

Note the contrast to the logit-type models, where n must increase at rate J for consistency and rate J^2 for the asymptotic normality result [when all shares are the same magnitude]. The difference between the models is due to the difference between localized and diffuse competition. In the models with sampling and simulation errors, the derivative of market share with respect to product quality is declining at the same rate as the shares. Therefore, the elements of the inverse derivative matrix $(d\sigma/d\xi)^{-1}$ are growing in J , and the number of simulation draws must increase at a faster rate to offset this. In the vertical model, competition is localized and the derivative of market share with respect to product quality does not decline in J , and so the elements of the inverse derivative matrix stay bounded. As a result our limit theorems can suffice with a lower rate of growth for n in the vertical model.

¹⁴Again what actually happens to these quantities as J grows large will depend on appropriate specification of the pricing and product placement equilibria. We suffice with these assumptions because they seem to be sufficiently general and are all we require for our results. However a similar argument could be used to establish rates under other assumptions, and it may be possible to improve on our rates.

5 Monte Carlo Results

In this section we discuss Monte Carlo results for simple versions of our models. We start with logit-type models. In particular we present results for a simple logit where the market shares are observed with sampling error, and then for a random coefficients logit with simulation error in the computed shares. Next we turn to the pure characteristic models. Here we start with a simple vertical model where market shares are observed with sampling error and then move to a multi-dimensional pure characteristics model with simulation error. We conclude with a brief comparison of alternative estimators for the standard errors and of the quality of the normal limiting approximation to the true distribution of our parameter estimates.

The monte-carlo results reinforce the theoretical discussion in the previous sections. That is to obtain a “well-behaved” estimator for the first class of models sample sizes and simulation draws must be quite large and increase rapidly in J . The sample sizes and number of simulation draws which seem to be necessary for estimating the versions of the pure-characteristic model can be much smaller, and, do not increase nearly as rapidly in J . Also the normal limiting approximation seems to fit the true distribution of the parameter estimates surprisingly well.

All of our examples involve data on a single-cross section of markets. Thus, as in our theoretical discussion, our focus is on how estimates behave as the number of products (J) within a market varies. In practice there are several reasons to prefer to estimate off data that features a cross section or time series of different market equilibria but, as noted above, we will typically still need insights into how estimates behave as J changes.

For the logit model, the deterministic part of utility is drawn as

$$\delta_j = x_j\beta + \xi_j, \tag{42}$$

where ξ_j is drawn from the standard normal distribution. The x 's are a constant and a standard normal, with a β coefficient on the constant of 3 and a slope coefficient of 1. (Except as noted, all random variables in the Monte Carlo exercises are i.i.d. draws.)

Table 1 gives the mean estimated value of β_2 across 1000 Monte Carlo datasets. Each column gives results for a different value of J , the number of products, while the panels running down the table vary the number of consumer draws used to calculate the market share of the sample (n). Note that zero shares are discarded from the dataset. The fourth panel gives results for n set proportional to J , while in the fifth panel n is set equal to J^2 . The last column uses the true expected shares (i.e., “ $n = \infty$ ”).

In the second row of each panel is the simulated standard deviation (the standard error of the estimate across the simulated samples) and the third row gives the standard error of the mean (the simulated standard error divided by $\sqrt{1000}$). Apart from the inversion, the simple logit model is linear in parameters. Thus, given no sampling error in the shares, we should get unbiased results even for small J (which seems consistent with the results for $n = \infty$).

We see that the results are particularly bad for small n relative to J , with a large apparent bias. This is in large part due to the sample selection bias that comes from throwing small share products out of the market.¹⁵ A good with a low value of x will tend to have a positive

¹⁵We did not deal with this problem in our theoretical analysis above, but it is likely to be a problem for datasets built from small samples of consumers.

market share only if it has a large value of ξ while a good with a high value of x will tend to have a positive share even for small ξ . This generates a negative correlation between x and ξ among goods with positive market shares.

Table 2 gives Monte Carlo results for a random coefficients logit. In this case (as in most of the empirical literature which as aggregate data), we assume that observed market shares have no sampling error (i.e. we are assuming that the observed shares are aggregated over a very large number of consumers). Since we always simulate positive predicted shares, there is no sample selection problem when there is data from a large consumer sample. So here we also consider smaller values of R .

Our random coefficients logit example once again sets $\delta_j = x_j\beta + \xi_j$, but now $\beta = (-5, 1)$. Utility of consumer i for product j is

$$u_{ij} = \delta_j + \theta_x \lambda_i x_{j2} + \epsilon_{ij}, \quad (43)$$

where λ is standard normal, the standard deviation of the random “taste for x ”, θ_x , is set to one and x_{j2} is the non-constant element of x . As usual, the ϵ ’s are i.i.d. extreme-value draws. The market shares are calculated by taking R draws from the distribution of the random coefficient λ . The “observed” market shares are set to their expected value at the true parameter values (i.e., we are assuming that the observed shares are aggregated over a very large number of consumers.)

Computation of the inverse shares follows BLP, but we do not use a variance reduction (importance sampling) scheme of sort used in that paper.

Table 2 summarizes the estimates of θ_x , the standard deviation of the random coefficient on the non-constant x . The results are consistent with the theory that suggest that the estimation routine will perform badly when the number of simulation draws is “small” relative to the number of products. In particular there seems to be a bias that increases in J holding R fixed (at least at low values of R); and for $R = 50$ the variance is actually increases as J goes from 50 to 100. Also when J gets large (take our $J = 100$), we need fairly large value of R for that bias to go away (for $J=100$ we probably need $R > 1000$). However, as the theory predicts when we set $R = J^2$, all estimates look to be close to their true values (relative to their standard deviations), and the standard errors decrease as does \sqrt{J} .

Table 3 has results for the vertical model. As in Table 1, the variance in observed shares is generated by small samples of consumers rather than from simulation error in the predicted shares. Once again, this can produce zero observed market shares, but in the vertical model the zero share products can be included in the estimation routine at little cost.¹⁶

The exact vertical model considers a utility function of

$$u_{ij} = \delta - \theta_p \lambda_i p_j, \quad (44)$$

where δ is “quality”, λ_i is consumer-specific part of the the marginal disutility of a price increase and θ_p is a parameter of the model. To keep the random coefficient in an easy

¹⁶In practice, the inversion for δ simply sets the δ of zero share products to the δ of the next lowest-priced good. Since zero shares occur in the vertical model when δ ’s are “close together”, this creates little bias. Note the contrast to the logit model, where zero share products have systematically low δ s and where the inversion routine cannot handle zero shares. We should note that this is the choice of δ for a zero market share product produced by our estimation algorithm; any δ below this value would also be consistent with a zero market share.

one-parameter family, we assume that λ_i is drawn from the unit exponential distribution, so that θ_p (set equal to one in the experiments) is the mean disutility of a price increase. In fact, θ_p is not separately identified from demand-side data and so is held fixed at one in the Monte Carlo experiments (this is just a normalization.)

Quality is modeled as $\delta_j = x_j\beta + \xi_j$, where the two components of x are a constant and a uniform drawn from $(0, 2)$. β is set equal to $(1.5, 1)$. The “unobserved” ξ_j is uniform on $(-1, 1)$. To insure that the expected shares are all positive, price is set equal to δ^2 .¹⁷

The results in Table 3 summarize the estimates of β_2 , the slope coefficient on x in the quality equation. These results are *very different* than those for the logit-type models in Table 2. Indeed, when we use the vertical model it is striking that there is no apparent inconsistency in the estimates anywhere in the table (even when $J = 200$ and $n = 50$). As expected, for fixed n the variance decreases in J . However for small n the decrease is almost imperceptible, while with large n the variance declines at very close to the rate of \sqrt{J} , which is the rate we would expect if simulation had no impact on the estimates at all.

Our final monte carlo for parameter estimates is an example of the computation of δ in a multidimensional pure characteristics- model. To the vertical model of Table 3, we add a random coefficient on the observed x ,

$$u_{ij} = \delta + \theta_x \lambda_{i1} x_j - \theta_p \lambda_{i2} p_j, \quad (45)$$

with

$$\delta = \beta_0 + \beta_1 x_j + \xi. \quad (46)$$

There are now two dimensions of the unobserved consumer tastes, related to x and p .

In the vertical model of Table 3, computation was not an issue and so we focused on small consumer samples as the source of “simulation” error. In the pure random coefficients model, the market shares must be simulated and so we focus on simulation error. We assume that the consumer sample is very large. This is fairly realistic in many datasets and avoids the problem of sample selection that arises in the small sample case.

As discussed in Berry and Pakes (2002), estimation of the parameters is computationally cumbersome as the number of products increases, making it difficult to estimate the model with many repetitions in a Monte Carlo exercise. For computational tractability, we therefore focus on the computation of δ at fixed values of parameters. In particular, we hold the parameters (θ_x, θ_p) of the random coefficients at their true values and then compute δ from the constructed data on x , p and expected market shares. The method of solving for ξ is the “exact” homotopy method of Berry and Pakes (2002). To summarize the relevant error in the computation, we regress the computed δ on x to obtain an estimate of (β_0, β_1) . In Table 4, we report the mean estimates of β_1 we obtained from repeating this procedure for different number of products and simulation draws.

The data for Table 4 were created via the following assumptions. The observed x_j is drawn as 1.5 times a random uniform on $(0, 1)$. The unobserved ξ_j is drawn as a random uniform on $(0, 1)$. (Note that somewhat more of the variance in δ comes from x as oppose to ξ , which will aid the estimation procedure.) The term δ is then constructed via the parameters $(\beta_0, \beta_1) = (2, 1)$. To ensure positive market shares, price is set equal to a convex function of δ , $e^\delta/10$. The random “taste” for x is standard normal, while the random term on price is modeled as a standard log-normal (with $\mu = 1$.)

¹⁷In the vertical model, all shares will be positive if price increases “fast enough” in quality.

The results in Table 4 are consistent with our conjecture that the multi-dimensional pure characteristics model behaves much as the single-dimensional (vertical) model. In particular there is no obvious bias in the estimates even when there are only a small number of simulation draws. The “ ∞ ” row of Table 4 uses the true δ that created the data (as this is the δ that would be recovered if both n and R were infinite). This row therefore gives the results from the model without any simulation error. It is apparent that at low values of J and high values of R very little of the standard error of the estimate is attributable to simulation error, but that fraction is still quite large when $J = R$ (note that throughout we keep R fairly small as that keeps the computational burden of estimating the model repeatedly on different simulated data sets manageable). Overall, however, the table seems consistent with the conjecture that the multidimensional pure characteristics model behaves similar to the unidimensional characteristic model; in particular we do not need R to grow faster than J for consistency and fairly precise estimates can be obtained from relatively small values of R .

We also wanted to evaluate the quality of the approximation to the distribution of parameter estimates that emanates from the normal limiting distributions and the asymptotic variance formulas given in the text. To evaluate this we re-ran, with different simulation draws, the results in Table 2 (random coefficients logit) and Table 3 (pure vertical model) for $J = 50$ and $R = 500$. For each dataset, we calculated both the parameter estimate and the estimate of the variance of the asymptotic normal distribution. In both models, to calculate the asymptotic variance we need to estimate the variance in the moment conditions due to the simulation (as well as due to the data). For the random coefficients logit, we simulated this variance (for each of the 1000 Monte Carlo experiments) by drawing 500 new datasets and re-calculating the moment conditions at each new dataset (holding the parameter fixed at its estimated value). For the pure vertical model, the contribution of the error in the observed shares can be evaluated using the usual asymptotic formula for the variance of multinomial draws.

As can be seen in Table 5 the average of the estimates of the standard deviations from the asymptotic expansions is almost identical to the true standard deviation as obtained from the monte carlo estimates. The fits of the asymptotic approximations are surprisingly good. I.e. our limiting distributions seem to provide a very good approximation to the actual distributions of the parameter estimates. Figures 1 and 2 provide comparisons of the actual distribution of the parameter estimates from the Monte Carlo experiments to the distribution obtained from the limiting approximation using the variances and parameter estimates reported in the table (for the random coefficients logit and the vertical model respectively). The plotted distribution is a kernel density estimate, calculated from the default values in the Stata statistical program. The smooth line is the normal distribution fit to the same data; the distribution of the calculated parameters does appear to be approximately normal.

A Appendix

PROOF OF THEOREM 1. We first show that the estimator defined as any sequence that satisfies

$$\|G_J(\hat{\theta}, s^0, P^0)\| = \inf_{\theta \in \Theta} \|G_J(\theta, s^0, P^0)\| + o_p(1)$$

is consistent. Note that assumption A1 together with the law of large numbers for triangular arrays [see, for example, Billingsley (1986, Theorem 6.2)] imply that $\|G_J(\theta^0, s^0, P^0)\| = o_p(1)$. Therefore,

by Theorem 3.1 of Pakes and Pollard (1989) it will suffice to show that for every $(\delta, \varepsilon) > (0, 0)$ there exists a $C^*(\delta) > 0$ and an $J(\varepsilon)$ such that for $J \geq J(\varepsilon)$

$$\Pr \left[\inf_{\theta \notin \mathcal{N}_{\theta^0}(\delta)} \|G_J(\theta)\| \geq C^*(\delta) \right] \geq 1 - \varepsilon,$$

where we have omitted indexing G_J by (s^0, P^0) for notational convenience. From the triangle inequality $\inf_{\theta \notin \mathcal{N}_{\theta^0}(\delta)} \|G_J(\theta) - G_J(\theta^0)\| \geq C(\delta)$ implies that

$$\inf_{\theta \notin \mathcal{N}_{\theta^0}(\delta)} \|G_J(\theta)\| \geq C(\delta) - \|G_J(\theta^0)\|.$$

Fix $\varepsilon > 0$, and let $\varepsilon^* = \min\{\varepsilon, C(\delta)\}$, so that $0 < \varepsilon^* \leq \varepsilon$. Since $\|G_J(\theta^0)\| = o_p(1)$, there exists $J_1(\varepsilon^*)$ such that for any $J \geq J_1(\varepsilon^*)$, $\Pr\{\|G_J(\theta^0)\| \geq \varepsilon^*/2\} \leq \varepsilon^*/2$. By assumption A1, there exists $J_2(\varepsilon^*)$ such that for $J \geq J_2(\varepsilon^*)$, $\Pr\{\inf_{\theta \notin \mathcal{N}_{\theta^0}(\delta)} \|G_J(\theta) - G_J(\theta^0)\| \geq C(\delta)\} \geq 1 - \varepsilon^*/2$. Consequently, (2) implies that for $J \geq \max\{J_1(\varepsilon^*), J_2(\varepsilon^*)\}$

$$\Pr \left[\inf_{\theta \notin \mathcal{N}_{\theta^0}(\delta)} \|G_J(\theta)\| \geq C(\delta) - \varepsilon^*/2 \right] \geq 1 - \varepsilon^* \geq 1 - \varepsilon.$$

To complete the proof let $C^*(\delta) = C(\delta) - \varepsilon^*/2 > 0$.

We now return to the actual estimator $\hat{\theta}$ and show that

$$\|G_J(\hat{\theta}, s^n, P^R)\| = \inf_{\theta \in \Theta} \|G_J(\theta, s^0, P^0)\| + o_p(1). \quad (47)$$

We show that

$$\sup_{\theta \in \Theta} \frac{1}{J} \|\xi(\theta, s^n, P^R) - \xi(\theta, s^0, P^0)\|^2 = o_p(1), \quad (48)$$

which implies that

$$\begin{aligned} \sup_{\theta \in \Theta} \left\| \frac{1}{J} z' \{\xi(\theta, s^n, P^R) - \xi(\theta, s^0, P^0)\} \right\|^2 &\leq \frac{1}{J} \|z' z\|^2 \times \frac{1}{J} \sup_{\theta \in \Theta} \|\xi(\theta, s^n, P^R) - \xi(\theta, s^0, P^0)\|^2 \\ &= o_p(1), \end{aligned}$$

i.e., that $\sup_{\theta \in \Theta} \|G_J(\theta, s^n, P^R) - G_J(\theta, s^0, P^0)\| = o_p(1)$. This in turn implies (47) by the triangle inequality.

The result (48) follows from the following argument. We show below that

$$\sup_{\theta \in \Theta} \|\psi_J(\xi(\theta, s^n, P^R), \theta, s^0, P^0)\| = o_p(1). \quad (49)$$

Then, by Assumption A5: when $\|\xi - \xi(\theta, s^0, P^0)\| \geq \delta\sqrt{J}$, we have $\inf_{\theta \in \Theta} \|\psi_J(\xi, \theta, s^0, P^0)\| \geq \varepsilon$. This implies that $\|\xi(\theta, s^n, P^R) - \xi(\theta, s^0, P^0)\|^2/J = o_p(1)$ by contradiction, which concludes the proof of (48) and hence (47). The result (49) follows because:

$$\begin{aligned} \sup_{\theta \in \Theta} \|\psi_J(\xi(\theta, s^n, P^R), \theta, s^0, P^0)\| &\leq \sup_{\theta \in \Theta} \|\psi_J(\xi(\theta, s^n, P^R), \theta, s^0, P^0) - \psi_J(\xi(\theta, s^n, P^R), \theta, s^n, P^R)\| \\ &\leq \sup_{\theta \in \Theta} \sup_{\xi} \|\psi_J(\xi, \theta, s^0, P^0) - \psi_J(\xi, \theta, s^n, P^R)\| \\ &\leq \|\tau_J(s^n) - \tau_J(s^0)\| \\ &\quad + \sup_{\theta \in \Theta} \sup_{\xi} \|\tau_J(\sigma(\xi, s^0, P^0) + \varepsilon^R(\theta)) - \tau_J(\sigma(\xi, s^0, P^0))\|. \end{aligned}$$

For some intermediate values \bar{s}_j we have by the mean value theorem

$$\begin{aligned}
\|\tau_J(s^n) - \tau_J(s^0)\|^2 &= \frac{1}{J^\alpha} \sum_{j=1}^J \left[\dot{\tau}_\alpha(\bar{s}_j)(s_j^n - s_j^0) \right]^2 \\
&\leq \max_{1 \leq j \leq J} \left| (s_j^0)^\alpha \dot{\tau}_\alpha(\bar{s}_j) \right|^2 \frac{1}{J^\alpha} \sum_{j=1}^J \left[\frac{s_j^n - s_j^0}{(s_j^0)^\alpha} \right]^2 \\
&\leq \max_{1 \leq j \leq J} \left| (s_j^0)^\alpha \dot{\tau}_\alpha(\bar{s}_j) \right|^2 \times \frac{1}{nJ^\alpha} \sum_{j=1}^J \frac{s_j^0(1 - s_j^0)}{(s_j^0)^{2\alpha}} \times (1 + o_p(1)) \\
&= o_p(1),
\end{aligned}$$

by assumption A3, while $\max_{1 \leq j \leq J} |(s_j^0)^\alpha \dot{\tau}_\alpha(\bar{s}_j)| \leq M$ with probability tending to one by assumptions A3. This is because

$$\begin{aligned}
M &\geq \max_{1 \leq j \leq J} \left| (\bar{s}_j)^\alpha \dot{\tau}_1(\bar{s}_j) \right| \\
&= \max_{1 \leq j \leq J} \left| \{ (s_j^0)^\alpha + \bar{s}_j^\alpha - (s_j^0)^\alpha \} \dot{\tau}_\alpha(\bar{s}_j) \right| \\
&\geq \max_{1 \leq j \leq J} \left| (s_j^0)^\alpha \dot{\tau}_\alpha(\bar{s}_j) \right| - \max_{1 \leq j \leq J} \left| \frac{\bar{s}_j^\alpha - (s_j^0)^\alpha}{\bar{s}_j^\alpha} \right| \max_{1 \leq j \leq J} \left| \bar{s}_j^\alpha \dot{\tau}_\alpha(\bar{s}_j) \right| \\
&= \max_{1 \leq j \leq J} \left| (s_j^0)^\alpha \dot{\tau}_\alpha(\bar{s}_j) \right| + o_p(1),
\end{aligned}$$

where the $o_p(1)$ term follows from A3(a) and (c). The result

$$\sup_{\theta \in \Theta} \|\tau_J(\sigma(\xi, s^0, P^0) + \varepsilon^R(\theta)) - \tau_J(\sigma(\xi, s^0, P^0))\| = o_p(1)$$

follows by similar arguments using A3(b) and (d) ■

PROOF OF THEOREM 2. As discussed in section 3, this will follow from Pakes and Pollard (1989, Theorem 2) provided our remainder terms are $o_p(1)$ and the leading terms satisfy a central limit theorem.

Leading Term Argument. We show that

$$\left[\text{var} \left(c' \sqrt{J} \mathcal{G}_J(\theta^0) \right) \right]^{-1/2} c' \sqrt{J} \mathcal{G}_J(\theta^0) \tag{50}$$

is asymptotically normally distributed with mean zero and variance one for any vector c . Since the three terms in $\sqrt{J} \mathcal{G}_J(\theta^0)$, denoted T_{J1}, T_{J2} , and T_{J3} , say, are mutually independent it suffices to show that $\text{var} (c' T_{J\ell})^{-1/2} c' T_{J\ell}$, $\ell = 1, 2, 3$, converge to standard normal random variables. Then, by the Cramér-Wold device [the fact that a multivariate random variable is normal if any linear combination of its elements are], we have the result.

A standard central limit theorem for mutually uncorrelated random variables establishes that

$$(c' E\{\text{var}(\xi|z) z z'\} c)^{-1/2} c' J^{-1/2} z' \xi(\theta^0, s^0, P^0) \implies N(0, 1).$$

Condition (29) enables us to apply the Lyapunov central limit theorem for triangular arrays [see for example, Billingsley (1986, Theorem 27.3)], which says that the random variables $c' \sum_{i=1}^n Y_{Ji}$ and $c' \sum_{r=1}^R Y_{J,r}^*$ are asymptotically normal.

We now turn to the remainder terms. For each fixed θ , we use a Taylor series approximation to $\xi(\theta, s^n, P^R) - \xi(\theta, s^0, P^R)$ and to $\xi(\theta, s^0, P^R) - \xi(\theta, s^0, P^0)$. Specifically, by the intermediate value theorem

$$\begin{aligned} 0 &= \sigma(\xi(\theta, s^n, P^R), \theta, P^R) - s^n \\ &= \sigma(\xi(\theta, s^0, P^R), \theta, P^R) - s^n + \frac{\partial \sigma(\bar{\xi}, \theta, P^R)}{\partial \xi'} \{ \xi(\theta, s^n, P^R) - \xi(\theta, s^0, P^R) \}, \end{aligned} \quad (51)$$

where $\bar{\xi}$ is intermediate between $\xi(\theta, s^n, P^R)$ and $\xi(\theta, s^0, P^R)$. In fact, there are different vectors $\bar{\xi}$ for each row, but we suppress this for notational convenience. Thus using the facts that $\sigma(\xi(\theta, s^0, P^R), \theta, P^R) = s^0$ and that for any $\xi \in \mathcal{N}_{\xi^0}(\theta; \epsilon)$ the matrix $\partial \sigma(\xi, \theta, P^R) / \partial \xi'$ is invertible with probability tending to one, we can write

$$\xi(\theta, s^n, P^R) - \xi(\theta, s^0, P^R) = - \left\{ \frac{\partial \sigma(\bar{\xi}, \theta, P^R)}{\partial \xi'} \right\}^{-1} \epsilon^n \quad (52)$$

with probability tending to one. Likewise,

$$\begin{aligned} 0 &= \sigma(\xi(\theta, s^0, P^R), \theta, P^R) - s^0 \\ &= \sigma(\xi(\theta, s^0, P^0), \theta, P^R) - s^0 + \frac{\partial \sigma(\underline{\xi}, \theta, P^R)}{\partial \xi'} \{ \xi(\theta, s^0, P^R) - \xi(\theta, s^0, P^0) \}, \end{aligned}$$

where $\underline{\xi}$ are intermediate between $\xi(\theta, s^0, P^R)$ and $\xi(\theta, s^0, P^0)$ as before. Then we use the fact that $\sigma(\xi(\theta, s^0, P^0), \theta, P^R) - s^0 = \sigma(\xi(\theta, s^0, P^0), \theta, P^R) - \sigma(\xi(\theta, s^0, P^0), \theta, P^0) = \epsilon^R(\theta)$ to obtain that with probability tending to one

$$\xi(\theta, s^0, P^R) - \xi(\theta, s^0, P^0) = - \left\{ \frac{\partial \sigma(\underline{\xi}, \theta, P^R)}{\partial \xi'} \right\}^{-1} \epsilon^R(\theta). \quad (53)$$

Therefore,

$$\begin{aligned} \sqrt{J}[\mathcal{G}_J(\theta) - G_J(\theta, s^n, P^R)] &= -\frac{1}{\sqrt{J}} z' [H(\bar{\xi}, \theta, P^R)^{-1} - H(\theta, s^0, P^0)^{-1}] \epsilon^n \\ &\quad -\frac{1}{\sqrt{J}} z' [H(\underline{\xi}, \theta, P^R)^{-1} - H(\theta, s^0, P^0)^{-1}] \epsilon^R(\theta). \end{aligned} \quad (54)$$

We must establish that $\sqrt{J}[\mathcal{G}_J(\theta) - G_J(\theta, s^n, P^R)] = o_p(1)$ uniformly in θ in a shrinking neighborhood of θ^0 . We just show that

$$\sup_{\|\theta - \theta^0\| \leq \epsilon_J} \left\| \frac{1}{\sqrt{J}} z' \{ H(\bar{\xi}, \theta, P^R)^{-1} - H(\theta^0, s^0, P^0)^{-1} \} \epsilon^n \right\| = o_p(1), \quad (55)$$

from which the result follows. The proof for the term (54) is similar and is omitted. Since $\bar{\xi}$ is intermediate between $\xi(\theta, s^n, P^R)$ and $\xi(\theta, s^0, P^R)$ it is also consistent in mean square, i.e., there

exists a sequence $\epsilon_J \rightarrow 0$ such that $\Pr[\bar{\xi} \notin \mathcal{N}_{\xi^0}(\theta^0; \epsilon_J)] \rightarrow 0$. Furthermore, for this ϵ_J we have $\Pr\{\rho_P(P^R, P^0) \geq \epsilon_J\} \rightarrow 0$ by the Glivenko-Cantelli theorem. Then, notice that for any $\eta > 0$,

$$\begin{aligned}
& \Pr \left[\sup_{\|\theta - \theta^0\| \leq \epsilon_J} \left\| \frac{1}{\sqrt{J}} z' \{H(\bar{\xi}, \theta, P^R)^{-1} - H(\theta, s^0, P^0)^{-1}\} \epsilon^n \right\| > \eta \right] \\
& \leq \Pr \left[\sup_{\|\theta - \theta^0\| \leq \epsilon_J} \|\nu_J(\bar{\xi}, P^R, \theta) - \nu_J(\xi(s^0, P^0), P^0, \theta)\| > \eta \right] \\
& \leq \Pr \left[\sup_{\|\theta - \theta^0\| \leq \epsilon_J} \sup_{(\xi, P) \in \mathcal{N}_{\xi^0}(\theta^0; \epsilon_J) \times \mathcal{N}_{P^0}(\epsilon_J)} \|\nu_J(\xi, P, \theta) - \nu_J(\xi(s^0, P^0), P^0, \theta)\| > \eta \right] \\
& \quad + \Pr \left[\bar{\xi} \notin \mathcal{N}_{\xi^0}(\theta^0; \epsilon_J) \right] + \Pr \left[P^R \notin \mathcal{N}_{P^0}(\epsilon_J) \right] \\
& = \Pr \left[\sup_{\|\theta - \theta^0\| \leq \epsilon_J} \sup_{(\xi, P) \in \mathcal{N}_{\xi^0}(\theta^0; \epsilon_J) \times \mathcal{N}_{P^0}(\epsilon_J)} \|\nu_J(\xi, P, \theta) - \nu_J(\xi(s^0, P^0), P^0, \theta)\| > \eta \right] + o(1) \\
& = o(1)
\end{aligned}$$

by the stochastic equicontinuity condition B5. ■

PROOF OF (35). It suffices to show that for any random sequence $s(n)$ converging to s^0 we have $\|\nu_J(s(n)) - \nu_J(s^0)\| \rightarrow_p 0$. We shall take $s(n) = s^n$ and show that $R_n = \nu_J(s^n) - \nu_J(s^0) = o_p(1)$, where

$$R_n = \frac{1}{\sqrt{J}} z' \{(S^n)^{-1} - S^{-1}\} (s^n - s) + \frac{1}{\sqrt{J}} z' i i' (s^n - s) \left\{ \frac{1}{1 - i' s^n} - \frac{1}{1 - i' s} \right\} \equiv R_{n21} + R_{n22}.$$

The following argument shows that under our conditions $R_{n21} = O_p(J^{3/2}/n)$ and $R_{n22} = O_p(J^{3/2}/n)$. We deal first with R_{n21} , which can be rewritten using a geometric series expansion as

$$|R_{n21}| \leq \max \|z_\ell\| \times \frac{1}{\sqrt{J}} \sum_{\ell=1}^J \frac{\delta_\ell^2}{1 + \delta_\ell},$$

where $\delta_\ell = (s_\ell^n - s_\ell)/s_\ell$. For any $\epsilon > 0$,

$$\begin{aligned}
\Pr[|R_{n21}| > \epsilon] & \leq \Pr \left[|R_{n21}| > \epsilon \text{ and } \max_{1 \leq \ell \leq J} |\delta_\ell| \leq 1/2 \right] + \Pr \left[\max_{1 \leq \ell \leq J} |\delta_\ell| > 1/2 \right] \\
& \leq \Pr \left[|R_{n21}| > \epsilon \text{ and } \max_{1 \leq \ell \leq J} |\delta_\ell| \leq 1/2 \right] + o(1)
\end{aligned}$$

by the uniform convergence of δ_ℓ assumed in A3. When $\max_{1 \leq \ell \leq J} |\delta_\ell| \leq 1/2$, $|R_{n21}| \leq \frac{2}{\sqrt{J}} \sum_{\ell=1}^J \delta_\ell^2$, and by the Markov inequality

$$\begin{aligned} \Pr \left[\frac{2}{\sqrt{J}} \sum_{\ell=1}^J \delta_{\ell}^2 > \epsilon \right] &\leq \frac{\frac{2}{\sqrt{J}} \sum_{\ell=1}^J E(\delta_{\ell}^2)}{\epsilon} \\ &= \frac{\frac{2}{n\sqrt{J}} \sum_{\ell=1}^J \frac{(1-s_{\ell})}{s_{\ell}}}{\epsilon} = O(J^{3/2}/n). \end{aligned}$$

Similar calculation applies to R_{n22} . ■

PROOF OF (39) AND (38). We show that for any vector z ,

$$\frac{z'H^{-1}e_k}{J^2} = \frac{\mu_z}{\bar{s}_0} + O(1/J), \quad k = 1, \dots, J, \quad (56)$$

where $\mu_z = \lim_{J \rightarrow \infty} J^{-1} \sum_{j=1}^J z_j$ and $\bar{s}_0 = \lim_{J \rightarrow \infty} J \int s_0(\lambda) dP(\lambda)$. The variance formula then follows from (33). Note that the matrix Γ is the same as in the fixed coefficient logit case.

The proof of (56) is quite long because we can't directly calculate the inverse of H in this case. Instead we approximate the continuous mixture by a sequence of finite mixture, each of whose inverse we can compute. Let $T, T_J : \mathbf{P} \rightarrow \mathbb{R}$, where

$$T_J(P) = \frac{z'H(P)^{-1}e_k}{J^2} \quad ; \quad T(P) = \frac{\mu_z}{\bar{s}_0(P)},$$

where the notation $H(P)$ emphasizes the dependence of the matrix H on the probability measure P . We must show that for all $\epsilon > 0$, there exists J_0 such that for all $J \geq J_0$,

$$|T_J(P) - T(P)| < \epsilon.$$

We shall work with a discrete mixture of fixed coefficient models indexed by m . By the triangle inequality

$$\begin{aligned} |T_J(P) - T(P)| &\leq |T_J(P) - T_J(P_m)| + |T_J(P_m) - T(P_m)| + |T(P_m) - T(P)| \\ &= I + II + III \end{aligned}$$

for any m . The proof that III is small follows directly from our assumptions and the strong law of large numbers. We show below that II converges to zero uniformly in m, J . What remains is to show that I is small, which follows from the crude inequality

$$\frac{1}{J^2} |z'H(P)^{-1}e_k - z'H(P_m)^{-1}e_k| \leq \frac{1}{J^2} \|z'H(P_m)^{-1}\| \|H(P) - H(P_m)\| \|H(P)^{-1}e_k\| \quad (57)$$

and the following bounds (obtained below)

$$\|z'H(P_m)^{-1}\| \leq O(J^{5/2}) \quad (58)$$

$$\|H(P)^{-1}e_k\| \leq O(J^2) \quad (59)$$

$$\|H(P) - H(P_m)\| \leq O(1/m^{(1-\eta)/2} J^{1/2}), \quad (60)$$

provided $J^{4+\eta}/m \rightarrow 0$.

PROOF OF (59). Writing $H(P)^{-1} = C^{-1}B$, we have that $H(P)^{-1}e_k = \left(\frac{b_{1k}}{c_1}, \dots, \frac{b_{Jk}}{c_J}\right)$ whose (squared) norm is

$$\sum_{j=1}^J \frac{b_{jk}^2}{c_j^2} \leq \frac{1}{\min_{1 \leq j \leq J} c_j^2} \left(\sum_{j=1}^J b_{jk} \right)^2 \leq \frac{\text{const} \tan t}{J^2 \{1 - \Delta(J)\}^2} = O(J^4)$$

because the elements of B and C are known to be positive. This establishes (59). The verification of (58) is given below.

PROOF OF (60). Specifically, we show that the matrix $H(P) = \int S(\lambda)dP(\lambda) - \int s(\lambda)s(\lambda)'dP(\lambda)$ can be well approximated by the matrix $H(P_m) = \int S(\lambda)dP_m(\lambda) - \int s(\lambda)s(\lambda)'dP_m(\lambda)$, where P_m is an empirical distribution of size m from the population governed by P , that is,

$$H(P_m) = \frac{1}{m} \sum_{\ell=1}^m \{S(\lambda_\ell) - s(\lambda_\ell)s(\lambda_\ell)'\}.$$

We work element by element. Since $J s_j(\lambda)$ is bounded away from both zero and infinity, we have that for positive finite constants c_1 and c_2 ,

$$\begin{aligned} \Pr \left[\left| J^2 \int s_j(\lambda)s_k(\lambda) \{dP_m(\lambda) - dP(\lambda)\} \right| > \frac{\kappa}{m} \right] &\leq \exp[-2\kappa^2/mc_1] \\ \Pr \left[\left| J \int s_j(\lambda)(1 - s_j(\lambda)) \{dP_m(\lambda) - dP(\lambda)\} \right| > \frac{\kappa}{m} \right] &\leq \exp[-2\kappa^2/mc_2], \end{aligned}$$

by Hoeffding's exponential inequality, see Pollard (1984, p191). Therefore taking $\kappa = cm^{1/2}(\log m)^r$, we have by the Bonferroni inequality,

$$\begin{aligned} &\Pr \left[\max_{1 \leq j \neq k \leq J} \left| J^2 \int s_j(\lambda)s_k(\lambda) \{dP_m(\lambda) - dP(\lambda)\} \right| > \frac{c(\log m)^r}{m^{1/2}} \right] \\ &\leq \sum_{j \neq k} \Pr \left[\left| J^2 \int s_j(\lambda)s_k(\lambda) \{dP_m(\lambda) - dP(\lambda)\} \right| > \frac{c(\log m)^r}{m^{1/2}} \right] \\ &= O(J^2) \exp[-c^*(\log m)^{2r}] \end{aligned} \tag{61}$$

for some constant c^* . Taking $m = J^\alpha$ for any $\alpha > 0$, we get that

$$\sum_{m=1}^{\infty} \Pr \left[\max_{1 \leq j \neq k \leq J} \left| J^2 \int s_j(\lambda)s_k(\lambda) \{dP_m(\lambda) - dP(\lambda)\} \right| > \frac{c(\log m)^r}{m^{1/2}} \right] < \infty$$

provided $r > 3/2c^*\alpha$, so that by the Borel-Cantelli lemma, we have for any $\eta > 0$,

$$m^{(1-\eta)/2} \max_{1 \leq j \neq k \leq J} \left| J^2 \int s_j(\lambda)s_k(\lambda) \{dP_m(\lambda) - dP(\lambda)\} \right| \longrightarrow 0 \tag{62}$$

with probability one. Similarly,

$$m^{(1-\eta)/2} \max_{1 \leq j \leq J} \left| J \int s_j(\lambda)(1 - s_j(\lambda)) \{dP_m(\lambda) - dP(\lambda)\} \right| \longrightarrow 0 \tag{63}$$

with probability one. In conclusion, the discrete mixture of logits well approximates any random coefficient logit matrix. Specifically, (60) follows because

$$\begin{aligned} \|H(P) - H(P_m)\|^2 &= \sum_{j=1}^J \{H(P) - H(P_m)\}_{j,j}^2 + \sum_{j=1}^J \sum_{\substack{k=1 \\ j \neq k}}^J \{H(P) - H(P_m)\}_{j,k}^2 \\ &\leq J \max_{1 \leq j \leq J} \{H(P) - H(P_m)\}_{j,j}^2 + J^2 \max_{1 \leq j \neq k \leq J} \{H(P) - H(P_m)\}_{j,k}^2 \\ &= O(1 / \sqrt{Jm^{1-\eta}}) \end{aligned}$$

with probability one for large m, J by (62) and (63).

PROOF OF II. Consider the discrete mixture

$$H = \frac{1}{m} \sum_{\ell=1}^m (S^\ell - s^\ell s^{\ell'}),$$

where $s^\ell = (s_1^\ell, \dots, s_J^\ell)'$, $\ell = 1, \dots, m$. We show that

$$\frac{1}{J^2} z' H^{-1} e_k = \frac{\frac{1}{J} \sum_{j=1}^J z_j}{J \frac{1}{m} \sum_{\ell=1}^m s_0^\ell} + O(1/J), \quad k = 1, \dots, J, \quad (64)$$

where $s_0^\ell = 1 - \sum_{j=1}^J s_j^\ell = O(1/J)$, $\ell = 1, \dots, m$.

Write $H = (D + UV')/m$, where $D = \sum_{\ell=1}^m S^\ell$ and $U = (s^1, \dots, s^J)$ and $V = -(s^1, \dots, s^J)$. We have

$$z' H^{-1} e_k = m \{z' D^{-1} e_k - z' D^{-1} U (I + V' D^{-1} U)^{-1} V' D^{-1} e_k\} \quad (65)$$

by the Sherman-Morrison-Woodbury formula [Golub and Van Loan (1989, p51)]. First note that

$$z' D^{-1} e_k = \frac{z_k}{d_k} = O(J/m),$$

where $d_j = \sum_{\ell=1}^m s_j^\ell = O(m/J)$, $j = 1, \dots, J$, so this term is of smaller order. We are going to establish that

$$[(I + V' D^{-1} U)^{-1}]_{ij} = \frac{1 + O(1/J)}{\sum_{\ell=1}^m s_0^\ell [1 + O(1/J)]} \quad (66)$$

for all $i, j = 1, \dots, m$. In this case,

$$\frac{m}{J^2} z' D^{-1} U (I + V' D^{-1} U)^{-1} V' D^{-1} e_k = \frac{1}{J^2 \frac{1}{m} \sum_{\ell=1}^m s_0^\ell} z' D^{-1} U i i' V' D^{-1} e_k + O(1/J),$$

where $i' V' D^{-1} e_k = 1$ and $z' D^{-1} U i = \sum_{j=1}^J z_j$, so we get the required result (64).

We have

$$z' D^{-1} U_{1 \times m} = \left(\sum_{j=1}^J \frac{z_j s_j^1}{d_j}, \dots, \sum_{j=1}^J \frac{z_j s_j^m}{d_j} \right) \quad ; \quad V' D^{-1} e_k = - \begin{pmatrix} \frac{s_k^1}{d_k} \\ \vdots \\ \frac{s_k^m}{d_k} \end{pmatrix}$$

and

$$I + V'D^{-1}U = \begin{pmatrix} 1 - \sum_{j=1}^J \frac{(s_j^1)^2}{d_j} & -\sum_{j=1}^J \frac{s_j^1 s_j^2}{d_j} & \cdots & -\sum_{j=1}^J \frac{s_j^1 s_j^m}{d_j} \\ -\sum_{j=1}^J \frac{s_j^2 s_j^1}{d_j} & 1 - \sum_{j=1}^J \frac{(s_j^2)^2}{d_j} & & -\sum_{j=1}^J \frac{s_j^2 s_j^m}{d_j} \\ \vdots & & \ddots & \vdots \\ -\sum_{j=1}^J \frac{s_j^m s_j^1}{d_j} & -\sum_{j=1}^J \frac{s_j^m s_j^2}{d_j} & \cdots & 1 - \sum_{j=1}^J \frac{(s_j^m)^2}{d_j} \end{pmatrix}. \quad (67)$$

Substitute $s_j^m = d_j - \sum_{\ell=1}^{m-1} s_j^\ell$ and use the fact that $\sum_{j=1}^J s_j^\ell = 1 - s_0^\ell$, to obtain

$$\begin{aligned} \sum_{j=1}^J \frac{s_j^m s_j^k}{d_j} &= 1 - s_0^k - \sum_{\ell=1}^{m-1} \left(\sum_{j=1}^J \frac{s_j^k s_j^\ell}{d_j} \right) \equiv 1 - s_0^k - \frac{1}{m} \sum_{\ell=1}^{m-1} a_{\ell k} \\ \sum_{j=1}^J \frac{(s_j^m)^2}{d_j} &= \sum_{j=1}^J d_j + \sum_{\ell=1}^{m-1} \sum_{k=1}^{m-1} \sum_{j=1}^J \frac{s_j^\ell s_j^k}{d_j} - 2 \sum_{\ell=1}^{m-1} \sum_{j=1}^J s_j^\ell \\ &\equiv \sum_{\ell=1}^{m-1} \sum_{k=1}^{m-1} a_{\ell k} + \sum_{\ell=1}^{m-1} s_0^\ell - s_0^m - (m-2), \end{aligned}$$

where $a_{\ell k} = \sum_{j=1}^J \frac{s_j^k s_j^\ell}{d_j}$. Therefore, we can write

$$I + V'D^{-1}U = \begin{bmatrix} A & a \\ a' & b \end{bmatrix} + \frac{1}{J} \begin{bmatrix} 0_{m-1, m-1} & \delta \\ \delta' & \phi \end{bmatrix} = X + \frac{E}{J},$$

where the $m-1 \times m-1$ matrix A is

$$A = \begin{pmatrix} 1 - a_{11} & -a_{12} & \cdots & -a_{1, m-1} \\ -a_{12} & 1 - a_{22} & \cdots & -a_{2, m-1} \\ \vdots & & \ddots & \vdots \\ -a_{1, m-1} & -a_{2, m-1} & \cdots & 1 - a_{m-1, m-1} \end{pmatrix},$$

while the $m-1 \times 1$ column vectors

$$a = \begin{bmatrix} -\left\{1 - \sum_{\ell=1}^{m-1} a_{1\ell}\right\} \\ \vdots \\ -\left\{1 - \sum_{\ell=1}^{m-1} a_{m-1, \ell}\right\} \end{bmatrix}; \quad \delta = \begin{bmatrix} J s_0^1 \\ \vdots \\ J s_0^{m-1} \end{bmatrix},$$

and the scalars $b = (m-1) - \sum_{\ell=1}^{m-1} \sum_{k=1}^{m-1} a_{\ell k}$ and $\phi = J(-\sum_{\ell=1}^{m-1} s_0^\ell + s_0^m)$.

Note that the matrix $X = (x_{jk})$ is singular, in fact the last column (row) is equal to minus the sum of the preceding $m-1$ columns (rows). Therefore, by Taylor expansion

$$\det\left(X + \frac{E}{J}\right) = \frac{1}{J} \sum_{j,k=1}^m \frac{\partial \det(X)}{\partial x_{jk}} e_{jk} + \frac{1}{2J^2} \sum_{j,k,l,r=1}^m \frac{\partial^2 \det(X)}{\partial x_{jk} \partial x_{lr}} e_{jk} e_{lr} + \dots \quad (68)$$

First, we have that $\partial \det(X) / \partial x_{jk} = x_{jk}^{Adj}$, where x_{jk}^{Adj} is the adjoint [i.e., the determinant of the matrix X_{jk} formed by deleting the j 'th row and k 'th column from X , see Anderson (1984, p598)] of x_{jk} . In fact, for all j, k

$$x_{jk}^{Adj} = \det(A), \quad (69)$$

as we show below. Since most of the matrix $E = (e_{jk})$ is zero, we only need the adjoints corresponding to the outer (right) border of the matrix X , which means there are only order m terms in the first summation in (68). Also, note that

$$\frac{\partial^2 \det(X)}{\partial x_{mj} \partial x_{mk}} = \frac{\partial^2 \det(X)}{\partial x_{jm} \partial x_{km}} = 0 \quad j, k = 1, \dots, m,$$

so there are only order m^2 terms in the second summation. Furthermore, since

$$\frac{\partial^2 \det(X)}{\partial x_{mj} \partial x_{km}} = \det(A_{jk}) = O(\det(A)/m),$$

the second term in (68) is of order m/J^2 and

$$\det(I + V'D^{-1}U) = \det(A) \sum_{\ell=1}^m s_0^\ell [1 + O(1/J)]. \quad (70)$$

Finally, we must show that the adjoints of the matrix $Z = X + E/J$ satisfy

$$z_{jk}^{Adj} = \det(A)[1 + O(1/J)], \quad j \neq k, \quad (71)$$

which implies (66) holds. ■

PROOF OF (69). We use the fact that determinants are invariant to certain linear transformations and also that the matrix X has the following property

$$x_{jm} = - \sum_{\ell=1}^{m-1} x_{j\ell} \quad ; \quad x_{mk} = - \sum_{\ell=1}^{m-1} x_{m\ell}, \quad j, k = 1, \dots, m,$$

to show that the determinant of the matrix

$$X_{mj} = \begin{bmatrix} x_{11} & \cdots & x_{1,j-1} & x_{1,j+1} & \cdots & x_{1,m} \\ \vdots & & \vdots & \vdots & & \vdots \\ x_{m-1,1} & \cdots & x_{m-1,j-1} & x_{m-1,j+1} & \cdots & x_{m-1,m} \end{bmatrix}$$

is the same as the determinant of the matrix A . Specifically, add columns 1 to $m-2$ to the $m-1$ 'th column and one gets the matrix A . For general X_{jk} a sequence of such transformations gives the result. ■

PROOF OF (71). Essentially the same as above. ■

References

- [1] ANDERSON, T.W. (1984): *An introduction to multivariate analysis*. Wiley.
- [2] ANDREWS, D.W.K. (1994): "Empirical process methods in econometrics," in *The Handbook of Econometrics* volume 4, pp 2247-2294, Eds R.F. Engle and D.L. McFadden.
- [3] BERRY, S., (1994): "Estimating discrete choice models of product differentiation," *RAND Journal of Economics*, **25** , 242-262.
- [4] BERRY, S., LEVINSOHN, J., AND A. PAKES, (1995): "Automobile prices in market equilibrium," *Econometrica* **63** , 841-890.
- [5] BERRY, S., LEVINSOHN, J., AND A. PAKES, (2001): "Differentiated Products Demand Systems from a Combination of Micro and Macro Data: The New Car Market," *Cowles Foundation Working Paper #1337* Forthcoming in *Journal of Political Economy*.
- [6] BERRY, S. AND A. PAKES, (2002): "The Pure Characteristics Discrete Choice Model," *Harvard University Working Paper*
- [7] BILLINGSLEY, P. (1986): *Probability and Measure* . 2nd Edition. John Wiley. New York.
- [8] CAPLIN, A., AND B. NALEBUFF (1987): "Aggregation and Imperfect Competition: On the Existence of Equilibrium," *Econometrica* 59, 26-61..
- [9] CHAMBERLAIN, G. (1987): "Asymptotic efficiency in estimation with conditional moment restrictions," *Journal of Econometrics* 34, 305-334.
- [10] COURT, A.T. (1939): "Hedonic Price Indexes with Automotive Examples," *The Dynamics of Automobile Demand*, pp.99-117, the General Motors Corporation.
- [11] DIEUDONNE, J. (1969): *Foundations of Modern Analysis*. Academic Press: New York.
- [12] FERNHOLZ, L.(1983): *Von Mises Calculus for Statistical Functionals*. Lecture Notes in Statistics Vol 19. Springer Verlag, New York.
- [13] FIEDLER, M., (1986): *Special Matrices and their applications in numerical mathematics*. Kluwer: Dordrecht.
- [14] GOLUB, G.H., AND C.F. VAN LOAN (1989): *Matrix Computations*. The Johns Hopkins University Press: Baltimore.
- [15] GRILICHES, Z. (1961): "Hedonic Price Indices for Automobiles: An Econometric Analysis of Quality Change", *The Price Statistics of the Federal Government*, New York, the National Bureau of Economic Research.
- [16] GROVES, T., AND T.J. ROTHENBERG (1969). A note on the expected value of an inverse matrix. *Biometrika* , 690-691.
- [17] HANSEN, L. (1982): "Large Sample Properties of Generalized Method of Moments Estimators", *Econometrica* 50, 1029-1054.

- [18] HAUSMAN, J.A. AND D. WISE, (1978). "A Conditional Probit Model for Qualitative Choice: Discrete Decisions Recognizing Interdependence and Heterogeneous Preferences," *Econometrica* 46, 403-426.
- [19] HOTELLING, H. (1929): "Stability in Competition," *Economic Journal* 39, 41-57.
- [20] HOTZ, V.J. AND R. MILLER (1993): "Conditional choice probabilities and the estimation of dynamic models," *Review of Economic Studies* 60, 497-529.
- [21] MAMMEN, E. (1989): "Asymptotics with increasing dimension for robust regression with applications to the bootstrap," *The Annals of Statistics* 17, 382-400.
- [22] MARCUS, M. AND H. MINC (1964): *A survey of matrix theory and matrix inequalities*. Dover, N.Y.
- [23] NEVO, A. (2001). "Measuring Market Power in the Ready to Eat Cereal Industry," *Econometrica*, 69(2), 307-342.
- [24] NEWEY, W.K., (1990): "Efficient instrumental variables estimation of nonlinear models," *Econometrica* 58, 809-837.
- [25] NEWEY, W.K., (1993): "Efficient estimation of models with conditional moment restrictions", in *Handbook of Statistics, vol. 11*, G.S. Maddala, C.R. Rao, and H.D. Vinod eds., Amsterdam: North-Holland.
- [26] ORTEGA, J.M. AND W.C. RHEINBOLDT (1970): *Iterative solution of nonlinear equation in several variables*. Academic Press, London.
- [27] PAKES, A., (1992): "Dynamic Structural Models; Problems and Prospects", in J.J. Laffont and C. Sims (ed.) *Advances in Econometrics; Proceedings of the Sixth World Congress of the Econometric Society*, Cambridge University Press.
- [28] PAKES, A. AND S. BERRY AND J. LEVINSOHN, (1993): "Some Applications and Limitations of Recent Advances in Empirical Industrial Organization: Price Indexes and the Analysis of Environmental Change," *American Economic Review Papers and Proceedings*, 83, 240-246.
- [29] PAKES, A., AND S. OLLEY (1995): "A Limit Theorem for a Smooth Class of Semiparametric Estimators," *Journal of Econometrics* 65, 295-332.
- [30] PAKES, A., AND D. POLLARD (1989): "Simulation and the asymptotics of optimization estimators," *Econometrica* 57, 1027-1057.
- [31] POLLARD, D. (1984): *Convergence of Stochastic Processes*. Springer-Verlag. New York.
- [32] PORTNOY, S. (1984): "Asymptotic behavior of M-estimators of p regression parameters when p^2/n is large. I. Consistency," *The Annals of Statistics* 12, 1298-1309.
- [33] PORTNOY, S. (1985): "Asymptotic behavior of M-estimators of p regression parameters when p^2/n is large. II. Normal approximation," *The Annals of Statistics* 13, 403-417.

- [34] SHAKED, A. AND J. SUTTON (1982). “Relaxing price competition through product differentiation,” *The Review of Economic Studies*, 49, 3–13.
- [35] SHAKED, A. AND J. SUTTON (1990). “Multiproduct Firms and Market Structure,” *Rand Journal of Economics*, 21, 45-62.
- [36] SUTTON, J. (1991). *Sunk costs and Market Structure: Price Competition, Advertising, and the Evolution of Concentration*, MIT Press: Boston, MA.

Table 1:
Monte Carlo Estimates for the Simple Logit Model
True Value of the Parm is 1
1000 Monte Carlo Repetitions

# Consumer Draws (n)	# of Products (J)				
	10	25	50	100	200
500	0.941 (0.362) [0.011]	0.798 (0.209) [0.007]	0.778 (0.137) [0.068]	0.633 (0.086) [0.004]	0.518 (0.076) [0.002]
1000	0.997 (0.426) [0.014]	1.013 (0.255) [0.008]	0.974 (0.149) [0.005]	0.934 (0.120) [0.004]	0.882 (0.077) [0.002]
2000	1.023 (0.500) [0.016]	1.046 (0.224) [0.007]	0.998 (0.138) [0.004]	0.976 (0.123) [0.004]	0.923 (0.089) [0.004]
10J	0.685 (0.406) [0.013]	0.728 (0.214) [0.007]	0.768 (0.132) [0.004]	0.921 (0.110) [0.004]	0.916 (0.088) [0.004]
J ²	0.615 (0.358) [0.011]	0.857 (0.200) [0.006]	1.021 (0.139) [0.004]	1.022 (0.101) [0.003]	1.015 (0.077) [0.002]
∞	1.027 (0.376) [0.012]	0.997 (0.242) [0.008]	0.995 (0.133) [0.004]	1.007 (0.094) [0.003]	1.008 (0.073) [0.002]

Notes: Simulated Standard Errors (empirical standard deviations across the repetitions) in (·) and Simulated Standard Error of the Estimated Mean in [·].

Table 2:
Monte Carlo Estimates for the Random Coefficients Logit
True Value of the Parm is 1
1000 Monte Carlo Repetitions

# Simulation draws (R)	# of Products (J)		
	10	50	100
10	1.194 (0.982) [.031]	1.218 (0.512) [0.016]	*
50	1.025 (0.645) [0.020]	1.039 (0.311) [0.010]	1.241 (0.495) [0.016]
100	0.982 (0.674) [0.021]	1.013 (0.271) [0.009]	1.037 (0.209) [0.007]
500	0.998 (0.633) [0.002]	1.008 (0.255) [0.008]	1.015 (0.181) [0.006]
10J	0.982 (0.674) [0.014]	1.008 (0.255) [0.008]	1.018 (0.181) [0.006]
J ²	0.982 (0.674) [0.021]	.998 (0.244) [.008]	1.004 (0.175) [.006]

Notes: Simulated Standard Errors (empirical standard deviations across the repetitions) in (·) and Simulated Standard Error of the Estimated Mean in [·].

*With 100 products and only 10 draws, we had numeric problems computing the estimates.

Table 3:
Monte Carlo Estimates for the Pure Vertical Model
True Value of the Parm is 1
1000 Monte Carlo Repetitions

# Consumer Draws (n)	# of Products (J)				
	10	25	50	100	200
50	1.023 (0.494) [0.016]	1.022 (0.373) [0.012]	1.011 (0.349) [0.011]	0.997 (0.321) [0.010]	1.013 (0.302) [0.010]
100	1.005 (0.426) [0.014]	1.010 (0.303) [0.010]	1.005 (0.257) [0.008]	1.002 (0.244) [0.008]	1.009 (0.217) [0.007]
500	0.993 (0.371) [0.012]	0.998 (0.223) [0.007]	1.001 (0.176) [0.006]	1.005 (0.142) [0.005]	1.007 (0.123) [0.004]
1000	1.01 (0.361) [0.011]	0.99 (0.227) [0.007]	1.00 (0.162) [0.006]	1.00 (0.118) [0.004]	1.00 (0.097) [0.003]
$10J$	1.018 (0.440) [0.014]	1.014 (0.253) [0.008]	1.008 (0.175) [0.006]	0.998 (0.120) [0.004]	0.996 (0.085) [0.003]
J^2	0.998 (0.423) [0.014]	0.998 (0.227) [0.007]	1.000 (0.153) [0.005]	1.002 (0.105) [0.003]	1.000 (0.074) [0.002]
∞	0.997 (0.364) [0.011]	0.999 (0.214) [0.007]	0.999 (0.141) [0.005]	1.001 (0.101) [0.003]	0.997 (0.072) [0.002]

Notes: Simulated Standard Errors (empirical standard deviations across the repetitions) in (·) and Simulated Standard Error of the Estimated Mean in [·].

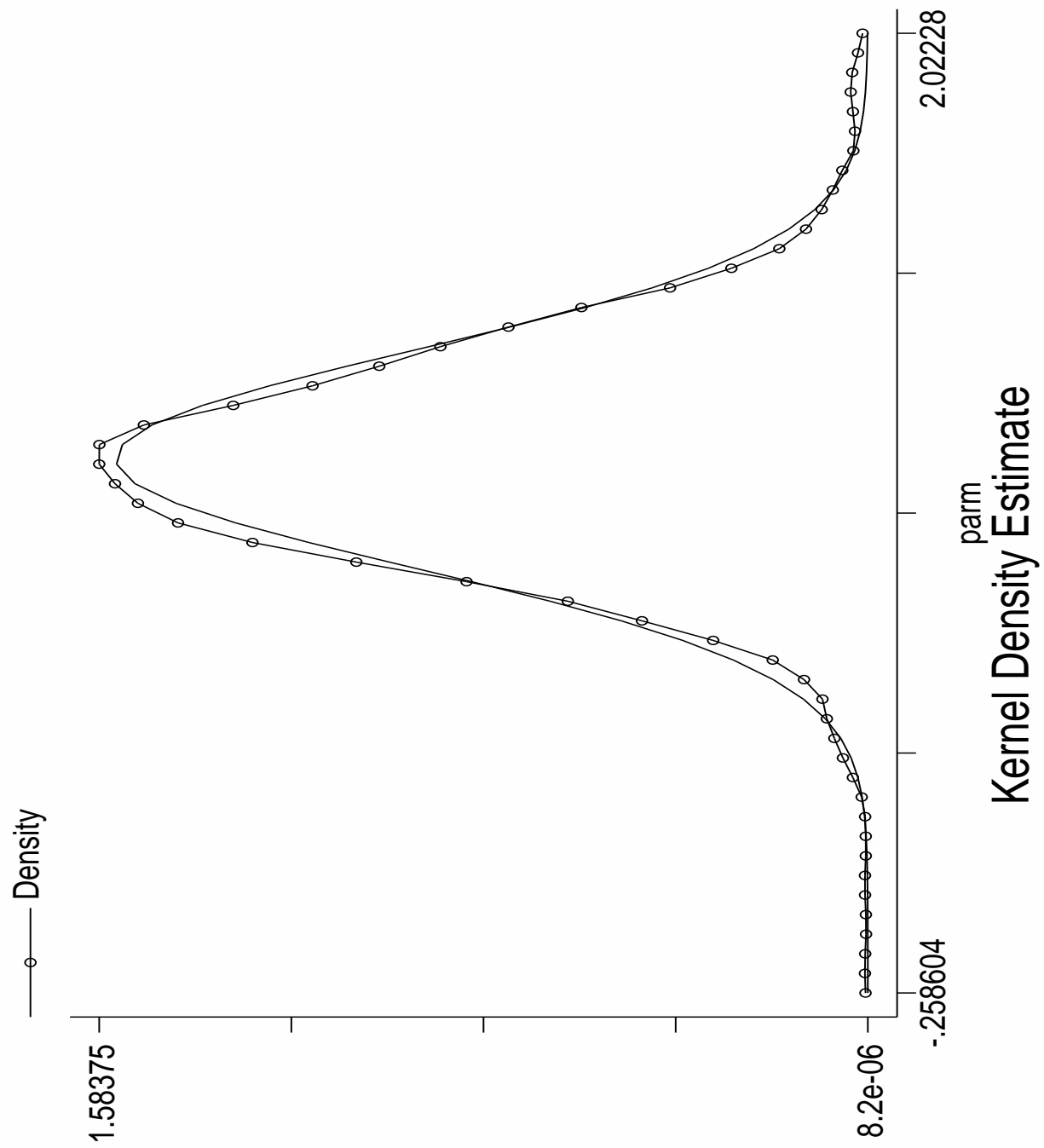
Table 4:
Monte Carlo Estimates for a Pure Characteristics Model
True Value of the Parm is 1
100 Monte Carlo Repetitions

# Simulation Draws (R)	# of Products, (J)			
	10	25	50	100
10	1.039 (0.370) [0.037]	0.999 (0.332) [0.033]	1.016 (0.311) [0.031]	1.021 (0.325) [0.033]
25	1.043 (0.279) [0.028]	0.993 (0.268) [0.027]	0.999 (0.235) [0.024]	1.010 (0.214) [0.021]
50	1.040 (0.243) [0.024]	1.006 (0.215) [0.021]	0.992 (0.187) [0.019]	1.024 (0.161) [0.016]
100	1.036 (0.224) [0.022]	1.023 (0.182) [0.018]	0.987 (0.143) [0.014]	1.012 (0.136) [0.014]
J	1.039 (0.370) [0.037]	0.993 (0.268) [0.027]	0.992 (0.187) [0.019]	1.012 (0.136) [0.014]
∞	1.030 (0.207) [0.021]	1.013 (0.164) [0.016]	0.986 (0.103) [0.010]	1.002 (0.061) [0.006]

Notes: Simulated Standard Errors (empirical standard deviations across the repetitions) in (\cdot) and Simulated Standard Error of the Estimated Mean in $[\cdot]$.

Table 5:
Monte Carlo and Estimated Standard Errors
($J = 50, R = 500, 1000$ Repetitions)

Model	Mean Parm	Monte Carlo Std. Dev.	Mean Asymp. Std. Dev.
R.C. Logit	1.010	0.2574	0.2201
Pure Vert	1.002	0.1720	0.1719



Density

Figure 1

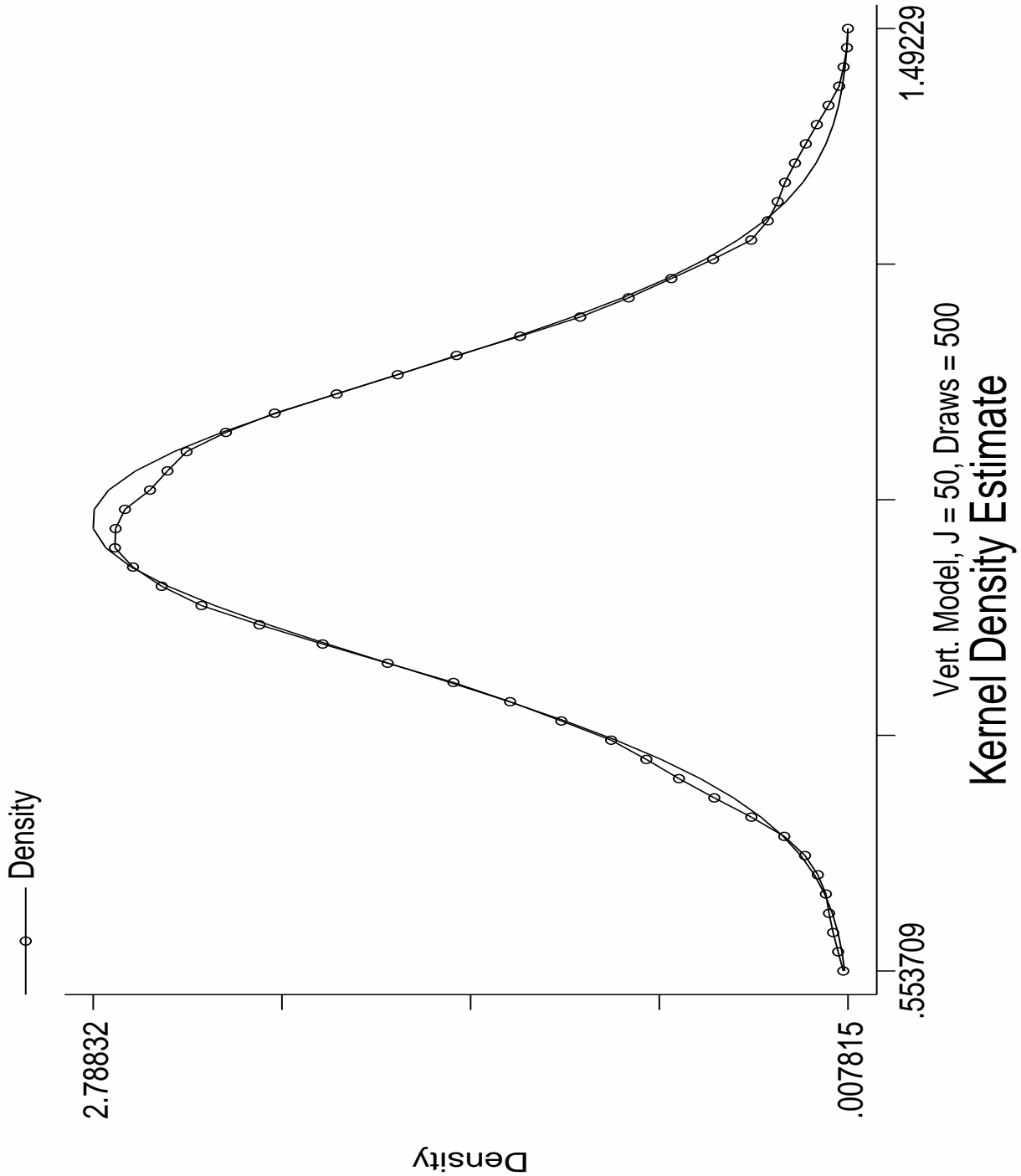


Figure 2