

# Best practices for differentiated products demand estimation with PyBLP

Christopher Conlon\*

and

Jeff Gortmaker\*\*

*Differentiated products demand systems are a workhorse for understanding the price effects of mergers, the value of new goods, and the contribution of products to seller networks. Berry, Levinsohn, and Pakes (1995) provide a flexible random coefficients logit model which accounts for the endogeneity of prices. This article reviews and combines several recent advances related to the estimation of BLP-type problems and implements an extensible generic interface via the PyBLP package. Monte Carlo experiments and replications suggest different conclusions than the prior literature: multiple local optima appear to be rare in well-identified problems; good performance is possible even in small samples, particularly when “optimal instruments” are employed along with supply-side restrictions. If Python is installed on your computer, PyBLP can be installed with the following command: `pip install pyblp`. Up-to-date documentation for the package is available at <https://pyblp.readthedocs.io>.*

## 1. Introduction

■ Empirical models of supply and demand for differentiated products are one of the most important achievements of the New Empirical Industrial Organization (NEIO) literature of the last 30 years. The workhorse model is the Berry, Levinsohn, and Pakes (1995) or BLP approach, which provides an estimator that allows for flexible substitution patterns across products, addresses the potential endogeneity of price, and also provides an algorithm for recovering that estimator. It has the advantage that it both scales well for large numbers of potential products, and can utilize both aggregated and dis-aggregated data. The BLP model and its variants have been used in a wide variety of applications: understanding the value of new goods (Petrin, 2002), evaluating the price effects of mergers (Nevo, 2001, 2000a), and studying two-sided markets (Fan, 2013; Lee, 2013). The BLP approach has been applied to a wide number of different questions and industries including hospital demand and negotiations with insurers (Ho and Pakes,

\* New York University, Stern School of Business; [cconlon@stern.nyu.edu](mailto:cconlon@stern.nyu.edu).

\*\* Harvard University; [jgortmaker@g.harvard.edu](mailto:jgortmaker@g.harvard.edu).

Thanks to Steve Berry, Jeremy Fox, Phil Haile, Mathias Reynaert, and Frank Verboven and seminar participants at NYU, Rochester, and the 2019 IIOC conference. Thanks to the editor Marc Rysman and to three anonymous referees. Daniel Stackman provided excellent research assistance. Any remaining errors are our own.

2014) and students' choice of schools (Bayer, Ferreira, and McMillan, 2007; Nielson, 2017). Moreover, the BLP approach has been extremely influential in the practice of prospective merger evaluation, particularly in recent years.

The model itself is both quite simple to understand, and quite challenging to estimate. At its core, it involves a nonlinear change of variables from the space of observed market shares to the space of mean utilities for products. After this nonlinear change of variables, the BLP problem is simply either a single linear instrumental variables (IV) regression problem (demand alone), or a two equation (supply and demand) linear IV regression problem. This means that a wide variety of tools for that class of problems are available to researchers.

The main disadvantage of the BLP estimator is that parameters governing the nonlinear change of variables are unknown. This results in a nonlinear, nonconvex optimization problem with a simulated (or approximated) objective function. The problem must be solved iteratively using nonlinear optimization software, and because of the nonconvexity, there is no mathematical guarantee that a solution will always be found. This has led to some frustration with the BLP approach (see Knittel and Metaxoglou, 2014). There is also the fear that when estimation is slow or complicated, researchers may cut corners in undesirable ways and sacrifice modeling richness for computational speed.

Despite its popularity, this literature lacks a standardized implementation that is sufficiently general to encompass a wide range of potential problems and use cases. Instead, nearly every researcher implements the estimator on their own with problem-specific tweaks and adjustments. This makes replication extremely challenging, and also makes it hard to evaluate different methodological and statistical improvements to the estimator.

The goal of this article is to present best practices for the estimation of BLP-type models, some of which are well-known in the literature, others of which are lesser known, and others still are novel to this article. In addition to presenting these best practices, we provide a common framework, PyBLP, which offers a general implementation of the BLP approach as a Python package.<sup>1</sup> We recommend installing PyBLP on top of an Anaconda distribution, which comes pre-packaged with PyBLP's primary dependencies.<sup>2</sup> Users of other languages such as MATLAB, Julia, and R can easily use PyBLP from their language of choice with packages that allow for between-language interoperability.<sup>3</sup> The PyBLP software is general, extensible, and open-source so that it can be modified and extended by scholars as new methodological improvements become available. The hope is that these best practices, along with this standardized and extensible software implementation, reduce some of the barriers to BLP-type estimation, making these techniques accessible to a wider range of researchers and facilitating replication of existing results.

This article and the accompanying package build upon a growing literature focused on methodological innovations and econometric properties of the BLP estimator. In Section 3, we discuss several such improvements and evaluate them with Monte Carlo studies in Section 5. Our objective is to compare practices in the literature and arrive at best practices suitable for a large number of use cases. We then implement these best practices as defaults in PyBLP. We organize the best practices around several of the tasks in the BLP estimator: solving the fixed point, optimization, integration, and solving counterfactual pricing equilibria.

---

<sup>1</sup> We chose Python over other languages because its popularity is growing, its package management systems are well-established, it has a mature scientific computing ecosystem, and as a general purpose language, it is conducive to the development of larger projects.

<sup>2</sup> PyBLP depends on standard packages in Python's scientific ecosystem: NumPy, SciPy, SymPy, and Patsy. It also depends on a companion package PyHDFE (Gortmaker and Tarascina, 2020), which implements algorithms for absorbing high-dimensional fixed effects.

<sup>3</sup> Python can be called from MATLAB with the `py` command, from Julia with `PyCall` (Johnson, 2019), and from R with `reticulate` (Allaire et al., 2017), which we give an example of in Section 6. Scientific computing in all of these high-level languages is backed by similar implementations of numerical linear algebra routines such as LAPACK.

In addition to best practices, we also provide some novel results. We provide a slightly different characterization of the BLP problem in Section 2 which facilitates estimation with simultaneous supply and demand restrictions. We show how this characterization can be made amenable to large numbers of fixed effects, and in Section 4 we characterize an approximation to optimal instruments in the spirit of Amemiya (1977) or Chamberlain (1987). Our characterization of the problem under optimal instruments allows us explore parametric identification with supply and demand in a way that parallels Berry and Haile (2014) and makes explicit cross equation and exclusion restrictions.

On the matter of instruments, our results generally coincide with and build on the existing literature. Gandhi and Houde (2019) construct what they refer to as *differentiation IV*, whereas Reynaert and Verboven (2014) evaluate a *feasible approximation to the optimal IV* in the sense of Amemiya (1977) or Chamberlain (1987).<sup>4</sup> We provide routines to construct both sets of instruments. Our results with respect to differentiation IV are mostly consistent with Gandhi and Houde (2019) in that they outperform other simple, yet commonly used forms of the *BLP instruments* (functions of exogenous product characteristics) such as sums or averages. Our results with respect to approximate optimal instruments are broadly similar to those of Reynaert and Verboven (2014) in that the performance gains under correctly specified supply models are substantial both in terms of bias and efficiency.

Our simulations indicate somewhat more positive results than previously observed in the literature when additional moments from correctly specified models of supply are also included. These findings are somewhat different from those of Reynaert and Verboven (2014), which suggest that once optimal demand-side instruments are included, the addition of a supply side has limited benefit.<sup>5</sup> The explanation for this phenomenon is directly related to our theoretical result in Section 4. Indeed, our simulations indicate that with both a correctly specified supply side and optimal instruments, the finite sample performance of the estimator is good even with relatively weak excluded cost-shifters, which supports the “folklore” around the original Berry, Levinsohn, and Pakes (1995) article: supply restrictions are valuable in improving the econometric performance of the estimator.

Employing best practices and the PyBLP software, we are able to revisit recent findings regarding methodological issues and innovations in BLP-type estimators in large-scale Monte Carlo experiments. Although many of our results confirm previous findings in the literature, we arrive at different conclusions on several occasions. The findings of Knittel and Metaxoglou (2014) suggest that the BLP problem is often characterized by a large number of local optima, and that these local optima can produce a wide range of potential elasticities and welfare effects. In contrast, our experience is that after implementing best practices, parameter estimates and counterfactual predictions are quite stable across starting values and choices of optimization software, both open-source and commercial. Likewise, Armstrong (2016) finds that absent strong cost shifting instruments, as the number of products increases, BLP instruments (characteristics of own and competing products) become weak, and it becomes difficult to reliably estimate demand parameters. Our findings suggest that with a correctly specified supply side and approximations to the optimal instruments, parameter estimates can still be estimated precisely. In general, we struggle to replicate some of the difficulties found in the previous literature, suggesting that the finite sample performance of BLP estimators may be better than previously thought.

There is also a recent literature of alternative approaches to BLP problems employing different algorithms or statistical estimators, which we do not directly address. This is not meant to suggest that there is anything wrong with these approaches, but merely that they are beyond the scope of this article. For example, Dubé, Fox, and Su (2012) propose an alternative estimation algorithm based on the mathematical programming with equilibrium constraints (MPEC)

---

<sup>4</sup> It is worth mentioning that the actual optimal IV are well-known to be infeasible. See Berry, Levinsohn, and Pakes (1995) or Berry, Levinsohn, and Pakes (1999).

<sup>5</sup> This is most likely because Reynaert and Verboven (2014) examine scenarios with strong cost-shifters.

TABLE 1 Model Notation

$j$	Products	$p_{jt}$	Price
$t$	Markets	$c_{jt}$	Marginal cost
$i$	"Individuals"	$x_{jt}$	Exogenous product characteristic
$f$	Firms	$v_{jt}$	Excluded demand-shifter
$h$	Nesting groups	$w_{jt}$	Excluded supply-shifter
$N$	Set/number of products across all markets	$U_{ijt}$	Indirect utility
$T$	Set/number of markets	$\delta_{jt}$	Mean utility
$J_t$	Set/number of products in market $t$	$\mu_{ijt}$	Heterogeneous utility
$I_t$	Set/number of individuals in market $t$	$\epsilon_{ijt}$	Idiosyncratic preference
$F_t$	Set/number of firms in market $t$	$d_{ijt}$	Choice indicator
$J_{ft}$	Set/number of products of firm $f$ in market $t$	$s_{ijt}$	Choice probability
$H$	Set/number of nesting groups	$s_{jt}$	Calculated market share
$J_{ht}$	Set/number of products in nest $h$ and market $t$	$S_{jt}$	Observed market share
		$\xi_{jt}$	Demand-side structural error
$\theta$	Parameters with dimension $K$	$\mathcal{H}_t$	Ownership or holdings matrix
$\theta_1$	Linear demand-side parameters with dimension $K_1$	$\Delta_t$	Intra-firm demand derivatives
$\theta_2$	Nonlinear common parameters with dimension $K_2$	$\eta_{jt}$	Multi-product Bertrand markup
$\theta_3$	Linear supply-side parameters with dimension $K_3$	$\omega_{jt}$	Supply-side structural error
$\beta$	$\theta_1$ excluding fixed effects	$Z$	Instruments
$\tilde{\theta}_2$	Parameters in $\theta_2$ governing heterogeneity	$W$	Weighting matrix
$\alpha$	Parameter in $\theta_2$ on price	$g$	Sample moments
$\gamma$	$\theta_3$ excluding fixed effects	$q$	Objective function
$\rho$	Nesting parameter in $\theta_2$		

method of Su and Judd (2012) and which Conlon (2017) extends to generalized empirical likelihood estimators. Although the MPEC approach has some advantages, we focus on the more popular nested fixed point problem. Lee and Seo (2015) provide an *approximate* BLP estimator, which is asymptotically similar to the BLP estimator though differs in finite sample. Hong, Li, and Li (Forthcoming) propose a Laplace-type estimator (LTE). Salanie and Wolak (2019) propose another approximate estimator that can be estimated with linear IV, and is helpful for constructing good starting values for optimization. Other common modifications to the BLP model that are beyond the scope of this article include the *pure characteristics* model of Berry and Pakes (2007) and the inclusion of *micro moments*, as in Petrin (2002) and Berry, Levinsohn, and Pakes (2004a).<sup>6</sup>

## 2. Model and estimation

■ In Table 1 we summarize the notation that we will introduce in this section. Bold font denotes a  $J_t \times 1$  vector for all products within a market. For example,  $\mathbf{s}_t$  and  $\mathbf{p}_t$  denote vectors of shares and prices for all products in market  $t$ . For clarity, we partition the parameter space  $\theta$  into three parts: the  $K_1 \times 1$  vector  $\theta_1$  contains the demand parameters  $\beta$ , the  $K_3 \times 1$  vector  $\theta_3$  contains the supply parameters  $\gamma$ , and the remaining parameters, including the price coefficient  $\alpha$  and parameters governing heterogeneous tastes  $\tilde{\theta}_2$ , are contained in the  $K_2 \times 1$  vector  $\theta_2$ .

□ **Demand.** Berry, Levinsohn, and Pakes (1995) begin with the following problem. An individual  $i$  in market  $t = 1, \dots, T$  receives indirect utility from selecting a particular product  $j$ :

$$U_{ijt} = \delta_{jt} + \mu_{ijt} + \epsilon_{ijt}. \quad (1)$$

<sup>6</sup> PyBLP supports an approximation to the pure characteristics model and common forms of micro moments, but we do not analyze econometric performance of these modifications in this article.

Consumers then choose among  $J_t = \{0, 1, \dots, J_t\}$  discrete alternatives—including the outside alternative, denoted  $j = 0$ , which gives  $U_{i0t} = \epsilon_{i0t}$  for all  $(i, t)$ —and select the option that provides the most utility:

$$d_{ijt} = \begin{cases} 1 & \text{if } U_{ijt} > U_{ikt} \text{ for all } k \neq j, \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$

Aggregate market shares are given by integrating over heterogeneous consumer choices:<sup>7</sup>

$$s_{jt} = \int d_{ijt}(\delta_t, \mu_{it}) d\mu_{it} d\epsilon_{it}.$$

When the  $\epsilon_{ijt}$  are i.i.d. with the type I extreme value distribution,<sup>8</sup>

$$s_{jt}(\delta_t, \tilde{\theta}_2) = \int \frac{\exp(\delta_{jt} + \mu_{ijt})}{\sum_{k \in J_t} \exp(\delta_{kt} + \mu_{ikt})} f(\mu_{it} | \tilde{\theta}_2) d\mu_{it}. \quad (3)$$

This is often referred to as a *mixed logit* or *random coefficients logit* because each individual  $i$ 's demands are given by a multinomial logit kernel where  $f(\mu_{it} | \tilde{\theta}_2)$  denotes the *mixing distribution* over the heterogeneous types  $i$  and  $\tilde{\theta}_2$  parameterizes this heterogeneity.<sup>9</sup> Indeed,  $\tilde{\theta}_2$  contains all of the parameters in the model related to the endogenous objects that are common to both supply and demand: the parameters  $\tilde{\theta}_2$  governing heterogeneous tastes  $\mu_{it}$  and the price coefficient  $\alpha$ .

The key insight of Berry (1994) or Berry, Levinsohn, and Pakes (1995) is that we can perform a nonlinear change of variables:  $\delta_t \equiv D_t^{-1}(\mathcal{S}_t, \tilde{\theta}_2)$  where  $\mathcal{S}_t$  denotes the  $J_t$  vector of observed market shares (see Berry and Haile, 2014). For each market  $t$ , (3) represents a system of  $J_t$  equations in  $J_t$  unknowns  $\delta_t$ . Given the  $\delta_t(\mathcal{S}_t, \tilde{\theta}_2)$  that solves that system of equations; along with some covariates  $x_{jt}$  and  $v_{jt}$ , prices  $p_{jt}$ , and a structural error  $\xi_{jt}$ ; under an additivity assumption, one can write the index:<sup>10</sup>

$$\delta_{jt}(\mathcal{S}_t, \tilde{\theta}_2) = [x_{jt}, v_{jt}] \beta - \alpha p_{jt} + \xi_{jt}. \quad (4)$$

With the addition of some instruments  $Z_{jt}^D$ , which include the exogenous regressors  $x_{jt}$  and  $v_{jt}$ , one can construct moment restriction conditions of the form  $E[\xi_{jt} Z_{jt}^D] = 0$ .

□ **Supply.** We can derive an additional set of supply moments from the first order conditions of firms. The conventional approach assumes that multi-product oligopoly firms simultaneously set prices independently for each market  $t$ .

Consider the profits of firm  $f$ , which for a single market  $t$  controls several products  $J_{ft}$  and sets prices  $p_{jt}$ . We take the first-order conditions of the profit function as follows:

$$\begin{aligned} \max_{p_{jt} : j \in J_{ft}} \sum_{j \in J_{ft}} s_{jt}(\mathbf{p}_t) \cdot (p_{jt} - c_{jt}), \\ s_{jt}(\mathbf{p}_t) + \sum_{k \in J_{ft}} \frac{\partial s_{kt}}{\partial p_{jt}}(\mathbf{p}_t) \cdot (p_{kt} - c_{kt}) = 0. \end{aligned}$$

It is helpful to write the first-order conditions in matrix form so that for a single market  $t$ ,

$$\begin{aligned} \mathbf{s}_t(\mathbf{p}_t) &= \Delta_t(\mathbf{p}_t) \cdot (\mathbf{p}_t - \mathbf{c}_t), \\ \underbrace{\Delta_t(\mathbf{p}_t)^{-1} \mathbf{s}_t(\mathbf{p}_t)}_{\eta_t(\mathbf{p}_t, \mathbf{s}_t, \theta_2)} &= \mathbf{p}_t - \mathbf{c}_t. \end{aligned} \quad (5)$$

<sup>7</sup> Here we define  $\delta_t$ ,  $\mu_{it}$ , and  $\epsilon_{it}$  as the  $J_t \times 1$  vectors with elements  $\delta_{jt}$ ,  $\mu_{ijt}$ , and  $\epsilon_{ijt}$ .

<sup>8</sup> Identification generally requires normalizing one of the options. As mentioned above, the typical choice is to normalize indirect utility from the outside option:  $U_{i0t} = \epsilon_{i0t}$ .

<sup>9</sup> McFadden and Train (2000) show that any random utility model (RUM) can be approximated with some mixture of multinomial logits with a sufficient basis of characteristics  $x_{jt}$ . The *mixed multinomial logit* itself dates back to at least Boyd and Mellman (1980) and Cardell and Dunbar (1980).

<sup>10</sup> Well-known special cases are logit without any parameters, in which  $D_t^{-1} = \log s_{jt} - \log s_{0t}$ , and nested logit with one parameter, in which  $D_t^{-1} = \log s_{jt} - \log s_{0t} - \rho \log s_{jht}$  where  $s_{jht}$  is the market share of  $j$  in its nest  $h$ .

Here the multi-product Bertrand markup  $\eta_t(\mathbf{p}_t, \mathbf{s}_t, \theta_2)$  depends on  $\Delta_t(\mathbf{p}_t, \mathbf{s}_t, \theta_2)$ , a  $J_t \times J_t$  matrix of intra-firm demand derivatives given by:

$$\Delta_t(\mathbf{p}_t) \equiv -\mathcal{H}_t \odot \frac{\partial \mathbf{s}_t}{\partial \mathbf{p}_t}(\mathbf{p}_t), \quad (6)$$

which is the element-wise Hadamard product of two  $J_t \times J_t$  matrices: the matrix of demand derivatives with each  $(j, k)$  entry given by  $\frac{\partial s_{jt}}{\partial p_{kt}}$  and the holdings or ownership matrix  $\mathcal{H}_t$  with each  $(j, k)$  entry indicating whether the same firm produces products  $j$  and  $k$ .<sup>11</sup>

This enables us to recover an estimate of marginal costs  $c_{jt} = p_{jt} - \eta_{jt}(\theta_2)$ , which in turn allows us to construct additional *supply side moments*. We can parametrize marginal cost as<sup>12</sup>

$$f_{MC}(p_{jt} - \eta_{jt}(\theta_2)) = f_{MC}(c_{jt}) = x_{jt}\gamma_1 + w_{jt}\gamma_2 + \omega_{jt} \quad (7)$$

and construct moment conditions of the form  $E[\omega_{jt}Z_{jt}^S] = 0$ . The idea is that we can use observed prices, along with information on demand derivatives and firm conduct, to recover markups  $\eta_{jt}$  and then marginal costs  $c_{jt}$ . This also imposes a functional form restriction on marginal costs, which depends on both product characteristics  $x_{jt}$  and the marginal cost shifters  $w_{jt}$  that are excluded from demand.

Some researchers may wish to allow marginal costs to depend on the quantity sold (either in market share or total units). This was true in the original Berry, Levinsohn, and Pakes (1995) article, which allowed  $\log c_{jt}$  to depend on  $\log q_{jt}$  to allow for returns to scale. PyBLP allows for this, but we do not explore quantity dependent marginal costs in our simulations. See Appendix A.3 for more information.

□ **The estimator.** We can construct a GMM estimator using our supply and demand moments. To do so, we stack their sample analogues to form

$$g(\theta) = \begin{bmatrix} g_D(\theta) \\ g_S(\theta) \end{bmatrix} = \begin{bmatrix} \frac{1}{N} \sum_{j,t} \xi_{jt} Z_{jt}^D \\ \frac{1}{N} \sum_{j,t} \omega_{jt} Z_{jt}^S \end{bmatrix}$$

and construct a nonlinear GMM estimator for  $\theta = [\beta, \alpha, \tilde{\theta}_2, \gamma]$  with some weighting matrix  $W$ :<sup>13</sup>

$$\min_{\theta} q(\theta) \equiv g(\theta)'Wg(\theta). \quad (8)$$

<sup>11</sup> Each  $(j, k)$  entry equals 1 if both  $j, k \in J_{ft}$  for some  $f \in F_t$ , and equals 0 otherwise. We can easily consider alternative forms of conduct such as Single- or Multi-Product Oligopoly, or Monopoly. Miller and Weinberg (2017) consider estimating a single parameter  $\mathcal{H}_t(\kappa)$  and Backus, Conlon, and Sinkinson (2020) use PyBLP test various forms of  $\mathcal{H}_t(\kappa)$ .

<sup>12</sup> The most common choice for  $f_{MC}(\cdot)$  is the identity function. Some authors also consider  $f_{MC}(\cdot) = \log(\cdot)$ . In practice this constrains marginal costs to be always positive.

<sup>13</sup> Some of the literature implicitly scales the objective by  $N^2$ . For example, Nevo (2000b) defines the objective for the demand-only problem as  $q(\theta) = \xi'Z_D W Z_D' \xi$ . Unless stated otherwise, we leave objective values unscaled throughout this article.

The  $\theta_2$  parameters are common to both supply and demand, govern the endogenous objects, and require at least one excluded instrument each.<sup>14</sup> To be explicit we write the entire program as follows:

$$\begin{aligned}
 \min_{\theta} q(\theta) &\equiv g(\theta)'Wg(\theta), \\
 g(\theta) &= \begin{bmatrix} \frac{1}{N} \sum_{j,t} \xi_{jt} Z_{jt}^D \\ \frac{1}{N} \sum_{j,t} \omega_{jt} Z_{jt}^S \end{bmatrix}, \\
 \xi_{jt} &= \delta_{jt} - [x_{jt}, v_{jt}] \beta + \alpha p_{jt}, \\
 \omega_{jt} &= f_{MC}(p_{jt} - \eta_{jt}) - [x_{jt}, w_{jt}] \gamma, \\
 \eta_t &= \Delta_t(\theta_2)^{-1} s_t, \\
 S_{jt} &= s_{jt}(\delta_t, \theta_2) \equiv \int \frac{\exp(\delta_{jt} + \mu_{ijt})}{\sum_{k \in J_t} \exp(\delta_{kt} + \mu_{ikt})} f(\mu_{it} | \tilde{\theta}_2) d\mu_{it}.
 \end{aligned} \tag{9}$$

This estimator and its econometric properties are discussed in Berry, Levinsohn, and Pakes (1995) and Berry, Linton, and Pakes (2004b). Our focus is not going to be on the econometric properties of  $\hat{\theta}$  but on rather various algorithms by which one might obtain  $\hat{\theta}$ . Technically, we need to solve this program twice. Once to obtain a consistent estimate for  $W$  and a second time to obtain the efficient GMM estimator.

Many applied articles omit the supply side from (9), and instead estimate  $\theta = [\theta_1, \theta_2]$  using demand moments alone, which is what PyBLP will do if the user does not provide a supply side. An important justification for not including the supply side is that it may be misspecified if the researcher does not know (or is not willing to assume) either the functional form of marginal costs  $f_{MC}(\cdot)$  or firm conduct  $\mathcal{H}_t$ .<sup>15</sup> The program without a supply side is as follows:

$$\begin{aligned}
 \min_{\theta} q_D(\theta) &\equiv g_D(\theta)'Wg_D(\theta), \\
 g_D(\theta) &= \frac{1}{N} \sum_{j,t} \xi_{j,t} Z_{jt}^D, \\
 \xi_{jt} &= \delta_{jt} - [x_{jt}, v_{jt}] \beta + \alpha p_{jt}, \\
 S_{jt} &= s_{jt}(\delta_t, \theta_2) \equiv \int \frac{\exp(\delta_{jt} + \mu_{ijt})}{\sum_{k \in J_t} \exp(\delta_{kt} + \mu_{ikt})} f(\mu_{it} | \tilde{\theta}_2) d\mu_{it}.
 \end{aligned} \tag{10}$$

□ **The nested fixed point algorithm.** In addition to providing an estimator, Berry, Levinsohn, and Pakes (1995) provide an algorithm for solving (9), which attempts to simplify the problem. Parameters on exogenous regressors enter the problem linearly; we concentrate out  $[\theta_1, \theta_3]$  and perform a nonlinear search over just  $\theta_2$  because  $[\hat{\theta}_1(\theta_2), \hat{\theta}_3(\theta_2)]$  are implicit functions of other parameters. Our modified algorithm is given below.

Our setup differs slightly from many BLP applications.<sup>16</sup> A key distinction occurs in step (d) where  $\alpha p_{jt}$  appears on the left-hand side of (11). The second distinction occurs in step (f), which requires that we stack the supply and demand equations appropriately. We provide a detailed

<sup>14</sup> Berry and Haile (2014) show that  $D_t^{-1}(\mathcal{S}_t, \tilde{\theta}_2)$  depends on the endogenous market shares of all products within the market, thus each parameter in  $\tilde{\theta}_2$  requires an additional instrument.

<sup>15</sup> There are other cases where the supply side may be misspecified. Many of these stem from misspecification around the functional form of  $\eta_t$  or  $\mathcal{H}_t$ . Important examples include: possible collusion, Cournot rather than Bertrand competition and double marginalization (Bonnet and Dubois, 2010; Villas-Boas, 2007).

<sup>16</sup> Arguably it is more in line with the original Berry, Levinsohn, and Pakes (1995) article.



**Algorithm 1 Nested Fixed Point**

For each guess of  $\theta_2$ :

- (a) For each market  $t$ , solve  $S_{jt} = s_{jt}(\delta_t, \theta_2)$  for  $\hat{\delta}_t(\mathbf{S}_t, \theta_2) \equiv \hat{\delta}_t(\theta_2)$ .
- (b) For each market  $t$ , use the  $J_t \times 1$  vector  $\hat{\delta}_t(\theta_2)$  to construct the  $J_t \times J_t$  matrix  $\Delta_t(\mathbf{p}_t, \hat{\delta}_t(\theta_2), \theta_2)$ .
- (c) For each market  $t$ , recover  $\hat{\eta}_t(\theta_2) = \Delta_t(\hat{\delta}_t(\theta_2), \theta_2)^{-1} \mathbf{S}_t$  by solving the  $J_t \times J_t$  linear system.
- (d) Stack up  $\hat{\delta}_t(\mathbf{S}_{jt}, \theta_2)$  and  $\hat{c}_{jt}(\hat{\delta}_t(\theta_2), \theta_2) = f_{MC}(p_{jt} - \hat{\eta}_{jt}(\hat{\delta}_t(\theta_2), \theta_2))$  and use linear IV-GMM to recover  $[\theta_1(\theta_2), \theta_3(\theta_2)]$  following the recipe in Appendix A.1. The following is our somewhat different formulation:

$$\begin{aligned}\hat{\delta}_{jt}(\mathbf{S}_t, \theta_2) + \alpha p_{jt} &= [x_{jt}, v_{jt}] \beta + \xi_{jt}, \\ f_{MC}(p_{jt} - \hat{\eta}_{jt}(\theta_2)) &= [x_{jt}, w_{jt}] \gamma + \omega_{jt}.\end{aligned}\quad (11)$$

- (e) Construct the residuals:

$$\begin{aligned}\hat{\xi}_{jt}(\theta_2) &= \hat{\delta}_{jt}(\theta_2) - [x_{jt}, v_{jt}] \hat{\beta}(\theta_2) + \alpha p_{jt}, \\ \hat{\omega}_{jt}(\theta_2) &= \hat{c}_{jt}(\theta_2) - [x_{jt}, w_{jt}] \hat{\gamma}(\theta_2).\end{aligned}\quad (12)$$

- (f) Stack the sample moments:

$$g(\theta_2) = \begin{bmatrix} \frac{1}{N} \sum_{jt} \hat{\xi}_{jt}(\theta_2) Z_{jt}^D \\ \frac{1}{N} \sum_{jt} \hat{\omega}_{jt}(\theta_2) Z_{jt}^S \end{bmatrix}.\quad (13)$$

- (g) Construct the GMM objective:  $q(\theta_2) = g(\theta_2)' W g(\theta_2)$ .

derivation in Appendix A.1. Later in Section 3, we show how this setup can be adapted to incorporate fixed effects when supply and demand are estimated simultaneously. This requires that the endogenous markup term  $\eta_{jt}(\mathbf{s}_t, \mathbf{p}_t, \delta_t(\theta_2), \theta_2)$  can be written as a function of only the  $\theta_2$  parameters and does not depend on  $[\theta_1, \theta_3]$ .<sup>17</sup>

We provide analytic gradients for the BLP problem with supply and demand in Appendix A.2. One advantage of the BLP algorithm is that it performs a nonlinear search over only  $K_2$  *nonlinear parameters*. Consequently, the Hessian matrix is only  $K_2 \times K_2$ . This implies relatively minimal memory requirements. Also, the IV-GMM regression in step (d) concentrates out the *linear* parameters  $[\theta_1, \theta_3]$ . This implies that large numbers of linear parameters can be estimated essentially for free, which is important if one includes a large number of fixed effects such as product- or market-level fixed effects.<sup>18</sup> In fact, other than (a) the remaining steps are computationally trivial. As is well known, (a)-(c) can be performed separately for each market across multiple processors.

The main disadvantage is that all parameters are implicit functions of other parameters, particularly of  $\theta_2$ . The objective is a complicated implicit function of  $\theta_2$ . Once we incorporate any heterogeneous tastes, the resulting optimization problem is nonconvex. In general the complexity of this problem grows rapidly with the number of nonlinear  $\theta_2$  parameters, whereas a high number of linear  $[\theta_1, \theta_3]$  parameters are more or less for free.

□ **Nested logit and RCNL variants.** The random coefficients nested logit (RCNL) model of Brenkers and Verboven (2006) instead assumes that  $\epsilon_{ijt}$  are not i.i.d. but rather follow the assumptions of a two-level nested logit. This model is popular in applications where the most important set of characteristics governing substitution is categorical. This has made it popular in studies of alcoholic beverages such as distilled spirits (Conlon and Rao, 2017; Miravete, Seim, and Thurk, 2018) and beer (Miller and Weinberg, 2017).

Much like the random coefficients model integrates over a heterogeneous distribution where each individual type follows a logit distribution, the RCNL model integrates over a heterogeneous distribution where each individual now follows a *nested logit*. Within group correlation is

<sup>17</sup> Why? We already know we can invert the shares to solve for  $\delta_t(\theta_2)$ , and the matrix of demand derivatives  $\frac{\partial s_{kt}}{\partial p_{jt}} = -\int \alpha_i s_{ikt}(\delta_t(\theta_2), \mu_{it}) [1(j=k) - s_{ijt}(\delta_t(\theta_2), \mu_{it})] f(\mu_{it}, \alpha_i | \theta_2) d\mu_{it}$  again depends only on  $\theta_2$ , which contains the price coefficient  $\alpha$ .

<sup>18</sup> For example, Nevo (2001, 2000b) includes product fixed effects.



governed by a new parameter  $\rho$ . We expand our definition of  $\theta_2$  to include the nesting parameter  $\rho$  so that  $\theta_2 \equiv [\alpha, \rho, \tilde{\theta}_2]$ .<sup>19</sup>

$$U_{ijt} = \delta_{jt} + \mu_{ijt}(\tilde{\theta}_2) + \epsilon_{ijt}(\rho).$$

The primary difference from the nested logit is that the inclusive value term for all products  $J_{ht}$  in nest  $h \in H = \{1, \dots, H\}$  now depends on the consumer's type  $i$ :

$$s_{jt}(\delta_t, \theta_2) = \int \frac{\exp[(\delta_{jt} + \mu_{ijt})/(1 - \rho)]}{\exp[IV_{iht}(\delta_t, \mu_{it})/(1 - \rho)]} \cdot \frac{\exp IV_{iht}(\delta_t, \mu_{it})}{1 + \sum_{h \in H} \exp IV_{iht}(\delta_t, \mu_{it})} f(\mu_{it} | \tilde{\theta}_2) d\mu_{it}, \quad (14)$$

in which the inclusive value term is

$$IV_{iht}(\delta_t, \mu_{it}) = (1 - \rho) \log \sum_{j \in J_{ht}} \exp \left( \frac{\delta_{jt} + \mu_{ijt}}{1 - \rho} \right).$$

A challenge for estimation is that  $\delta_t \leftarrow \delta_t + \log S_t - \log s_t(\delta_t, \theta_2)$  is no longer a contraction. Instead, the contraction must be dampened:<sup>20</sup>

$$\delta_t \leftarrow \delta_t + (1 - \rho)[\log S_t - \log s_t(\delta_t, \theta_2)]. \quad (15)$$

This creates an additional challenge because the rate of convergence for the contraction in (15) can become arbitrarily slow as  $\rho \rightarrow 1$ .<sup>21</sup> Thus as more consumers substitute within the nest, this model becomes much harder to estimate. Our simulations will demonstrate that this can be problematic.

As is well known, the relation in (15) has an analytic solution in the absence of random coefficients when  $\mu_{ijt} = 0$ . This model reduces to the regular nested logit for which the following expression was derived in Berry (1994):

$$\delta_{jt} = \log S_{jt} - \log S_{0t} - \rho \log S_{j|ht}$$

where  $S_{j|ht}$  is the market share of  $j$  in its nest  $h$ .

### 3. Algorithmic improvements

■ In each subsection, we review several methods from the literature, including some of our own design. Although we make many of these methods available as options in PyBLP, our focus is on finding best practices that are fastest and most reliable for most users. Later, we provide support for these decisions with our Monte Carlo studies.

□ **Incorporating many fixed effects.** There is a long tradition of extending the demand side utility to incorporate product or market fixed effects. For example, Nevo (2001, 2000a) allows for product fixed effects  $\xi_j$  so that

$$\delta_{jt} = [x_{jt}, v_{jt}] \beta - \alpha p_{jt} + \xi_j + \Delta \xi_{jt}.$$

These can manageably be incorporated as dummy variables in the linear part of the regression as there are only  $J = 24$  products in these articles.

With weekly UPC-store level Nielsen scanner data it is not uncommon for there to be  $J_t > 3500$  products (true in both distilled spirits and ready-to-eat cereal). There are approximately  $T = 500$  weeks  $t$  of Nielsen scanner data from 2006 to 2016. Incorporating store-week

<sup>19</sup> The nesting parameter can also be indexed as  $\rho_h$  so as to vary by group  $h$ . We support both types of nesting parameters in PyBLP.

<sup>20</sup> See Grigolon and Verboven (2014) for a derivation. The expression in (15) does not precisely match Grigolon and Verboven (2014) because of a minor typesetting error.

<sup>21</sup> This can be formalized in terms of the modulus of the contraction mapping or the Lipschitz constant. See Dubé, Fox, and Su (2012) for more details.

fixed effects  $\xi_{st}$  for only 100 stores could reach the order of 50,000 such fixed effects. Allowing for UPC-store fixed effects  $\xi_{js}$  can imply 100,000 or more fixed effects.<sup>22</sup> Clearly, the *least squares dummy variable* (LSDV) approach will not scale with tens or hundreds of thousands of fixed effects. We might consider the *within transformation* to remove the fixed effects, though we cannot directly incorporate both a within transformation and a supply side without re-writing the problem because of endogenous prices  $p_{jt}$ . We show how to re-write the problem in Appendix A.1. Define  $Y_{jt}^D$ ,  $Y_{jt}^S$ ,  $X_{jt}^D$ , and  $X_{jt}^S$  as follows:

$$\begin{aligned} Y_{jt}^D &\equiv \hat{\delta}_{jt}(\theta_2) + \alpha p_{jt} = [x_{jt}, v_{jt}]\beta + \xi_{jt} \equiv X_{jt}^D \beta + \xi_{jt}, \\ Y_{jt}^S &\equiv p_{jt} - \hat{\eta}_{jt}(\theta_2) = [x_{jt}, w_{jt}]\gamma + \omega_{jt} \equiv X_{jt}^S \gamma + \omega_{jt}. \end{aligned}$$

Stacking the system yields

$$\underbrace{\begin{bmatrix} Y_D \\ Y_S \end{bmatrix}}_Y = \underbrace{\begin{bmatrix} X_D & 0 \\ 0 & X_S \end{bmatrix}}_X \underbrace{\begin{bmatrix} \beta \\ \gamma \end{bmatrix}}_X + \underbrace{\begin{bmatrix} \xi \\ \omega \end{bmatrix}}_X. \quad (16)$$

After re-arranging terms and re-stacking, this is just a conventional linear IV problem in terms of  $(Y, X)$  where the endogenous parameters have been incorporated into  $Y$ . This means that the *within transform* can be used to absorb a single-dimensional fixed effect. Consider two dimensions of fixed effects  $j = 1, \dots, J$  and  $t = 1, \dots, T$ :

$$\begin{aligned} \tilde{Y}_{jt} &= Y_{jt} - \bar{Y}_{j\cdot} - \bar{Y}_{\cdot t}, \\ \tilde{X}_{jt} &= X_{jt} - \bar{X}_{j\cdot} - \bar{X}_{\cdot t}. \end{aligned}$$

The simplest approach might be to *iteratively demean*: remove the product mean  $\bar{X}_{j\cdot}$ , update  $X_{jt}$ , remove the market mean  $\bar{X}_{\cdot t}$ , and repeat this process until  $\bar{X}_{j\cdot} = \bar{X}_{\cdot t} = 0$ . This *method of alternating projections* (MAP) can be done in a single iteration if  $\text{Cov}(\bar{X}_{\cdot t}, \bar{X}_{j\cdot}) = 0$ . However, if the two dimensions of fixed effects are correlated this can require many iterations and can be quite burdensome.

The LSDV approach handles the burden of correlation but requires constructing the annihilator matrix to remove all the fixed effects. This approach requires inverting a  $(J + T) \times (J + T)$  matrix. Constructing and inverting such a large matrix is often infeasible because of memory requirements. Several algorithms have been proposed to deal with this problem. The most popular algorithm is perhaps that of Correia (2016), based on the accelerated MAP approach of Guimarães and Portugal (2010).<sup>23</sup>

The BLP application is unusual in that we re-run regressions using the same  $X$  variables many times. However, the left hand side  $Y$  variables  $\hat{\delta}_{jt}(\theta_2) + \alpha p_{jt}$  and  $\hat{c}(\theta_2)$  change with  $\theta_2$ , which means the entire procedure needs to be repeated for each guess of  $\theta_2$ . For more than a single dimension of fixed effects, PyBLP supports a number of different algorithms.<sup>24</sup> These algorithms have the advantage that the linear explanatory variables in  $X$  can be residualized only once, whereas the left hand side variables  $Y$  still need to be residualized for each guess of  $\theta_2$ . The savings are particularly large if the dimensions of  $X_D$  or  $X_S$  are large.

□ **Solving for the shares.** The main challenge of the Nested Fixed Point (NFXP) algorithm is solving the system of market shares:  $S_{jt} = s_{jt}(\delta_t, \theta_2)$ . Because the NFXP approach holds  $\theta_2$

<sup>22</sup> See Backus, Conlon, and Sinkinson (2020) and Conlon and Rao (2017) for examples of product-chain or store-week fixed effects.

<sup>23</sup> The Correia (2016) algorithm is implemented for the linear IV model in the Stata command `ivreghdfe`.

<sup>24</sup> PyBLP supports a number of different MAP acceleration schemes, the LSMR solver of Fong and Saunders (2011), and for two-dimensional fixed effects, the algorithm of Somaini and Wolak (2016). Comparison of these different approaches is beyond the scope of this article—for a more in-depth discussion, refer to Correia (2016).

fixed, rather than solve a system of  $N$  nonlinear equations and  $N$  unknowns, we solve  $T$  systems of  $J_t$  equations and  $J_t$  unknowns in  $\delta_t$  in parallel.<sup>25</sup>

Consider a single market  $t$  where we search for the  $J_t$  vector  $\delta_t$  which satisfies:

$$\mathcal{S}_{jt} = s_{jt}(\delta_t | \theta_2) = \int \frac{\exp(\delta_{jt} + \mu_{ijt})}{\sum_{k \in J_t} \exp(\delta_{kt} + \mu_{ikt})} f(\mu_{it} | \tilde{\theta}_2) d\mu_{it}. \quad (17)$$

Although mathematically there is a unique solution, it is impossible, numerically speaking, to choose a vector  $\delta_t$  that solves (17) exactly. Instead, we must solve the system of equations to some tolerance. We express the tolerance in terms of the log difference in shares:

$$\|\log \mathcal{S}_{jt} - \log s_{jt}(\delta_t, \theta_2)\|_{\infty} \leq \epsilon^{tol}. \quad (18)$$

There is a tradeoff with regard to the tolerance of the *inner loop* in (18). If the tolerance is too loose, the numerical error propagates to the rest of the estimation routine.<sup>26</sup> It is also possible to set a tolerance which is too tight and thus can never be satisfied. This is particularly problematic when summing over a large number of elements. We prefer to set  $\epsilon^{tol}$  between 1E-14 and 1E-12 as the *machine epsilon* or detectable difference between two double precision floating point numbers is around 1E-16 (on 64-bit architectures).

*Jacobian based approaches.* A direct approach would be to solve the system of  $J_t$  equations and  $J_t$  unknowns using Newton-type methods. Consider the following Newton-Raphson iteration:<sup>27</sup>

$$\delta_t^{h+1} \leftarrow \delta_t^h - \lambda \cdot \Psi_t^{-1}(\delta_t^h, \tilde{\theta}_2) \cdot s_t(\delta_t^h, \tilde{\theta}_2). \quad (19)$$

Each Newton-Raphson iteration would require computation of the  $J_t$  vector of market shares  $s_t(\delta_t^h, \tilde{\theta}_2)$ , the  $J_t \times J_t$  Jacobian matrix  $\Psi_t(\delta_t^h, \tilde{\theta}_2) = \frac{\partial s_t}{\partial \delta_t}(\delta_t^h, \tilde{\theta}_2)$ , and its inverse  $\Psi_t^{-1}(\delta_t^h, \tilde{\theta}_2)$ .

There are alternative *quasi-Newton* methods which solve variants of (19). These variants generally involve modifying the step-size  $\lambda$  or approximating  $\Psi_t^{-1}(\delta_t^h, \tilde{\theta}_2)$  in ways that avoid calculating the inverse Jacobian at each step. Quasi-Newton methods are often relatively fast when they work, though they need not converge (e.g., they may oscillate or not reach a resting point) and may be sensitive to starting values.<sup>28</sup>

Our experience indicates that the Levenberg–Marquardt (LM) algorithm is the fastest and most reliable Jacobian-based solution method.<sup>29</sup> LM minimizes the following least-squares problem in order to solve the following  $J_t \times J_t$  system of nonlinear equations:

$$\min_{\delta_t} \sum_{j \in J_t} [\mathcal{S}_{jt} - s_{jt}(\delta_t, \tilde{\theta}_2)]^2.$$

The idea is to update our guess of  $\delta_t^h$  to  $\delta_t^h + \mathbf{x}_t$  where  $\mathbf{x}_t$  is a  $J_t \times 1$  vector. The LM update is given by the solution to the following linear system of equations:

$$[\Psi_t' \Psi_t + \lambda \text{diag}(\Psi_t' \Psi_t)] \cdot \mathbf{x}_t = \Psi_t' [\mathcal{S}_t - s_t(\delta_t, \theta_2)] \quad (20)$$

<sup>25</sup> This same idea provides the sparsity of share constraints in the MPEC approach of Dubé, Fox, and Su (2012).

<sup>26</sup> Dubé, Fox, and Su (2012) show how error propagates from (18) to the estimates of  $\hat{\theta}$ . Lee and Seo (2016) provide a more precise characterization of this error when using Newton's method.

<sup>27</sup> In practice it is generally faster to solve the linear system:  $\Psi_t(\delta_t^h, \tilde{\theta}_2)(\delta_t^{h+1} - \delta_t^h) = -s_t(\delta_t^h, \tilde{\theta}_2)$ .

<sup>28</sup> Though the BLP problem is a nonconvex system of nonlinear equations, there are some properties which make it amenable to quasi-Newton methods. The market share function  $s_{jt}(\delta_t | \tilde{\theta}_2)$  is  $\mathbb{C}^\infty$  with respect to  $\delta_t$ , is bounded between (0,1), and agrees with its Taylor approximation at any  $\delta_t$ . Fox et al. (2012) and Iaria and Wang (2019) establish the real analyticity of the mixtures of logits under different conditions on the mixing distribution. Within some *basin of attraction*, quasi-Newton methods will be quadratically convergent. For an example of a quasi-Newton solution to (19), see Houde (2012). Another useful property is that with a strictly positive outside good share,  $\frac{\partial s_{jt}}{\partial \mu_{kt}} > 0$  for all  $k$ , which guarantees strict diagonal dominance of  $\Psi_t$  and hence that it is always nonsingular.

<sup>29</sup> Specifically we use the `lm` option of `scipy.optimize.root`, which calls the LM routine from MINPACK (More, Garbow, and Hillstom, 1980).

where  $\mathbf{x}_t$  is multiplied by an approximation to the Hessian. This has the advantage that for  $\lambda = 0$  the algorithm takes a full Gauss-Newton step and for  $\lambda$  large it takes a step in the direction of the gradient. The additional diagonal term also guarantees the invertibility of the approximate Hessian, even as it becomes nearly singular.

As in all iterative solution methods there are two primary costs: the cost per iteration and the number of iterations until convergence. The cost per iteration is driven primarily by the cost of computing (rather than inverting) the Jacobian matrix which involves  $J_t \times J_t$  numerical integrals.<sup>30</sup>

*Fixed point approaches.* Berry, Levinsohn, and Pakes (1995) also propose a fixed point approach to solve the  $J_t \times J_t$  system of equations in (17). They show that the following fixed point relation  $f(\delta_t) = \delta_t$  is a contraction mapping:<sup>31</sup>

$$f : \delta_t^{h+1} \leftarrow \delta_t^h + \log \mathcal{S}_t - \log s_t(\delta_t^h, \tilde{\theta}_2). \quad (21)$$

This kind of contraction mapping is linearly convergent with a rate of convergence that is proportional to  $L(\tilde{\theta}_2)/[1 - L(\tilde{\theta}_2)]$  where  $L(\tilde{\theta}_2)$  is the Lipschitz constant. Because (21) is a contraction, we know that  $L(\tilde{\theta}_2) < 1$ . Dubé, Fox, and Su (2012) show that for the BLP contraction the Lipschitz constant is  $L(\tilde{\theta}_2) = \max_{\delta_t} \|I_{J_t} - \frac{\partial \log s_t}{\partial \delta_t}(\delta_t, \tilde{\theta}_2)\|_{\infty}$ .

A smaller Lipschitz constant implies that (21) converges more rapidly. Dubé, Fox, and Su (2012) show in simulation that all else being equal, a larger outside good share generally implies a smaller Lipschitz constant.<sup>32</sup> Conversely, as the outside good share becomes smaller, the convergence of the fixed point relationship takes increasingly many steps.

*Accelerated fixed points.* Given a fixed point relationship there may be faster ways to obtain a solution to  $f(\delta_t) = \delta_t$  than direct iteration on the fixed point relation as in (21). There is a large literature on acceleration methods for fixed points. Most of these methods use information from multiple iterations  $(\delta_t^h, \delta_t^{h+1}, \delta_t^{h+2}, f(\delta_t^h), f(f(\delta_t^h)))$  to approximate  $\Psi_t$  or its inverse.<sup>33</sup>

Reynaerts, Varadhan, and Nash (2012) conduct extensive testing of various fixed point acceleration methods and find that the SQUAREM algorithm of Varadhan and Roland (2008) works well on the BLP contraction in (21). The intuition is to form a residual  $\mathbf{r}^h$  which is determined by the change between the current iteration  $\delta_t^h$  and the next iteration  $f(\delta_t^h)$ , as well as the change in the residual from this iteration to the next  $\mathbf{v}^h$  to form an estimate of the Jacobian. The residual and the curvature can also be used to construct a step-size  $\alpha^h$ . The exact algorithm is described below:<sup>34</sup>

$$\begin{aligned} \delta_t^{h+1} &\leftarrow \delta_t^h - 2\alpha^h \mathbf{r}^h + (\alpha^h)^2 \mathbf{v}^h, \\ \alpha^h &= \frac{(\mathbf{v}^h)' \mathbf{r}^h}{(\mathbf{v}^h)' \mathbf{v}^h}, \quad \mathbf{r}^h = f(\delta_t^h) - \delta_t^h, \quad \mathbf{v}^h = f(f(\delta_t^h)) - 2f(\delta_t^h) + \delta_t^h. \end{aligned} \quad (22)$$

In general the SQUAREM method is 3 to 6 times faster than direct iteration on the BLP contraction in (21). The idea is to take roughly the same number of steps as Newton-Raphson iteration, but to reduce the cost of steps by avoiding calculating the Jacobian directly. In fact, all of the terms in (22) are computed as a matter of course, because these are just iterations of  $\delta_t^h$  and  $f(\delta_t^h)$ . Unlike direct iteration on (21), there is technically no convergence guarantee as the iteration on (22) is no longer a contraction.

<sup>30</sup> A typical entry is  $\frac{\partial s_{jt}}{\partial \theta_t} = \int [1(j=k)s_{ijt}(\boldsymbol{\mu}_{it}) - s_{ijt}(\boldsymbol{\mu}_{it})s_{ikt}(\boldsymbol{\mu}_{it})]f(\boldsymbol{\mu}_{it}|\tilde{\theta}_2)d\boldsymbol{\mu}_{it}$ . The primary cost arises from numerical integration over heterogeneous types. Even for a large market with  $J_t = 1,000$  products, inverting a  $1,000 \times 1,000$  matrix is easy relative to numerically computing  $J_t^2$  integrals.

<sup>31</sup> Here  $f(\cdot)$  defines a contraction iteration and is not to be confused with the functional form for marginal costs  $f_{MC}(\cdot)$  or the distribution of heterogeneity  $f(\boldsymbol{\mu}_{it}|\tilde{\theta}_2)$ .

<sup>32</sup> A simple but novel derivation. Consider the matrix  $\frac{\partial \log s_t}{\partial \delta_t} = I_{J_t} - \text{diag}^{-1}(s_t)\Gamma_t(\tilde{\theta}_2)$  in which element  $(j, k)$  is  $1(j=k) - s_{jt}^{-1} \int s_{ijt}(\boldsymbol{\mu}_{it})s_{ikt}(\boldsymbol{\mu}_{it})f(\boldsymbol{\mu}_{it}|\tilde{\theta}_2)d\boldsymbol{\mu}_{it}$ . This implies that  $L(\tilde{\theta}_2) = \max_{\delta_t} [\max_j s_{jt}^{-1} \sum_k |\Gamma_{jkt}(\delta_t, \tilde{\theta}_2)|]$ . A rough approximation to the term inside the square braces is  $\max_j \sum_k |s_{jt}| \cdot |\text{Corr}(s_{ijt}, s_{ikt})| < 1 - s_{0t}$ .

<sup>33</sup> Many of these algorithms are special cases of *Anderson Mixing*.

<sup>34</sup> PyBLP includes a Python port of the SQUAREM package from R.

There are alternative acceleration methods in addition to SQUAREM. Reynaerts, Varadhan, and Nash (2012) also consider DF-SANE which takes the form  $\delta_t^{h+1} \leftarrow \delta_t^h - \alpha^h f(\delta_t^h)$  with a different choice of the step-size  $\alpha^h$ . They find performance is similar to SQUAREM though it can be slightly slower and less robust. Consistent with Reynaerts, Varadhan, and Nash (2012), we find the SQUAREM to be the fastest and most reliable accelerated fixed point approach. We find that the Jacobian-based Levenberg–Marquardt approach gives similar reliability and slightly better performance, though comparisons can be problem dependent.

□ **Optimization.** Optimization of the GMM objective function for the BLP problem can be challenging. The BLP problem in (9) represents a nonconvex optimization problem.<sup>35</sup> This has some important implications. First, the Hessian need not be positive semidefinite at all values of  $\theta_2$ . This also means that no optimization routine is guaranteed to find a global minimum in a fixed amount of time. This critique applies both to derivative-based quasi-Newton approaches and to simplex-based Nelder-Mead type approaches. Well-known recommendations involve considering a number of different starting values and optimization routines, verifying that  $\hat{\theta}$  satisfies both the first order conditions (the gradient is within a tolerance of zero) and second order conditions (the Hessian matrix has all positive eigenvalues).<sup>36</sup> Both of these are reported by default in PyBLP.

The PyBLP package has built-in support for the optimization routines implemented in the open-source SciPy library and can also interface with commercial routines such as Artleys Knitro. In Section 5, we compare different routines implemented in SciPy and Knitro. The optimization interface to PyBLP is generic in the sense that any optimizer can be used if it is implemented as a Python function, or, with the help of packages that allow for between-language interoperability, a function in most other languages.<sup>37</sup> Though nonderivative based routines such as Nelder-Mead have been frequently used in previous articles, they are not recommended.<sup>38</sup> Indeed, the most important aspect of PyBLP optimization is that it calculates the analytic gradient for any user-provided model, including models that incorporate both supply and demand moments or fixed effects.<sup>39</sup> Analytic gradients provide a major speedup in computational time and also generally improve the chances that the algorithm converges to a valid minimum.<sup>40</sup>

Some optimization routines allow the user to input constraints on parameters, which can speed up estimation and prevent the optimization routine from considering “unreasonable” values. An important example are *box constraints* on the parameter space, which restrict components of  $\theta_2$  to  $\theta_2^{(\ell)} \in [\underline{\theta}_2^{(\ell)}, \bar{\theta}_2^{(\ell)}]$ . Some typical constraints are that demand slopes down and that random coefficients have nonnegative but bounded variances.<sup>41</sup> This is particularly helpful because large values for random coefficients can generate some of the numerical issues that we discuss below.

<sup>35</sup> Absent unobserved heterogeneity or random coefficients the problem is globally convex.

<sup>36</sup> When there are parameter bounds, these conditions are based on the projected gradient and the reduced Hessian instead.

<sup>37</sup> In PyBLP’s documentation we give an example of such a “custom” routine by constructing a brute force solver that searches over a grid of parameter values. This flexibility should allow users to experiment with routines without much difficulty or “upgrade” if better routines are developed. For more on interoperability, see Footnote 3.

<sup>38</sup> Both Dubé, Fox, and Su (2012) and Knittel and Metaxoglou (2014) find that derivative-based routines outperform the simplex-based Nelder-Mead routine both in terms of speed and reliability. Our own tests concur with this prior opinion.

<sup>39</sup> We could not find any examples of simultaneous estimation of supply and demand with analytic gradients in the literature. A likely explanation is that the derivatives of the markup term  $\frac{\partial \mu}{\partial \theta_2}$  are quite complicated. See Appendix A.2 for a derivation.

<sup>40</sup> A promising alternative is automatic differentiation (AD). We chose to implement analytic gradients because ceding control to an AD library can limit one’s ability to handle numerical errors, which are pervasive in the BLP problem. However, AD is a promising technique for structural estimation and we are optimistic going forward.

<sup>41</sup> Requiring nonnegative variances on random coefficients is not absolutely necessary because the objective function should be symmetric about zero so that  $g(\sigma) = g(-\sigma)$ . This is because we optimize over the Cholesky root of the covariance for the random coefficients  $LL' = \Sigma$  rather than the covariance matrix itself.

In practice, a common issue is that the default termination tolerances in optimization software can be relatively loose. We provide some examples in Section 6. For termination conditions that are sensitive to the scale of the GMM objective value, this problem is made worse when  $N$  is large.<sup>42</sup> This highlights the importance of trying different optimizer configurations, particularly if the algorithm terminates surprisingly early.

Consistent with the prior literature, we recommend trying multiple optimizers and starting values to check for agreement. Our recommendation is that researchers try Knitro's Interior/Direct algorithm first if it is available,<sup>43</sup> and then try a BFGS-based routine, ideally with constraints on the parameter space such as SciPy's L-BFGS-B solver. An advantage of commercial solvers is that they work well out of the box, and do not require much in the way of configuration. However, our simulations indicate that when properly configured, most optimizers arrive at the same parameter estimates, satisfying both first and second order conditions. For the NFXP algorithm, the choice of solver seems much less important than proper configuration of parameter bounds, analytic gradients, and tight termination tolerances, which are all implemented by default in PyBLP.<sup>44</sup>

□ **Numerical issues and tricks.** There are several numerical challenges posed by the BLP problem, most of which are related to the exponentiated indirect utilities and the logit denominator:  $\sum_j \exp(\delta_{jt} + \mu_{ijt})$ . If some values in this summation are on the order of  $\exp(-5) \approx 0.0067$  and others are  $\exp(30) > 10^{13}$ , their sum may be rounded. This rounding occurs because most computers follow the IEEE-754 standard for floating point arithmetic and on most 64-bit architectures this means that floating point operations have around 15 significant digits. A *loss of precision* arises when taking summations of many numbers with different scales, and depending on the situation, may or may not be problematic. A related problem is *overflow*, which is likely to arise when attempting to compute large values such as  $\exp(800)$ . This can mean that  $s_{jt}(\delta_t, \theta_2) \rightarrow 1$  whereas  $s_{kt}(\delta_t, \theta_2) \rightarrow 0$  for  $k \neq j$ , leading optimization routines to fail.<sup>45</sup>

There are a number of solutions to these problems. One solution is to avoid large values in the argument of  $\exp(\cdot)$  by limiting the magnitudes of random coefficients through the *box constraints* described above. Another simple solution involves working market by market and avoiding very large summations. One additional method would be to use a summation algorithm that attempts to preserve precision such as *Kahan Summation*.<sup>46</sup> We found this to be substantially slower and thus do not implement it by default. As suggested in Skrainka (2012b), yet another approach is to use *extended precision arithmetic*. Again, we found this to be substantially slower without improving our statistical performance.<sup>47</sup>

Another way to guarantee overflow safety is to use a protected version of the log-sum-exp function:  $\text{LSE}(x) = \log \sum_k \exp x_k = a + \log \sum_k \exp(x_k - a)$ . By choosing  $a = \max\{0, \max_k x_k\}$ , use of this function helps ensure overflow safety when evaluating the multinomial logit function. This is a well-known trick from the applied mathematics and machine learning literatures, but there does not appear to be evidence of it in the literature on BLP models. This is implemented by default, as the additional computational cost appears to be trivial.

<sup>42</sup> In theory, the scale of the objective  $q(\theta) = g(\theta)'Wg(\theta)$  should be unimportant. By default, PyBLP scales by  $N$  so that after two-step GMM the objective equals the Hansen (1982)  $J$  statistic, the scale of which is invariant to problem size. We leave objective values unscaled throughout this article unless stated otherwise.

<sup>43</sup> The main disadvantage of Knitro is that it is not freely available—it must be purchased and installed by end-users.

<sup>44</sup> Commercial solvers likely have other advantages for other formulations of the problem, such as the MPEC approach of Dubé, Fox, and Su (2012), for which they recommend using Knitro.

<sup>45</sup> In this case  $\delta_t$  cannot be solved for in  $s_{jt}(\delta_t, \theta_2) = S_{jt}$ , so the inversion for the mean utilities fails.

<sup>46</sup> In Python this is implemented as `math.fsum`.

<sup>47</sup> PyBLP supports `numpy.longdouble` with the `pyblp.options.dtype` setting. On most Unix platforms, it is implemented with a 128-bit long double. For problems very with very small shares, or for problems with very large numbers of products, additional precision might still be valuable.



After implementing these suggestions, numerical issues are less common but can still occur. For example, the log-sum-exp function does not guard against occasional *underflow*, and in edge cases the weighting matrix  $W$  or the intra-firm matrix of demand derivatives  $\Delta_i$  may be nearly singular. If not dealt with, issues like these can cause the optimization routine to fail, which highlights the importance of robust and verbose error handling.<sup>48</sup>

Implementations of the BLP method have used a number of “tricks” to speed up the algorithm. One well-known example involves working with  $\exp(\delta_{jt})$  in place of  $\delta_{jt}$  in the contraction mapping to avoid transcendental functions like  $\exp(\cdot)$ . This has little benefit on modern architectures as transcendental functions are highly optimized and  $\exp(\cdot)$  of a billion numbers takes less than a second. Another trick is using  $\delta_i$  from a previous guess of  $\theta_2$  as a starting value. The SQUAREM and LM algorithms we suggest for solving the system of equations are relatively insensitive to starting values, so these sorts of speedups are not particularly useful. A “vectorization” speedup used in Nevo (2000b) and Knittel and Metaxoglou (2014) is to stack all markets together and construct cumulative sums of  $\exp(\delta_{jt} + \mu_{ijt})$ . We find that this approach is dominated by parallelization across markets  $T$ ; furthermore, this approach will result in a loss of precision as  $T \rightarrow \infty$ .

□ **Heterogeneity and integration.** An important aspect of the BLP model is that it incorporates heterogeneity via random coefficients. There are several ways to redefine the integral in (3):<sup>49</sup>

$$s_{jt}(\delta_t, \tilde{\theta}_2) = \int \frac{\exp(\delta_{jt} + \mu_{ijt})}{\sum_{k \in J_t} \exp(\delta_{kt} + \mu_{ikt})} f(\mu_{it} | \tilde{\theta}_2) d\mu_{it} \quad (23a)$$

$$= \int \frac{\exp[\delta_{jt} + \mu_{ijt}(\tilde{v}_{it}(v_{it}), \tilde{\theta}_2)]}{\sum_{k \in J_t} \exp[\delta_{kt} + \mu_{ikt}(\tilde{v}_{it}(v_{it}), \tilde{\theta}_2)]} \phi(v_{it}) dv_{it} \quad (23b)$$

$$= \int_{[0,1]^{K_2}} \frac{\exp[\delta_{jt} + \mu_{ijt}(\tilde{v}_{it}(u_{it}), \tilde{\theta}_2)]}{\sum_{k \in J_t} \exp[\delta_{kt} + \mu_{ikt}(\tilde{v}_{it}(u_{it}), \tilde{\theta}_2)]} du_{it}. \quad (23c)$$

The first approach in (23a) draws directly from the distribution of unobserved tastes  $f(\mu_{it} | \tilde{\theta}_2)$ . The second approach in (23b) integrates over the  $K_2$ -dimensional standard normal  $\phi(v_{it})$  and transforms the draws into a correlated normal using the Cholesky root  $L(\tilde{\theta}_2)$  of the covariance matrix so that  $\tilde{v}_{it} = L(\tilde{\theta}_2) \cdot v_{it}$  and  $\mu_{ijt} = \sum_{k \in K_2} x_{jt}^{(k)} \cdot \tilde{v}_{it}^{(k)}$ . The third method in (23c) integrates over the  $K_2$ -dimensional hypercube  $[0, 1]^{K_2}$  and transforms the draws using the inverse CDF of the mixing distribution:  $\tilde{v}_{it} = F^{-1}(u_{it} | \tilde{\theta}_2)$ . Under all three formulations the integral is “well behaved” in the sense that the integrand is bounded on  $[0, 1]$ , smooth, and infinitely continuously differentiable on  $\mathbb{C}^\infty$ .<sup>50</sup>

We can approximate the integral at a finite set of  $I_t$  nodes  $v_{it}$  and weights  $w_{it}$ . Together the nodes and weights  $(v_{it}, w_{it})$  define an *integration rule*. We use  $s_{jt}(\delta_t, \tilde{\theta}_2)$  to denote the true value of the integral and use  $s_{jt}(\delta_t, \tilde{\theta}_2; I_t)$  to denote the approximation with  $I_t$  nodes:

$$s_{jt}(\delta_t, \tilde{\theta}_2) \approx s_{jt}(\delta_t, \tilde{\theta}_2; I_t) \equiv \sum_{i \in I_t} w_{it} \cdot s_{ijt}(\delta_t, \mu_{it}(v_{it}, \tilde{\theta}_2)). \quad (24)$$

<sup>48</sup> When PyBLP encounters such a numerical error, it replaces the problematic object with a “reasonable” counterpart (e.g., corresponding values from the last optimization iteration or the Moore-Penrose pseudo-inverse of a near-singular matrix) and provides an informative warning or error message.

<sup>49</sup> We assume that heterogeneity is normally distributed to simplify some expressions. PyBLP currently accommodates normal and lognormal distributions of parameters. Others such as half-normal or exponential distributions can be implemented with a different change of variables.

<sup>50</sup> For a proof that all the derivatives are bounded see Iaria and Wang (2019).



*Pseudo-Monte Carlo.* The simplest integration rule is pseudo-Monte Carlo (pMC) integration. Here pseudo-random draws  $v_{it}$  are taken from one of the three candidate distributions in (23a)–(23c) and are used to calculate  $s_{ijt}(\delta_t, \mu_{it}(v_{it}, \tilde{\theta}_2))$ .<sup>51</sup> Each draw  $v_{it}$  is equally weighted:  $w_{it} = I_t^{-1}$ .

How accurate is the approximation in (24) under the pMC approach? Define the simulation error as  $\epsilon_t^{pMC} = s_{jt}(\delta_t, \theta_2) - s_{jt}(\delta_t, \tilde{\theta}_2; I_t)$ . A straightforward application of the Central Limit Theorem shows that:<sup>52</sup>

$$\epsilon_t^{pMC} \xrightarrow{d} N\left(0, \sqrt{V(s_{jt})/I_t}\right) \quad \text{with} \quad V(s_{jt}) = \int (s_{ijt}(\delta_t, \mu_{it}) - s_{jt}(\delta_t, \tilde{\theta}_2))^2 f(\mu_{it}|\tilde{\theta}_2) d\mu_{it}.$$

The main advantage of pMC integration is that it avoids the *curse of dimensionality* because the simulation error is not directly related to the dimension of the integral. Its main disadvantage is that the simulation error declines slowly in the number of draws at an  $O(I_t^{-1/2})$  rate. For a more extensive discussion of the pMC approach for BLP problems including bias correction and standard error adjustments, consult Freyberger (2015).

*Quasi-Monte Carlo.* Quasi-Monte Carlo (qMC) rules use deterministic sequences to more evenly cover the unit hypercube  $[0, 1]^{K_2}$  in (23c). Integration error can be reduced by choosing *low discrepancy sequences*.<sup>53</sup> A popular choice are Halton (1960) sequences. Train (2000), Bhat (2001), and Train (2009) report success with qMC sequences (Halton draws in particular) for maximum simulated likelihood (MSL) estimators of mixed logit models. Nevo (2001) reports success with Halton sequences for the BLP problem.

For fixed dimension of integration  $K_2$ , the  $O(I_t^{-1} \cdot (\log I_t)^{K_2})$  rate of qMC is expected to outperform the  $O(I_t^{-1/2})$  rate of pMC as the number of nodes  $I_t$  becomes large.<sup>54</sup> Owen (1997) shows that randomizing qMC sequences in a particular manner can give a  $O(I_t^{-3/2} \cdot (\log I_t)^{K_2})$  rate for smooth integrands as in (23c). For these reasons, it is common to *scramble* qMC sequences. Other common practices involve *discarding* the first few points and *skipping* others in-between selected points rather than taking points sequentially.<sup>55</sup>

*Variance reduction and importance sampling.* There are a number of additional techniques that have been used to reduce simulation error or variance. For example, Hess, Train, and Polak (2006) and Brunner et al. (2017) employ Modified Latin Hypercube Sampling (MLHS) for MSL and BLP problems, respectively. The MLHS approach takes a stratified sample by cutting the unit hypercube into smaller hypercubes and then sampling within each hypercube with pMC. Another common variance reduction technique is *antithetic sampling* which exploits the symmetry of  $\phi(v_{it}) = \phi(-v_{it})$  by taking each standard normal draw and reflecting it through the origin.

Importance sampling is also common. Rather than drawing from the standard multivariate normal distribution  $\phi(v_{it})$ , importance sampling instead proposes to draw from a biased distribution  $q(v_{it})$ . The idea is to draw more points where  $s_{ijt}(\delta_t, v_{it}, \tilde{\theta}_2) = s_{ijt}(\delta_t, \mu_{it}(v_{it}, \tilde{\theta}_2))$  is large and fewer where it is close to zero. To keep the value of the integral the same, the oversampled part of the distribution places less weight  $w(v_{it})$  on each point, and more on the undersampled

<sup>51</sup> By default PyBLP follows (23b), which is popular because it allows one to fix the draws  $v_{it}$  and simply rescale  $\tilde{v}_{it}$  as  $\tilde{\theta}_2$  changes without adding additional randomness to the estimation procedure. If the draws change with each estimation iteration this can create “chattering” because the simulated objective is not identical at every evaluation.

<sup>52</sup> See Chapter 7 of Judd (1998) for the derivation. It is easy to bound the variance term  $V(s_{jt}) < 1$ .

<sup>53</sup> These are *low discrepancy* in that  $D_h^*(u_{1,t}, \dots, u_{h,t}) \leq c \cdot I_t^{-1} \cdot (\log I_t)^{K_2}$  where the *star discrepancy*  $D_h^*$  bounds integration error in the Koksma-Hlawka inequality:  $|s_{jt}(\delta_t, \tilde{\theta}_2) - s_{jt}(\delta_t, \tilde{\theta}_2; I_t)| \leq V_{HK}(s_{jt}) \cdot D_h^*(u_{1,t}, \dots, u_{h,t})$ . Here  $V_{HK}$  denotes bounded variation in the sense of Hardy-Krause. See Judd (1998) and Owen (2005) for more details.

<sup>54</sup> Because the constant term can be large, this is by no means a guarantee. Also one requires that  $2^{K_2} < I_t$ .

<sup>55</sup> By default, PyBLP scrambles Halton sequences with the recipe in Owen (2017) and discards the first 1000 points in each dimension. Otherwise the first few points can exhibit a high degree of correlation between dimensions.

parts:

$$s_{jt}(\delta_t, \tilde{\theta}_2) = \int s_{ijt}(\delta_t, v_{it}, \tilde{\theta}_2) \cdot \frac{\phi(v_{it})}{q(v_{it})} \cdot q(v_{it}) dv_{it}.$$

One challenge is that the optimal  $q_{jt}(v_{it})$  is product specific.<sup>56</sup> Berry, Levinsohn, and Pakes (1995) oversample consumers with a small outside good by drawing from  $q_t(v_{it}) = \phi(v_{it}) \cdot \frac{1-s_{0t}(\delta_t, v_{it}, \tilde{\theta}_2)}{1-s_{0t}}$  where  $\tilde{\theta}_2$  and  $\delta_t$  are replaced by consistent estimates.<sup>57</sup>

*Quadrature.* The other approach is *Gaussian quadrature*. Theoretically, the integrand is approximated with a polynomial and then integrated exactly as polynomial integration is trivial. This amounts to a weighted sum in (24) over a particular choice of  $(v_{it}, w_{it})$ . The main choice to make is the *polynomial order* of the rule. As the order grows, more nodes are required but the accuracy of the approximation improves.

Gaussian quadrature works best when certain conditions are met: that the integrand is bounded and continuously differentiable. Thankfully, the logit kernel in (24) is always bounded between (0, 1) and is infinitely continuously differentiable. There are a number of different flavors of quadrature rules designed to best approximate integrals under different weighting schemes. The *Gauss-Hermite* family of rules work best when  $f(v_{it}) \propto \exp(-v_{it}^2)$ , which (with a change of variables) includes integrals over a normal density. *Nested* rules offer an alternative where for a given polynomial order  $p$ , they reuse the set of nodes from the rule with order  $p - 1$ .<sup>58</sup> Both the advantage and disadvantage of Gaussian quadrature rules is that they do a better job covering the “tail” of the probability distribution. Although this increases the accuracy of the approximation, it can also lead to very large values which create overflow issues.<sup>59</sup> We prefer to use quadrature rules, and to be careful of potential numerical issues when computing shares.

The Gaussian quadrature rules apply only to a single dimension. One way to estimate higher dimensional integrals is to construct the product of single dimensional integrals. The disadvantage of *product rules* is that if one needs  $I_t$  points to approximate the integral in dimension one, then one needs  $I_t^d$  points to approximate the integral in dimension  $d$ . This is the so-called *curse of dimensionality*.

The curse of dimensionality is a well-known problem in numerical analysis and several off-the-shelf solutions exist. There are several clever algorithms for improving upon the product rule for higher dimensional integration. Judd and Skrainka (2011) explore *monomial cubature rules* whereas Heiss and Winschel (2008) use *sparse grid* methods to selectively eliminate nodes from the product rules. One disadvantage of these methods is that they often involve negative weights  $w_{it} < 0$ , which can create problems during estimation or when trying to decompose the distribution of heterogeneity (particularly for counterfactuals).

*Integration in PyBLP.* Though PyBLP allows for flexible, user-supplied distributions of random coefficients, by far the most commonly employed choices in the literature are the independent normal, correlated normal, and lognormal distributions for  $f(\mu_{it}|\theta_2)$ . PyBLP supports all of these distributions and provides some specialized routines to handle these integrals with limited user

<sup>56</sup> Heiss (2010) proposes an “adaptive importance sampler” for probit-type models. Brunner (2017) adapts this to the BLP-type problem and reports success. This adaptive, optimal importance sampling requires updating the weights  $w_{it}$  as  $\theta_2$  varies. For this reason, we don’t currently implement this approach as part of PyBLP.

<sup>57</sup> At the consistent estimate for  $\tilde{\theta}_2$ , the  $\delta_t$  needs to be computed only once, so it is feasible to use a higher-precision integration rule than is used for estimation. Berry, Levinsohn, and Pakes (1995) draw from  $q_t(v_{it})$  with rejection sampling, and this same procedure is implemented in PyBLP.

<sup>58</sup> In theory, this allows for *adaptive* accuracy without wasting calculations. When the error from numerical integration is large, the polynomial degree can be expanded to reduce the error.

<sup>59</sup> Essentially for some simulated individual  $i$  we have that  $s_{ijt} \rightarrow 1$  and  $s_{ikt} \rightarrow 0$ . This problem has been previously documented by Judd and Skrainka (2011) and Skrainka (2012b).

intervention. There are a number of different methods one can use to generate  $(v_{it}, w_{it})$  for the normal case: pMC, Halton sequences, MLHS, product rules, and sparse grids.

In general, we find that the best practice in low dimensions is to use product rules to a relatively high degree of polynomial accuracy. In higher dimensions, Halton sequences and in particular sparse grids appear to scale the best, both in our own Monte Carlo studies and in those of Judd and Skrainka (2011) and Heiss and Winschel (2008).

□ **Solving for pricing equilibria.** Many counterfactuals of BLP-type problems involve perturbing either the market structure, marginal costs, or both, and solving for counterfactual equilibrium prices. Being able to solve for equilibrium prices quickly and accurately is also crucial to generating the optimal instruments in the next section. The Bertrand-Nash first order conditions are defined by (5) for each market  $t$ :<sup>60</sup>

$$\mathbf{p}_t = \mathbf{c}_t + \underbrace{\Delta_t(\mathbf{p}_t, \mathcal{H}_t)^{-1} \mathbf{s}(\mathbf{p}_t)}_{\boldsymbol{\eta}_t(\mathbf{p}_t, \mathcal{H}_t)}$$

In order to recover marginal costs during estimation, one need only invert the  $J_t \times J_t$  matrix  $\Delta_t(\mathbf{p}_t, \mathcal{H}_t)$ . Solving for counterfactual pricing equilibria is more difficult as it requires solving the  $J_t \times J_t$  nonlinear system of equations, often after replacing the ownership matrix  $\mathcal{H}_t$  with a post-merger counterpart  $\mathcal{H}_t^*$ :

$$\mathbf{p}_t = \mathbf{c}_t + \boldsymbol{\eta}_t(\mathbf{p}_t, \mathcal{H}_t^*). \quad (25)$$

In general, solving this problem is difficult because it represents a nonconvex, nonlinear system of equations where one must simulate in order to compute  $\boldsymbol{\eta}_t(\mathbf{p}_t, \mathcal{H}_t^*)$  and its derivatives. Once one incorporates both multi-product firms and arbitrary coefficients into the problem, both existence and uniqueness of an equilibrium become challenging to establish.<sup>61</sup>

One approach might be to solve the system using Newton's method, which requires calculating the  $J_t \times J_t$  Jacobian  $\frac{\partial \boldsymbol{\eta}_t}{\partial \mathbf{p}_t}$ . The expression for the Jacobian involves the Hessian matrix of demand  $\frac{\partial^2 s_{kt}}{\partial p_{jt}^2}$  as well as tensor products, and can be computationally challenging.<sup>62</sup>

The second, and perhaps most common approach in the literature is treating (25) as a fixed point and iterating on  $\mathbf{p}_t \leftarrow \mathbf{c}_t + \boldsymbol{\eta}_t(\mathbf{p}_t, \mathcal{H}_t^*)$ .<sup>63</sup> The problem is that although a fixed point of (25) may represent the Bertrand-Nash equilibrium of (7), it is not necessarily a contraction. In fact, as part of Monte Carlo experiments conducted in Armstrong (2016), the author finds that iterating on (25) does not always lead to a solution and at least some fraction of the time leads to cycles. We were able to replicate this finding for similarly constructed Monte Carlo experiments between 1-5% of the time. Our preferred approach borrows from the engineering literature and does not appear to be well known in Industrial Organization, but in our experiments appears to be highly effective. We follow Morrow and Skerlos (2011) who reformulate the solution to (5) by breaking up the matrix of demand derivatives into two parts: a  $J_t \times J_t$  diagonal matrix  $\Lambda_t$ , and a  $J_t \times J_t$

<sup>60</sup> We suppress the dependence on the parameters  $\theta_2$  as we hold everything fixed and solve for  $\mathbf{p}_t$ .

<sup>61</sup> Caplin and Nalebuff (1991) and Gallego et al. (2006) have results that apply to single product firms and linear in price utility under logit demands. Kononov and Sandor (2010) generalizes these results to logit demands with linear in price utility and multi-product firms. With the addition of random coefficients, it is possible that the resulting model will violate the quasi-concavity of the profit function that these results require. Morrow and Skerlos (2010) avoid some of these restrictions but place other restrictions on indirect utilities. Existence and uniqueness are beyond the scope of this article—we instead focus on calculating solutions to the system of first order conditions, assuming such solutions exist.

<sup>62</sup> For example, Knittel and Metaxoglou (2014) do not update  $s_t(\mathbf{p}_t)$  and thus avoid fully solving the system of equations. An Newton-type alternative with a finite-differenced Jacobian would be slow and ill-advised as there are  $J_t$  derivatives for  $J_t$  markups and every derivative involves integration in order to compute both  $\Delta_t(\mathbf{p}_t)$  and  $s_t(\mathbf{p}_t)$ .

<sup>63</sup> The “folklore” solution is to dampen this expression with some  $\rho$  and consider  $\mathbf{p}_t \leftarrow \rho \cdot \mathbf{p}_t + (1 - \rho) \cdot [\mathbf{c}_t + \boldsymbol{\eta}_t^*(\mathbf{p}_t, \mathcal{H}_t^*)]$ . This tends to be slower and more reliable, though we cannot find any theoretical justification for convergence.

dense matrix  $\Gamma_t$ :

$$\begin{aligned}\frac{\partial s_t}{\partial \mathbf{p}_t}(\mathbf{p}_t) &= \Lambda_t(\mathbf{p}_t) - \Gamma_t(\mathbf{p}_t), \\ \Lambda_{jj,t} &= \int \alpha_i s_{ijt}(\boldsymbol{\mu}_{it}) f(\boldsymbol{\mu}_{it} | \tilde{\theta}_2) d\boldsymbol{\mu}_{it}, \\ \Gamma_{jk,t} &= \int \alpha_i s_{ijt}(\boldsymbol{\mu}_{it}) s_{ikt}(\boldsymbol{\mu}_{it}) f(\boldsymbol{\mu}_{it} | \tilde{\theta}_2) d\boldsymbol{\mu}_{it},\end{aligned}\tag{26}$$

in which  $\alpha_i = \frac{\partial u_{ijt}}{\partial p_{jt}}$  is the marginal dis-utility of price. Morrow and Skerlos (2011) then reformulate the problem as a different fixed point that is specific to mixed logit demands:<sup>64</sup>

$$\mathbf{p}_t \leftarrow \mathbf{c}_t + \boldsymbol{\zeta}_t(\mathbf{p}_t) \quad \text{where} \quad \boldsymbol{\zeta}_t(\mathbf{p}_t) = \Lambda_t(\mathbf{p}_t)^{-1} [\mathcal{H}_t^* \odot \Gamma_t(\mathbf{p}_t)] (\mathbf{p}_t - \mathbf{c}_t) - \Lambda_t(\mathbf{p}_t)^{-1} \mathbf{s}_t(\mathbf{p}_t). \tag{27}$$

The fixed point in (27) is entirely different from that in (25) and coincides only at resting points. Consistent with results reported in Morrow and Skerlos (2011), we find that (27) is around 3-12 times faster than Newton-type approaches and reliably finds an equilibrium.<sup>65</sup>

Perhaps most consequentially, the ability to solve for a pricing equilibrium rapidly and reliably makes it possible to generate the Amemiya (1977) or Chamberlain (1987) *feasible approximation to the optimal instruments*.

#### 4. Supply and demand: optimal instruments and overidentifying restrictions

■ *In this section, we focus on joint estimation of supply and demand under optimal instruments in order to clarify the precise role of overidentifying restrictions in the parametric identification of parameters. Absent information on firm conduct or the precise functional form for marginal costs, many researchers estimate a demand side only.*

As a way to improve performance, we can construct approximations to the optimal instruments in the spirit of Amemiya (1977) or Chamberlain (1987). Approximations to the optimal instruments were featured in both Berry, Levinsohn, and Pakes (1995) and Berry, Levinsohn, and Pakes (1999) but are not commonly employed in many subsequent studies using the BLP approach in part because they are challenging to construct. Reynaert and Verboven (2014) show that approximations to the optimal instruments can improve the econometric performance of the estimator, particularly with respect to the  $\theta_2$  parameters. The form we derive is somewhat different from their expression, and arrives at different instruments for supply and demand. Although the procedure itself is quite involved, the good news is that it does not require much in the way of user input, and is fully implemented by the PyBLP software.

□ **Derivation.** Recall the GMM moment conditions are given by  $E[\xi_{jt}|Z_{jt}^D] = 0$  and  $E[\omega_{jt}|Z_{jt}^S] = 0$  and the asymptotic GMM variance depends on  $(D' \Omega^{-1} D)$  where the expressions are given below:

$$D = E \left[ \left( \frac{\partial \xi_{jt}}{\partial \theta}, \frac{\partial \omega_{jt}}{\partial \theta} \right) \middle| Z_t \right], \quad \Omega = E \left[ \begin{pmatrix} \xi_{jt} \\ \omega_{jt} \end{pmatrix} \begin{pmatrix} \xi_{jt} & \omega_{jt} \end{pmatrix} \middle| Z_t \right].$$

Chamberlain (1987) showed that the approximation to the optimal instruments are given by the expected Jacobian contribution for each observation  $(j, t)$ :  $E[D_{jt}(Z_t) \Omega_{jt}^{-1} | Z_t]$ . We use the word

<sup>64</sup> This resembles a well known “folklore” solution to the pricing problem, which is to rescale each equation by its own share  $s_{jt}$  (see Skrainka, 2012a). For the plain logit,  $\Lambda_{jj,t} = \alpha s_{jt}$ .

<sup>65</sup> In PyBLP, iteration is terminated when the firms’ first order condition  $\|\Lambda(\mathbf{p}_t)(\mathbf{p}_t - \mathbf{c}_t - \boldsymbol{\zeta}_t(\mathbf{p}_t))\|_\infty$  is less than a fixed tolerance. Morrow and Skerlos (2011) refer to this as the problem’s *numerical simultaneous stationarity condition*.

“approximation” because the aforementioned expectation over the unobserved  $(\xi_{jt}, \omega_{jt})$  lacks a closed form. We derive the components of the approximation below:<sup>66</sup>

$$D_{jt} \equiv \underbrace{\begin{bmatrix} \frac{\partial \xi_{jt}}{\partial \beta} & \frac{\partial \omega_{jt}}{\partial \beta} \\ \frac{\partial \xi_{jt}}{\partial \alpha} & \frac{\partial \omega_{jt}}{\partial \alpha} \\ \frac{\partial \xi_{jt}}{\partial \theta_2} & \frac{\partial \omega_{jt}}{\partial \theta_2} \\ \frac{\partial \xi_{jt}}{\partial \gamma} & \frac{\partial \omega_{jt}}{\partial \gamma} \end{bmatrix}}_{(K_1+K_2+K_3) \times 2} = \begin{bmatrix} -x_{jt} & 0 \\ -v_{jt} & 0 \\ \frac{\partial \xi_{jt}}{\partial \alpha} & \frac{\partial \omega_{jt}}{\partial \alpha} \\ \frac{\partial \xi_{jt}}{\partial \theta_2} & \frac{\partial \omega_{jt}}{\partial \theta_2} \\ 0 & -x_{jt} \\ 0 & -w_{jt} \end{bmatrix}, \quad \Omega_t \equiv \underbrace{\begin{bmatrix} \sigma_{\xi_t}^2 & \sigma_{\xi_t \omega_t} \\ \sigma_{\xi_t \omega_t} & \sigma_{\omega_t}^2 \end{bmatrix}}_{2 \times 2}. \quad (28)$$

A little calculation shows that for each market  $t$  and product  $j$ ,

$$D_{jt} \Omega_t^{-1} = \frac{1}{\sigma_{\xi}^2 \sigma_{\omega}^2 - \sigma_{\xi \omega}^2} \cdot \begin{bmatrix} -\sigma_{\omega}^2 x_{jt} & \sigma_{\xi \omega} x_{jt} \\ -\sigma_{\omega}^2 v_{jt} & \sigma_{\xi \omega} v_{jt} \\ \sigma_{\omega}^2 \frac{\partial \xi_{jt}}{\partial \alpha} - \sigma_{\xi \omega} \frac{\partial \omega_{jt}}{\partial \alpha} & \sigma_{\xi}^2 \frac{\partial \omega_{jt}}{\partial \alpha} - \sigma_{\xi \omega} \frac{\partial \xi_{jt}}{\partial \alpha} \\ \sigma_{\omega}^2 \frac{\partial \xi_{jt}}{\partial \theta_2} - \sigma_{\xi \omega} \frac{\partial \omega_{jt}}{\partial \theta_2} & \sigma_{\xi}^2 \frac{\partial \omega_{jt}}{\partial \theta_2} - \sigma_{\xi \omega} \frac{\partial \xi_{jt}}{\partial \theta_2} \\ \sigma_{\xi \omega} x_{jt} & -\sigma_{\xi}^2 x_{jt} \\ \sigma_{\xi \omega} w_{jt} & -\sigma_{\xi}^2 w_{jt} \end{bmatrix}. \quad (29)$$

For each column vector, the first and fifth entries are linear functions of  $x_{jt}$ . Let  $\Theta$  be a conformable matrix of zeros and ones such that

$$(D_{jt} \Omega_t^{-1}) \odot \Theta = \frac{1}{\sigma_{\xi}^2 \sigma_{\omega}^2 - \sigma_{\xi \omega}^2} \cdot \begin{bmatrix} -\sigma_{\omega}^2 x_{jt} & 0 \\ -\sigma_{\omega}^2 v_{jt} & \sigma_{\xi \omega} v_{jt} \\ \sigma_{\omega}^2 \frac{\partial \xi_{jt}}{\partial \alpha} - \sigma_{\xi \omega} \frac{\partial \omega_{jt}}{\partial \alpha} & \sigma_{\xi}^2 \frac{\partial \omega_{jt}}{\partial \alpha} - \sigma_{\xi \omega} \frac{\partial \xi_{jt}}{\partial \alpha} \\ \sigma_{\omega}^2 \frac{\partial \xi_{jt}}{\partial \theta_2} - \sigma_{\xi \omega} \frac{\partial \omega_{jt}}{\partial \theta_2} & \sigma_{\xi}^2 \frac{\partial \omega_{jt}}{\partial \theta_2} - \sigma_{\xi \omega} \frac{\partial \xi_{jt}}{\partial \theta_2} \\ 0 & -\sigma_{\xi}^2 x_{jt} \\ \sigma_{\xi \omega} w_{jt} & -\sigma_{\xi}^2 w_{jt} \end{bmatrix}. \quad (30)$$

We can partition our instrument set by column into “demand” and “supply” instruments:<sup>67</sup>

$$Z_{jt}^{Opt,D} \equiv \underbrace{E[(D_{jt}(Z_t) \Omega_t^{-1} \odot \Theta)_{\cdot 1} | Z_t]}_{K_1+K_2+(K_3-K_x)}, \quad Z_{jt}^{Opt,S} \equiv \underbrace{E[(D_{jt}(Z_t) \Omega_t^{-1} \odot \Theta)_{\cdot 2} | Z_t]}_{K_2+K_3+(K_1-K_x)}. \quad (31)$$

Here, we have  $K - K_x$  (where  $K_x$  denotes the dimension of common exogenous parameters  $x_{jt}$ ) instruments for both supply and demand, though it is evident from (30) that the instruments for the  $\theta_2$  parameters are not the same.<sup>68</sup> The optimal instruments from the *linear* portions of demand and supply are simply exogenous regressors re-scaled by covariances, whereas the optimal instruments from the  $\theta_2$  parameters are *nonlinear* functions of the data.

Two sets of overidentifying restrictions arise from *exclusion restrictions*, which are made explicit in (30) where  $w_{jt}$  and  $v_{jt}$  show up in one equation but not the other. There are  $K_3 - K_x$  overidentifying restrictions from cost shifters  $w_{jt}$  that are excluded from demand and  $K_1 - K_x$  demand shifters  $v_{jt}$  that are excluded from supply. When we include simultaneous supply and demand moments we also have *cross equation* restrictions. As shown in (30), we have

<sup>66</sup> In our Monte Carlo exercises we assume that  $(\xi_{jt}, \omega_{jt})$  are jointly i.i.d. across all  $j$  and  $t$  so that  $\Omega_{jt} = \Omega$ . This is merely a matter of convenient notation, as extensions to heteroskedastic or clustered covariances are straightforward.

<sup>67</sup> Here  $(\cdot)_1$  and  $(\cdot)_2$  denote the first and second column of the matrix in (30) and  $\odot$  is the elementwise Hadamard product.

<sup>68</sup> This is true except for the knife-edge case where  $\frac{\partial \xi_{jt}}{\partial \theta_2} / \frac{\partial \omega_{jt}}{\partial \theta_2} \propto \frac{\sigma_{\xi}^2 + \sigma_{\xi \omega}}{\sigma_{\omega}^2 + \sigma_{\xi \omega}}$ . In fact, (30) highlights that the set of instruments  $Z_D$  and  $Z_S$  should never be the same because of the different excluded variables.

$2 \times (K - K_x)$  restrictions and  $K$  parameters. This gives us  $K - 2K_x$  overidentifying restrictions. The additional  $K_2$  overidentifying restrictions in (30) come from the fact that we have two restrictions for each of the  $\theta_2$  parameters, including the price coefficient  $\alpha$ .

*Remark 1.* Our version of the optimal instruments makes explicit the exclusion restrictions in the BLP model. Perhaps most importantly, this tells us precisely where to find exclusion restrictions: *something that enters the other equation*. Both the supply and demand moments are informative for the  $\theta_2$  parameters. The role of the exogenous cost-shifters  $w_{jt}$  is now explicit: they provide overidentifying restrictions for demand that are informative about  $\theta_2$  (including the term on price  $\alpha$ ). Likewise, the role of the exogenous demand-shifters  $v_{jt}$  is also made explicit: they provide overidentifying restrictions for supply which are informative about  $\theta_2$  (and markups).<sup>69</sup> The link between the supply and demand side is the endogenous markup  $\eta_{jt}(\theta_2, \xi_t, \omega_t)$ , which depends on the common  $\theta_2$  parameters. This has led  $(\frac{\partial \xi_{jt}}{\partial \theta_2}, \frac{\partial \omega_{jt}}{\partial \theta_2})$ -type instruments to be described as *quantity shifters* or *markup shifters*.

*Remark 2.* It is worth pointing out that our derivation of the optimal instruments appears to vary from derivations in the prior literature. Some of the prior literature using optimal instruments for BLP-type problems suggests the resulting problem is *just identified* rather than *over identified* and relies on the same set of instruments for both supply and demand. Reynaert and Verboven (2014) appear to construct their version of optimal instruments by summing across the rows of (29) and excluding either the first or third row.<sup>70</sup> This gives  $K = K_1 + K_2 + K_3$  instruments and  $K$  unknowns so that the model is just identified. However, because they stack  $(\xi_t, \omega_t)$  they effectively have  $2N$  rather than  $N$  observations. Conceptually, one way to view their formulation is that it imposes  $E[\xi_{jt}Z_{jt}^D] + E[\omega_{jt}Z_{jt}^S] = 0$  rather than separately imposing  $E[\xi_{jt}Z_{jt}^D] = 0$  and  $E[\omega_{jt}Z_{jt}^S] = 0$ .

*Remark 3.* One alternative to the Chamberlain (1987) approximation to the optimal instruments is to instead construct a semiparametric basis that spans the same vector space as the optimal instruments. This approach was suggested by Newey (1990) and applied to conditional moment restriction models in Ai and Chen (2003) and Donald, Imbens, and Newey (2009). Given a conditional moment condition of the form  $E[\xi_{jt}|Z_t] = 0$ , one can instead write  $E[\xi_{jt}A(Z_t)] = 0$  for some choice of  $A(\cdot)$ . One could impose that the moments hold at several quantiles of  $Z_t$  or one can construct a polynomial sieve basis in  $Z_t$ . This approach also suffers from a curse of dimensionality as in high dimensions the number of needed interaction terms explodes. The bases must also be chosen carefully as it is easy to generate new instruments  $A(Z_t)$  that are highly correlated with one another, which leads to the “many moments” problem of Newey and Smith (2004).<sup>71</sup> To highlight this approach, we consider all quadratic interactions of  $(x_{jt}, w_{jt})$  when deriving instruments in “Own” characteristics. Gandhi and Houde (2019) derive a second-order polynomial basis in the *differences* of product characteristics  $d_{jkt} = x_{kt} - x_{jt}$ , and show that this basis has desirable properties and avoids several of the aforementioned problems.

□ **Constructing feasible instruments.** The main challenge with implementing the optimal instruments is that we must take the expectation of the Jacobian over the joint distribution of unobservables  $(\xi_t, \omega_t)$  for all products within a market. We need to compute the following

<sup>69</sup> The idea of using demand shifters as overidentifying restrictions to identify conduct has a long history in industrial organization going back at least as far as Bresnahan (1982), and was treated nonparametrically in Berry and Haile (2014). In related work, Backus, Conlon, and Sinkinson (2020) show how to use (30) to test for firm conduct.

<sup>70</sup> We should mention that Reynaert and Verboven (2014) do not consider joint supply and demand as their main specification.

<sup>71</sup> For example, the polynomial basis  $(1, x^2, x^3, x^4)$  exhibits a high degree of correlation.



expectation:<sup>72</sup>

$$E\left[\frac{\partial \xi_{jt}}{\partial \theta_2}, \frac{\partial \omega_{jt}}{\partial \theta_2} \middle| Z_t\right] = \int \left[\frac{\partial \xi_{jt}}{\partial \theta_2}, \frac{\partial \omega_{jt}}{\partial \theta_2}\right](\xi_t, \omega_t, Z_t; \theta_2) f(\xi_{1,t}, \dots, \xi_{J_t,t}, \omega_{1,t}, \dots, \omega_{J_t,t} \middle| Z_t, \theta_2) d\xi_t d\omega_t.$$

One challenge is that this is an integral over  $2J_t$  dimensions. A second challenge is that without additional assumptions,  $f(\xi_t, \omega_t | Z_t, \theta_2)$  is unknown. The third is that  $(s_t, p_t)$  are endogenous in that they depend on  $(\xi_t, \omega_t)$ . This means we must construct estimates of  $E[p_{jt} | Z_t]$  and  $E[s_{jt} | Z_t]$  in order to calculate  $\frac{\partial \xi_{jt}}{\partial \theta_2}$ . One approach would be to construct  $E[p_{jt} | Z_t] = \hat{p}_{jt}$  by regressing the endogenous prices  $p_{jt}$  on a series of exogenous regressors in a “first stage” regression (either linearly or nonlinearly) and evaluating  $\hat{s}_t(\hat{p}_t, \theta_2)$ .<sup>73</sup> The other approach is to solve the nonlinear system of equations for  $(\hat{p}_t, \hat{s}_t)$  directly. Here is what Berry, Levinsohn, and Pakes (1995) say about optimal instruments:

*Unfortunately  $D_j(z)$  is typically very difficult, if not impossible, to compute. To calculate  $D_j(z)$  we would have to calculate the pricing equilibrium for different  $(\xi_j, \omega_j)$  sequences, take derivatives at the equilibrium prices, and then integrate out over the distribution of such sequences. In addition, this would require an assumption that chooses among multiple equilibria when they exist, and either additional assumptions on the joint distribution of  $(\xi, \omega)$ , or a method for estimating that distribution.*

The appendix of the NBER working paper version of Berry, Levinsohn, and Pakes (1999) is even less positive:

*Calculating a good estimate of  $E[p|z]$  then requires (i) knowing or estimating the density of the unobservables and (ii) solving at some initial guess for  $\theta$  the fixed point that defines equilibrium prices for each  $(\xi, \omega)$  and then integrating this implicit function with respect to the density of the unknown parameters. This process is too complicated to be practical.*

In Algorithm 2, we follow the possibly more accurate but costly recipe proposed by Berry, Levinsohn, and Pakes (1999) and show that with other computational advances in PyBLP it is feasible to implement.

In general, “asymptotic” and “empirical” approaches are not believed to be computationally feasible, particularly when there are large numbers of draws for  $(\xi_t^*, \omega_t^*)$ . The costly step is Item 4 above, which involves solving for a new equilibrium  $(\hat{p}_t, \hat{s}_t)$  for each set of draws. These improved optimal instruments are feasible primarily because of the advances we describe at the end of Section 3, which drastically reduce the amount of time it takes to solve for equilibria. For relatively large problems, constructing optimal instruments may take several minutes. For smaller problems such as Berry, Levinsohn, and Pakes (1995) or Nevo (2000b) it takes only several seconds. Our simulations indicate that “approximate” performs as well as the more expensive options. Updating results with optimal instruments in PyBLP requires only the last two lines of code in Figure 6.

This approach is not without its limitations. It requires both a way to generate  $(\xi_t^*, \omega_t^*)$  and an estimate of its covariance. This is straightforward if  $(\xi_t^*, \omega_t^*)$  are i.i.d. or follow some other known structure (e.g., clustered at the product level), but this requires additional assumptions. Furthermore, the resulting  $2J_t$  dimensional integral may be hard to approximate accurately for distributions of  $(\xi_t^*, \omega_t^*)$  that are highly skewed or otherwise do not lend well to approximation.

<sup>72</sup> Here  $f(\cdot)$  denotes the joint distribution of unobservables  $(\xi_t, \omega_t)$  and is not to be confused with the distribution of heterogeneity  $f(\mu_{jt}, \tilde{\theta}_2)$ .

<sup>73</sup> Reynaert and Verboven (2014) suggest computing the approximation to the optimal instruments under an additional assumption of perfect competition, so that  $E[p_{jt} | Z_t] = E[c_{jt} | Z_t] = [x_{jt}, w_{jt}] \gamma + \hat{\omega}_{jt}$ . Likewise Gandhi and Houde (2019) suggest using  $E[p_{jt} | Z_t] = \hat{p}_{jt}$  via linear or nonlinear regression in order to construct  $d_{jkt} = \hat{p}_{kt} - \hat{p}_{jt}$ .



**Algorithm 2 Feasible Approximation to Optimal IV (Berry, Levinsohn, and Pakes, 1999)**

After obtaining an initial estimate  $\hat{\theta} = [\hat{\beta}, \hat{\alpha}, \hat{\theta}_2, \hat{\gamma}]$ , for each market  $t$  we can:

1. Obtain an initial estimate of  $\hat{\Omega}_{jt}^{-1}$  by computing covariances of  $(\hat{\xi}_{jt}, \hat{\omega}_{jt})$ . This can be i.i.d. or clustered at any desired level.
2. Draw the  $J_t \times 2$  matrix of structural errors  $(\xi_t^*, \omega_t^*)$  according to one of the below options.
3. Compute  $\hat{Y}_{jt}^S = \hat{c}_{jt} = [x_{jt}, w_{jt}]\hat{\gamma} + \omega_{jt}^*$  and the exogenous portion of utility  $\hat{Y}_{jt}^D = [x_{jt}, v_{jt}]\hat{\beta} + \xi_{jt}^*$ .
4. Use  $(\hat{Y}_t^D, \hat{Y}_t^S, \mathbf{x}_t, \mathbf{w}_t, \mathbf{v}_t)$  to solve for equilibrium prices and quantities  $(\hat{p}_t, \hat{s}_t)$  with the  $\zeta$ -markup approach in (27). Note that this does not involve any endogenous quantities.
5. Treating  $(\hat{p}_t, \hat{s}_t, \mathbf{x}_t, \mathbf{w}_t)$  as data, solve for  $\hat{\xi}_t = \xi_t(\hat{p}_t, \hat{s}_t, \mathbf{x}_t, \mathbf{w}_t, \mathbf{v}_t, \hat{\theta}_2)$  and  $\hat{\omega}_t = \omega_t(\hat{p}_t, \hat{s}_t, \mathbf{x}_t, \mathbf{w}_t, \mathbf{v}_t, \hat{\theta}_2)$ .
6. Construct the Jacobian terms  $\frac{\partial \hat{\xi}_{jt}}{\partial \theta_2}(\xi_t^*, \omega_t^*)$ ,  $\frac{\partial \hat{\omega}_{jt}}{\partial \theta_2}(\xi_t^*, \omega_t^*)$ , and  $\hat{D}_{jt}(\xi_t^*, \omega_t^*)$  using the analytic formulas in Appendix A.2.
7. Average  $\hat{D}_{jt}(\xi_t^*, \omega_t^*)$  over several draws of  $(\xi_t^*, \omega_t^*)$  to construct an estimate of  $E[\hat{D}_{jt}|Z_t]$ .

There are three options for generating  $(\xi_t^*, \omega_t^*)$  suggested by Berry, Levinsohn, and Pakes (1999) and PyBLP makes all three available:

- (a) “Approximate”: Replace  $(\xi_t^*, \omega_t^*)$  with their expectation:  $(E[\xi_t], E[\omega_t]) = (0, 0)$ . This is what Berry, Levinsohn, and Pakes (1999) do. Because the function is nonlinear, this is a good approximation only if  $(\xi_t, \omega_t)$  are very small.
- (b) “Asymptotic”: Estimate an asymptotic normal distribution for  $(\xi_{jt}, \omega_{jt}) \sim N(0, \hat{\Omega})$  and then draw  $(\xi_t^*, \omega_t^*)$  from that distribution. This assumes that a normal approximation is reasonable.
- (c) “Empirical”: Draw  $(\xi_t^*, \omega_t^*)$  from the joint empirical distribution of  $(\xi_{jt}, \omega_{jt})$ . This requires an assumption of *exchangeability*.

Finally, this approach leverages the fact that we have correctly specified the supply side. This may be problematic if we assume a multi-product oligopoly price setting game and the true equilibrium is described by collusive pricing, for example.

□ **Demand side only.** Most empirical applications of the BLP approach do not include a supply side, but estimate demand alone. This has some advantages and some disadvantages. One important implication is that we lose the *cross-equation* overidentifying restrictions, though we manage to retain the  $K_3 - K_x$  *exclusion restrictions* for demand from the cost-shifters  $w_{jt}$ . Absent the supply side, we no longer have a model for marginal costs and cannot solve for equilibrium  $(\hat{p}_t, \hat{s}_t)$ . Instead, the user can supply a vector of expected prices  $E[\mathbf{p}_t|Z_t^D]$  or allow PyBLP to construct the vector with a first-stage regression. Item 4 reduces to computing the market shares at the expected prices. The value of the approximate optimal IV then depends on how well price is explained by the exogenous instruments  $Z_t^D$  in the first stage.

One reason to omit the supply side is that including a misspecified supply side may be worse than no supply side at all. The most controversial aspect of the supply side is often the *conduct assumption*  $\mathcal{H}_t$  used to recover the markup  $\eta_{jt}(\mathcal{H}_t, \theta_2)$ , which may not be known to the researcher prior to estimation. The good news is that testing the validity of the supply side moments amounts to a test of over-identifying restrictions. The simplest test involves estimating the full model with supply and demand to obtain  $\hat{\theta}$ , re-estimating the model using only demand moments  $g_D$  to obtain  $g_D(\hat{\theta}_D)$  along with the optimal demand-only weighting matrix  $W_D$ , and then comparing GMM objectives in a Hausman manner (see Newey, 1985):

$$\text{LR} = N[g(\hat{\theta})'Wg(\hat{\theta}) - g_D(\hat{\theta}_D)'W_Dg_D(\hat{\theta}_D)] \sim \chi_{K-K_x}^2. \quad (32)$$

There are of course alternatives based on the LM (Score) and Wald tests, which are also supported by PyBLP.

*Remark 4.* In our Monte Carlo study, we find a substantial additional benefit when incorporating both supply side moments as well as the approximation to the optimal IV. These benefits substantially exceed those of the demand moments with optimal IV or demand and supply moments with other instruments. We offer two explanations: first, the difference between the nonlinear form of  $E[\mathbf{p}_t|Z_t^D]$  under the approximation to the optimal IV and its linear projection; and second, the

value of *cross equation restrictions* in (30), which use different expressions for  $(\frac{\partial \xi_{jt}}{\partial \theta_2}, \frac{\partial \omega_{jt}}{\partial \theta_2})$ . In our Monte Carlo exercises, we are generally able to reject misspecified conduct assumptions, but not correctly specified ones.

## 5. Monte Carlo experiments

■ Here we provide some Monte Carlo experiments to illustrate some of the best practices laid out in Sections 3 and 4.

□ **Monte Carlo configuration.** Our simulation configurations are loosely based on those of Armstrong (2016). The distinguishing feature is that we randomly draw utilities and costs, but solve for equilibrium prices and shares.<sup>74</sup> Below, we describe our baseline configurations, and in the following sections describe how we modify these configurations to compare different aspects of the problem.

For each configuration, we construct and solve 1000 different synthetic datasets. In each of  $T = 20$  markets we randomly choose the number of firms from  $F_t \in \{2, 5, 10\}$  and have each firm produce a number of products chosen randomly from  $J_{ft} \in \{3, 4, 5\}$ . This procedure generates variation in the number of firms and products across markets, which provides variation in our instruments. Sample sizes are generally between  $200 < N < 600$ .

We draw the structural error terms  $(\xi_{jt}, \omega_{jt})$  from a mean-zero bivariate normal distribution with variances  $\sigma_\xi^2 = \sigma_\omega^2 = 0.2$  and covariance  $\sigma_{\xi\omega} = 0.1$ . Linear demand characteristics are  $[1, x_{jt}, p_{jt}]$  and supply characteristics are  $[1, x_{jt}, w_{jt}]$ . The one nonlinear characteristic is  $x_{jt}$  and heterogeneity is parameterized by  $\mu_{ijt} = \sigma_x x_{jt} v_{it}$  where we draw  $v_{it}$  from the standard normal distribution for 1000 different individuals in each market. We draw the two exogenous characteristics  $(x_{jt}, w_{jt})$  from the standard uniform distribution and compute the endogenous  $(p_{jt}, s_{jt})$  with the  $\zeta$ -markup approach in (27). Demand-side parameters,  $[\beta_0, \beta_x, \alpha] = [-7, 6, -1]$  and  $\sigma_x = 3$ , generate realistic outside shares generally between  $0.8 < s_{0t} < 0.9$ . Supply-side parameters,  $[\gamma_0, \gamma_x, \gamma_w] = [2, 1, 0.2]$ , enter into a linear functional form for marginal costs:  $c_{jt} = [1, x_{jt}, w_{jt}]\gamma + \omega_{jt}$ . For our baseline specification,  $\text{Corr}(p_{jt}, w_{jt}) \approx 0.2$  is relatively low, which implies that our cost-shifting instruments are relatively weak.<sup>75</sup>

In our different Monte Carlo experiments we modify this baseline problem in a number of ways. In most experiments, we consider three variants:

- (a) “Simple” is the baseline problem described above.
- (b) “Complex” adds a random coefficient on price:  $\sigma_p = 0.2$ . Nonlinear coefficients are  $[x_{jt}, p_{jt}]$ .
- (c) “RCNL” adds a nesting parameter:  $\rho = 0.5$ . Each of the  $J_{ft}$  products produced by a firm is randomly assigned to one of  $H = 2$  nesting groups.

We estimate two broad classes of models: demand-only ones, which we estimate with single equation GMM, and models that also include the supply side, which we estimate with multiple equation GMM.

To numerically integrate choice probabilities, we use Gauss-Hermite product rules that exactly integrate polynomials of degree 17 or less.<sup>76</sup> In some specifications, we compare quadrature with the other integration methods described in Section 3.

<sup>74</sup> The first Monte Carlo studies to solve for price and quantity as an equilibrium for BLP-type models are likely Skrainka (2012a), Armstrong (2016), and Conlon (2017). Without solving for equilibrium prices and quantities when generating the data, markups are not “endogenous” and the relevance condition for many (BLP) IV are violated.

<sup>75</sup> With very strong cost shifters or when the variance of  $(\xi_{jt}, \omega_{jt})$  is very small, the estimator always performs well and we can obtain very low bias and median absolute error for nearly any choice of instruments.

<sup>76</sup> In the Simple specification when there is only one dimension of integration, the product rule has  $(17 + 1)/2 = 9$  nodes. In the Complex specification there are  $9^2 = 81$  nodes. If our simulations were of higher dimension, it would be more efficient to use sparse grid integration.

TABLE 2 Various Forms of Instrumental Variables

$$\begin{aligned}
Z_{jt}^{Own} &= \{1, x_{jt}, w_{jt}, x_{jt}^2, w_{jt}^2, x_{jt} \cdot w_{jt}\} \\
Z_{jt}^{Sums} &= \{Z_{jt}^{Own}, \sum_{k \in J_{ft} \setminus \{j\}} 1, \sum_{k \in J_{ft} \setminus \{j\}} x_{kt}, \sum_{k \in J_{ft} \setminus \{j\}} x_{kt}^2\} \\
Z_{jt}^{Local} &= \{Z_{jt}^{Own}, \sum_{k \in J_{ft} \setminus \{j\}} 1(|d_{jkt}| < SD(d)), \sum_{k \in J_{ft} \setminus \{j\}} 1(|d_{jkt}| < SD(d))\} \\
Z_{jt}^{Quad} &= \{Z_{jt}^{Own}, \sum_{k \in J_{ft} \setminus \{j\}} d_{jkt}^2, \sum_{k \in J_{ft} \setminus \{j\}} d_{jkt}^4\}
\end{aligned}$$

Note: Optimal instruments,  $Z_{jt}^{Opt,D}$  and  $Z_{jt}^{Opt,D}$  in (31), are approximated with Algorithm 2 using initial estimates from a single GMM step under  $Z_{jt}^{Sums}$ . We define the difference in characteristic space from product  $j$  as  $d_{jkt} = d_{kt} - d_{jt}$  for each characteristic in  $x_{jt}$ . For the Complex simulation we also include a measure of expected price  $E[p_{jt}|Z_t]$  as an additional  $x_{jt}$ : fitted values from a linear regression of endogenous prices onto all exogenous variables, including the above instruments. For the RCNL simulation, we also include counts of products within the same nest in each market.

To solve the standard fixed point for  $\delta_t$  in each market, we use the SQUAREM acceleration method of Varadhan and Roland (2008) with a  $L^\infty$  tolerance of  $1E-14$  and limit the number of contraction evaluations to 1000. When evaluating the multinomial logit function, we use the log-sum-exp function described in Section 3 to improve numerical stability. During the first GMM step, SQUAREM starts at the solution to the simple logit (or nested logit) model; in the second step, it starts at the estimated first-stage  $\hat{\delta}_t$ .

To optimize, we supply objective values and analytic gradients to a bounded limited-memory BFGS routine (L-BFGS-B), which is made available by the open-source SciPy library. We use an  $L^\infty$  projected gradient-based tolerance of  $1E-5$  and limit the number of major iterations to 1000.<sup>77</sup> Drawing different starting values from a uniform distribution with support 50% above and below the true parameter values, we solve each simulation three times and keep the solution with the smallest objective value. For each simulation, we use box constraints 1000% above and below the true values with the following exceptions:  $\sigma_x \geq 0$ ,  $\sigma_p \geq 0$ ,  $\rho \in [0, 0.95]$ , and when a supply side is estimated,  $\alpha \leq -0.001$ .<sup>78</sup>

Instruments are constructed in two stages. We list the various sets of instruments in Table 2. We begin by exclusively using “own product” characteristics and consider all quadratic interactions of  $(x_{jt}, w_{jt})$ .<sup>79</sup> The “BLP instruments” expand this set to include characteristics of other products. These are meant to be correlated both with the endogenous markup  $\eta_{jt}(\theta_2, x_t, w_t)$ , as well as the inverse mean utility  $D_t^{-1}(\mathcal{S}_t, \tilde{\theta}_2)$ , both of which depend on characteristics of all products in market  $t$ . The second set of instruments,  $Z_{jt}^{Sums}$ , incorporates sums of characteristics of other products, and separates products owned by the same brand  $f$  from products owned by competing brands in the same market  $t$ . The next two sets of instruments consider the *differentiation IV* of Gandhi and Houde (2019). These are variants on the BLP instruments as they represent different functions of own and rival product characteristics. Pooling products across all markets, for each pair of products  $(j, k)$  we construct the difference  $d_{jkt} = x_{kt} - x_{jt}$  of the exogenous regressor. The *differentiation IV* come in two flavors: *Local* and *Quadratic*. The Local measure,  $Z_{jt}^{Local}$ , counts the number of products within a standard deviation of product  $j$ , whereas the Quadratic measure,  $Z_{jt}^{Quad}$ , sums up the aggregate distance between  $j$  and other products. Following Gandhi and Houde (2019), for the Complex simulation where there is a random coefficient on price, we construct an additional instrument using fitted values from a regression of endogenous prices onto all exogenous variables, including the constructed instruments above. For the RCNL simulation, we follow the suggestions of Berry (1994) or Gandhi and Houde (2019) and include the number of products within the same nest and market.

<sup>77</sup> MATLAB solvers report absolute tolerances and SciPy solvers report relative ones. This means magnitudes are not directly comparable. We discuss termination tolerances in more detail in Section 6. In practice, we recommend trying different tolerances to ensure that the optimization routine is not terminating prematurely.

<sup>78</sup> When the optimization software considers  $\alpha = 0$ , the matrix of intra-firm demand derivatives  $\Delta_t$  becomes singular.

<sup>79</sup> One way to view this approach is as a sieve basis to approximate the optimal IV (Ai and Chen, 2003).

TABLE 3 Fixed Effect Absorption

Dimensions	Levels	Absorbed	Seconds	Megabytes
1	216	No	56	721
1	216	Yes	20	39
2	$216 \times 216$	No	112	1414
2	$216 \times 216$	Yes	25	43
2	$36 \times 1296$	No	330	2498
2	$36 \times 1296$	Yes	25	43
3	$36 \times 36 \times 36$	No	46	375
3	$36 \times 36 \times 36$	Yes	24	47

Note: This table documents the impact of fixed effect (FE) absorption on estimation speed and memory usage during one GMM step. Reported values are medians across 100 different simulations. When not absorbed, FEs are included as dummy variables. A single FE is absorbed with de-meaning; multiple, with iterative de-meaning, also known as the Method of Alternating Projections (MAP). To accommodate FEs in the Simple simulation, we first set  $T = 6^4$  and  $F_i = J_f = 6$  so that there are  $N = 6^6 = 46,656$  products. We then randomly assign each  $n \in \{1, \dots, N\}$  to a product and add FEs drawn from the standard uniform distribution. The 1-D case has  $\beta_{n \bmod 216}$ . For the “square” 2-D case, we add  $\beta_{n \div 216}$ . The “uneven” 2-D case has  $\beta_{n \bmod 36}$  and  $\beta_{n \div 36}$ . The 3-D case has  $\beta_{n \bmod 36}$ ,  $\beta_{(n \div 36) \bmod 36}$ , and  $\beta_{n \div 216}$ .

Using parameter estimates from a single GMM step under  $Z_{jt}^{Sums}$ , we construct feasible optimal instruments with the “approximate” variant described in Algorithm 2 and compare with the “asymptotic” and “empirical” alternatives. When a supply side is included, we compute the vector of expected prices with the  $\zeta$ -markup approach in (27). Absent a supply side, we estimate  $E[p_t | Z_t^D]$  with a regression of endogenous prices onto all exogenous variables, including constructed instruments.

□ **Monte Carlo results.** When reporting our Monte Carlo results, we focus on the median bias and median absolute error (MAE) of the parameter estimates. In the online appendix we provide additional results measuring the performance of standard errors and various counterfactual predictions, such as elasticities, merger effects, and welfare estimates. Computation was done on the NYU HPC cluster.<sup>80</sup>

□ **Fixed effects.** In Table 3, we add fixed effects to the Simple simulation and compare computational time and memory usage for experiments where we either absorb the fixed effects or include them as dummy variables. As one might expect, absorbing fixed effects dramatically reduces memory requirements by multiple orders of magnitude and can speed up computation by as much as 10 times. As expected, the largest improvements for the two-dimensional case are when one dimension is much larger than the other. Even absorbing relatively small numbers of fixed effects (216 in a single dimension) leads to a substantial reduction in memory usage and computational time. This is a likely advantage of PyBLP going forward.

□ **Fixed point iteration algorithms.** To highlight the effects of different iteration schemes from Section 3, we explore several methods of solving the system of share equations for  $\delta_t(\theta_2)$ . We compare the conventional fixed-point iteration scheme proposed by Berry, Levinsohn, and Pakes (1995) to two alternative methods of solving the system of nonlinear equations without the Jacobian: DF-SANE and SQUAREM. We also compare with two Jacobian-based methods: Powell’s method and Levenberg–Marquardt (LM).<sup>81</sup>

<sup>80</sup> We uniformly distribute our computation across five types of Intel Xeon processors: 71% on a E5-2690 v2 @ 3.00GHz, 13% on a E5-2660 v3 @ 2.60GHz, 10% on a E5-2690 v4 @ 2.60GHz, 5% on a Gold 6248 @ 2.50GHz, and 1% on a Gold 6148 @ 2.40GHz.

<sup>81</sup> Powell’s method and LM are implemented in MINPACK (More, Garbow, and Hillstom, 1980) as HYBRJ and LMDER, respectively. We do not report results for other SciPy root-finding methods such as Broyden’s Method or Anderson acceleration because we found them too slow and unreliable to be worth considering.

We focus mainly on computational burden and report the results in Table 4. For the Simple and Complex simulations we vary the coefficient on the constant term  $\beta_0$ . A smaller value of  $\beta_0$  leads to a larger share for the outside good. As shown by Dubé, Fox, and Su (2012), a smaller outside good share implies a larger value for the contraction's Lipschitz constant, which generally increases the number of required iterations. Several patterns emerge. As we shrink the outside good share from  $s_{0t} = 0.91 \rightarrow 0.27$ , thereby increasing the Lipschitz constant, the number of required iterations increases approximately by a factor of five times for the Simple and Complex simulations. As expected, the Jacobian-based routines are unaffected by variation of the Lipschitz constant.<sup>82</sup> With the SQUAREM iteration scheme the number of iterations increases, but by only 110%.

In general, we find Powell's method reliable for relatively loose tolerances, but it struggled with a tolerance of  $1\text{E-}14$ , and thus we do not recommend it. The DF-SANE algorithm performs substantially better than standard fixed point iteration, but less well than our two preferred algorithms: the SQUAREM fixed point acceleration scheme and the Jacobian-based LM algorithm.

In terms of other speedups, the effects of the SQUAREM iteration scheme are substantial. It reduces the number of iterations between 3-8 times without substantially increasing the cost per iteration. This is because it approximates the Newton step without actually computing the costly Jacobian. It performs well across all settings, but does particularly well in the Berry, Levinsohn, and Pakes (1995, 1999) example which includes a supply side.

The LM algorithm reduces the number of iterations, but at a higher cost per iteration. On some of the more difficult problems with a small outside good share, it performs as much as 10 times faster than fixed point iteration, and at other times roughly as fast as SQUAREM. The speedup is particularly large for the RCNL model, where the contraction is dampened by  $(1 - \rho)$ . When  $\rho \rightarrow 1$  the contraction becomes arbitrarily slow. In our experiments, LM performs somewhat better than quasi-Newton routines were reported to perform in Reynaerts, Varadhan, and Nash (2012).<sup>83</sup>

□ **Fixed point iteration tricks.** We also consider some commonly employed tricks in the literature to speed up the contraction mapping and report our results in Table B1. In all cases, we use the SQUAREM algorithm. We consider two common "tweaks" from the literature: working with  $\exp(\delta_{jt})$  rather than  $\delta_{jt}$  to avoid taking logs and eliminating a division in the share computation, and using a "hot start" where the starting value for the iterative procedure is the  $\delta_{jt}^{n-1}$  that solved the system of equations for the previous guess of  $\theta_2$ . In addition, we consider the potential cost of our overflow safe modification to the log-sum-exp function.

The results in Table B1 are largely underwhelming. Once we adopt SQUAREM acceleration, additional tricks to speed up the problem seem to have little benefit. The "hot start" approach was able to reduce the number of iterations between by between 10-20%, so it is worth considering.<sup>84</sup> One clear recommendation is the log-sum-exp trick described in Section 3, which reduces the chance of overflow problems. This seems relatively low-cost and reduces the possibility that the iteration routine fails to find a solution to the system of equations.

□ **Numerical integration.** Given extensive past work by Heiss and Winschel (2008), Judd and Skrainka (2011), Skrainka (2012b), and Freyberger (2015), it is not surprising that the choice of integration method has a striking effect. In Figure 1 we add random coefficients to the Simple simulation and compare integration error under pseudo-Monte Carlo

<sup>82</sup> This is consistent with the findings of Dubé, Fox, and Su (2012) using the MPEC method.

<sup>83</sup> One possible explanation is that Reynaerts, Varadhan, and Nash (2012) employ a standard *Newton-Raphson* solver, whereas the MINPACK implementation of Levenberg-Marquardt, LMDER, is designed to be more robust to poor starting values.

<sup>84</sup> This introduces a potential numerical problem where the GMM objective of iteration  $n$  need not evaluate to precisely the same quantity depending on the last iteration's  $\theta^{n-1}$ , although for sufficient fixed point tolerances this should not be an issue. In practice, we still recommend a tight tolerance of  $1\text{E-}14$ .

TABLE 4 Fixed Point Algorithms

Problem	Median $s_{\theta}$	Algorithm	Jacobian	Termination	Mean Milliseconds	Mean Evaluations	Percent Converged
Simple simulation ( $\beta_0 = -7$ )	0.91	Iteration	No	Absolute $L^\infty$	5.95	41.85	100.00%
Simple simulation ( $\beta_0 = -7$ )	0.91	DF-SANE	No	Absolute $L^\infty$	3.43	16.27	100.00%
Simple simulation ( $\beta_0 = -7$ )	0.91	SQUAREM	No	Absolute $L^\infty$	2.56	15.94	100.00%
Simple simulation ( $\beta_0 = -7$ )	0.91	SQUAREM	No	Relative $L^2$	2.75	15.26	100.00%
Simple simulation ( $\beta_0 = -7$ )	0.91	Powell	Yes	Relative $L^2$	3.48	16.33	28.56%
Simple simulation ( $\beta_0 = -7$ )	0.91	LM	Yes	Relative $L^2$	2.31	8.91	100.00%
Simple simulation ( $\beta_0 = -1$ )	0.27	Iteration	No	Absolute $L^\infty$	29.35	212.09	100.00%
Simple simulation ( $\beta_0 = -1$ )	0.27	DF-SANE	No	Absolute $L^\infty$	7.10	35.28	100.00%
Simple simulation ( $\beta_0 = -1$ )	0.27	SQUAREM	No	Absolute $L^\infty$	5.40	34.58	100.00%
Simple simulation ( $\beta_0 = -1$ )	0.27	SQUAREM	No	Relative $L^2$	5.73	33.91	100.00%
Simple simulation ( $\beta_0 = -1$ )	0.27	Powell	Yes	Relative $L^2$	3.67	17.21	11.38%
Simple simulation ( $\beta_0 = -1$ )	0.27	LM	Yes	Relative $L^2$	2.35	8.92	100.00%
Complex simulation ( $\beta_0 = -7$ )	0.91	Iteration	No	Absolute $L^\infty$	8.35	45.02	100.00%
Complex simulation ( $\beta_0 = -7$ )	0.91	DF-SANE	No	Absolute $L^\infty$	4.52	17.67	100.00%
Complex simulation ( $\beta_0 = -7$ )	0.91	SQUAREM	No	Absolute $L^\infty$	3.32	16.14	100.00%
Complex simulation ( $\beta_0 = -7$ )	0.91	SQUAREM	No	Relative $L^2$	3.50	15.50	100.00%
Complex simulation ( $\beta_0 = -7$ )	0.91	Powell	Yes	Relative $L^2$	4.03	14.89	32.29%
Complex simulation ( $\beta_0 = -7$ )	0.91	LM	Yes	Relative $L^2$	2.85	8.98	100.00%
Complex simulation ( $\beta_0 = -1$ )	0.28	Iteration	No	Absolute $L^\infty$	39.35	216.80	100.00%
Complex simulation ( $\beta_0 = -1$ )	0.28	DF-SANE	No	Absolute $L^\infty$	8.92	36.23	100.00%
Complex simulation ( $\beta_0 = -1$ )	0.28	SQUAREM	No	Absolute $L^\infty$	7.03	34.75	100.00%
Complex simulation ( $\beta_0 = -1$ )	0.28	SQUAREM	No	Relative $L^2$	7.51	34.09	100.00%
Complex simulation ( $\beta_0 = -1$ )	0.28	Powell	Yes	Relative $L^2$	4.80	17.70	9.71%
Complex simulation ( $\beta_0 = -1$ )	0.28	LM	Yes	Relative $L^2$	2.90	8.93	100.00%
RCNL simulation ( $\rho = 0.5$ )	0.92	Iteration	No	Absolute $L^\infty$	21.63	93.40	100.00%
RCNL simulation ( $\rho = 0.5$ )	0.92	DF-SANE	No	Absolute $L^\infty$	9.62	32.30	100.00%
RCNL simulation ( $\rho = 0.5$ )	0.92	SQUAREM	No	Absolute $L^\infty$	8.45	33.49	100.00%
RCNL simulation ( $\rho = 0.5$ )	0.92	SQUAREM	No	Relative $L^2$	8.53	31.33	100.00%
RCNL simulation ( $\rho = 0.5$ )	0.92	Powell	Yes	Relative $L^2$	4.67	14.44	60.65%
RCNL simulation ( $\rho = 0.5$ )	0.92	LM	Yes	Relative $L^2$	3.27	8.89	100.00%
RCNL simulation ( $\rho = 0.8$ )	0.92	Iteration	No	Absolute $L^\infty$	58.01	250.34	100.00%
RCNL simulation ( $\rho = 0.8$ )	0.92	DF-SANE	No	Absolute $L^\infty$	16.30	55.50	99.92%

(Continued)

TABLE 4 Continued

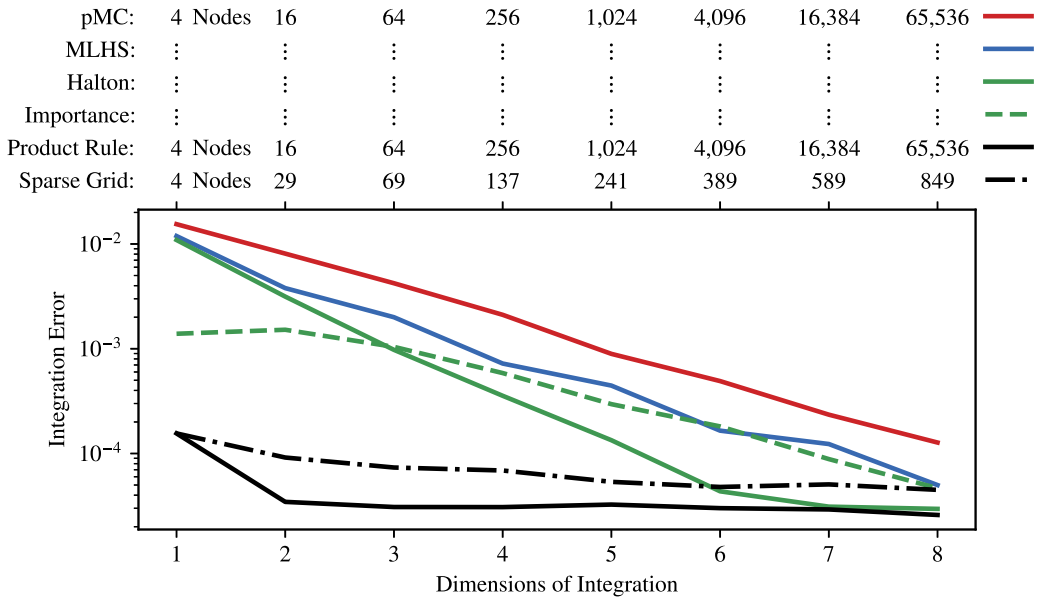
Problem	Median $s_{0t}$	Algorithm	Jacobian	Termination	Mean Milliseconds	Mean Evaluations	Percent Converged
RCNL simulation ( $\rho = 0.8$ )	0.92	SQUAREM	No	Absolute $L^\infty$	14.04	55.54	100.00%
RCNL simulation ( $\rho = 0.8$ )	0.92	SQUAREM	No	Relative $L^2$	13.95	51.70	100.00%
RCNL simulation ( $\rho = 0.8$ )	0.92	Powell	Yes	Relative $L^2$	6.48	20.16	61.48%
RCNL simulation ( $\rho = 0.8$ )	0.92	LM	Yes	Relative $L^2$	3.62	9.64	100.00%
NEVO example	0.54	Iteration	No	Absolute $L^\infty$	12.95	86.64	100.00%
NEVO example	0.54	DF-SANE	No	Absolute $L^\infty$	5.50	25.40	100.00%
NEVO example	0.54	SQUAREM	No	Absolute $L^\infty$	4.05	24.06	100.00%
NEVO example	0.54	SQUAREM	No	Relative $L^2$	4.26	22.80	100.00%
NEVO example	0.54	Powell	Yes	Relative $L^2$	3.84	17.06	29.20%
NEVO example	0.54	LM	Yes	Relative $L^2$	2.62	9.34	100.00%
BLP example	0.89	Iteration	No	Absolute $L^\infty$	152.72	203.04	100.00%
BLP example	0.89	DF-SANE	No	Absolute $L^\infty$	34.93	42.37	100.00%
BLP example	0.89	SQUAREM	No	Absolute $L^\infty$	31.54	40.61	100.00%
BLP example	0.89	SQUAREM	No	Relative $L^2$	31.43	39.38	100.00%
BLP example	0.89	Powell	Yes	Relative $L^2$	26.40	19.31	9.34%
BLP example	0.89	LM	Yes	Relative $L^2$	17.64	8.71	100.00%

Note: This table documents the impact of algorithm choice on solving the nested fixed point. Reported values are medians across 100 different simulations and 10 identical runs of the example problems. We report the number of milliseconds and contraction evaluations needed to solve the fixed point, averaged across all markets and one GMM step's objective evaluations. We also report each algorithm's convergence rate: the percent of times no numerical errors were encountered and a limit of 1000 iterations was not reached. We configure simple iteration with an absolute  $L^\infty$  norm tolerance of  $1E-14$  and compare it with DF-SANE and SQUAREM. The MINPACK (More, Garbow, and Hillstom, 1980) implementations of algorithms that do require a Jacobian—a modification of the Powell hybrid method and Levenberg-Marquardt (LM)—only support relative  $L^2$  norm tolerances, so for comparison's sake we also include SQUAREM with the same termination condition. Simulations are configured as described in Section 5, except for the coefficient on the constant term,  $\beta_0$ , and the nesting parameter,  $\rho$ , which we vary to document the effects of decreasing the outside share  $s_{0t}$  and of dampening the contraction in the RCNL model. The example problems from Nevo (2000b) and Berry, Levinsohn, and Pakes (1995, 1999) are the replications described in Section 6.



FIGURE 1

INTEGRATION ERROR [Color figure can be viewed at wileyonlinelibrary.com]



This plot documents the performance of different numerical integration methods in terms of root mean square error of market shares. Reported values are medians across 100 different markets and are on a logarithmic scale. To document the curse of dimensionality, as we increase the number of random coefficients  $K_2$ , we also increase the number of integration nodes  $I_i$  to match that of a Gauss-Hermite product rule that exactly integrates polynomials of degree 7 or less:  $I_i = 4^{K_2}$ . To accommodate  $K_2 > 1$  dimensions of integration in the Simple simulation, we draw additional exogenous characteristics from the standard uniform distribution. On these new characteristics, we add uncorrelated and mean-zero random coefficients. For comparability's sake, we use the unaltered Simple simulation's market shares  $\mathcal{S}_j$  for all  $K_2$ , and set each random coefficient's variance to  $1/K_2$  so that the distribution of utility is invariant to  $K_2$ . We use 1 million pseudo-Monte Carlo (pMC) draws to precisely compute  $\delta_i = D_i^{-1}(\mathcal{S}_i, \tilde{\theta}_2)$ . With this, we compute the integration error  $\|\mathcal{S}_i - s_i(\delta_i, \tilde{\theta}_2; I_i)\|_2$  of  $s_{ji}(\delta_i, \tilde{\theta}_2; I_i)$  from (24) under various integration methods and nodes  $I_i$ . In addition to a less-precise pMC rule, we also consider Modified Latin Hypercube Sampling (MLHS); Halton draws where in each dimension we use a different prime (2, 3, etc.), discard the first 1000 points, and scramble the sequence according to the recipe in Owen (2017); importance sampling of Halton draws according to the Berry, Levinsohn, and Pakes (1995) procedure described in Section 3 with  $\tilde{\theta}_2$  and  $\delta_i$  are replaced by their true values; a Gauss-Hermite product rule that exactly integrates polynomials of degree 7 or less; and sparse grids with the same polynomial order (Heiss and Winschel, 2008).

(pMC), Modified Latin Hypercube Sampling (MLHS), scrambled Halton draws, importance sampling, and Gauss-Hermite quadrature that exactly integrates polynomials of degree 7 or less.<sup>85</sup> As we increase the dimension of integration  $K_2$ , we decrease the variance of the random coefficients so that the variance of  $\mu_{ijt}(\tilde{\theta}_2)$  remains fixed. As a measure of integration error for  $I_i$  nodes, we compute  $\|\mathcal{S}_i - s_i(\delta_i, \tilde{\theta}_2; I_i)\|_2$  where  $\mathcal{S}_i$  is the vector of true market shares and the vector of mean utilities  $\delta_i = D_i^{-1}(\mathcal{S}_i, \tilde{\theta}_2)$  is precisely computed with one million pMC draws.<sup>86</sup>

To document the curse of dimensionality (CoD), as we add random coefficients we also increase the number of integration nodes  $I_i$  to match the size of the product rule. For the same number of draws, the quadrature routines outperform the qMC and pMC routines. As  $I_i$  increases, the CoD kicks in and the relative performance of pMC, MLHS, Halton draws, and importance

<sup>85</sup> For each dimension of Halton draws we use a different prime (2, 3, etc.), discard the first 1000 points, and scramble the sequence (Owen, 2017). When doing importance sampling, we use Halton draws. In the online appendix we find little difference between importance sampling procedures based on pMC, MLHS, or Halton draws.

<sup>86</sup> In the online appendix we find similar results for a relative measure of integration error.

sampling improves relative to quadrature. Of the nonquadrature methods, scrambled Halton draws perform the best, particularly for a larger number of random coefficients. In our setting, MLHS and importance sampling seem largely underwhelming.<sup>87</sup>

Particularly for a small number of random coefficients, the product rule provides the most accuracy. For more than  $K_2 = 5$  random coefficients, sparse grids provide similar accuracy and require fewer than 10% as many nodes. In Table B2 we confirm that these conclusions translate to performance of parameter estimates. When compared to the product rule, 10 times as many pMC, MLHS, Halton, or importance sampling draws provide similar accuracy, but take 3-20 times as long to estimate.

Our recommendation is to use a product rule for estimating models with only a few random coefficients. For larger numbers of random coefficients we recommend either sparse grids or scrambled Halton draws. However, performance is likely to be context-dependent. After obtaining an estimate of  $\tilde{\theta}_2$ , we recommend verifying that the chosen integration rule performs well relative to feasible alternatives in that setting. This can be done with a procedure like the one used to construct Figure 1: precisely compute  $\delta_i = D_i^{-1}(\mathcal{S}_i, \tilde{\theta}_2)$  using an integration rule with many more nodes than would be feasible during estimation, compute approximations  $s_i(\delta_i, \tilde{\theta}_2; I_i)$  with feasible numbers of integration nodes  $I_i$  for various integration rules, and compare integration errors  $\|\mathcal{S}_i - s_i(\delta_i, \tilde{\theta}_2; I_i)\|_2$ .<sup>88</sup>

□ **Instruments and supply moments.** We also consider several different choices of instruments as well as models which include both supply and demand moments or demand moments only. In Section 4, we present a slightly different construction of the Chamberlain (1987) optimal instruments for the BLP problem than the construction proposed in Reynaert and Verboven (2014). In Table 5, we present simulation results using own characteristics only, sums of characteristics for other products, both the Local and Quadratic forms of the Gandhi and Houde (2019) differentiation IV, and the “approximate” version of the feasible optimal instruments from equation (31).

We find that in most settings the feasible approximation to the optimal instruments performs best, which is consistent with the findings of Reynaert and Verboven (2014). We also find that the differentiation IV outperform the sums of characteristics BLP instruments, as Gandhi and Houde (2019) suggest.

Including moments from a correctly specified supply side substantially improves performance for most sets of instruments with the exception of the “Own” instruments, where the performance gains are limited. We also find that the feasible approximation to the optimal instruments performs much better when a correctly specified supply side is included. In fact, with both optimal instruments and a supply side, the bias is all but eliminated in most of our Monte Carlo experiments, and the MAE is substantially reduced, particularly for the random coefficients. This does not match Reynaert and Verboven (2014) who find that including the supply side has little effect once feasible optimal instruments are used.<sup>89</sup> We highlight the main theoretical difference in Section 4.<sup>90</sup> The gains from including a correctly specified supply side also translate to nonlinear functions of model parameters such as average elasticities and counterfactual simulations such as price effects of mergers. We provide additional details in the online appendix.

<sup>87</sup> In the online appendix we find that importance sampling performs slightly better when there is a larger outside share and a random coefficient on the constant term, although this is sensitive to how we measure integration error. Importance sampling is expected to perform worse when done at an estimate of  $\tilde{\theta}_2$  instead of the true parameters.

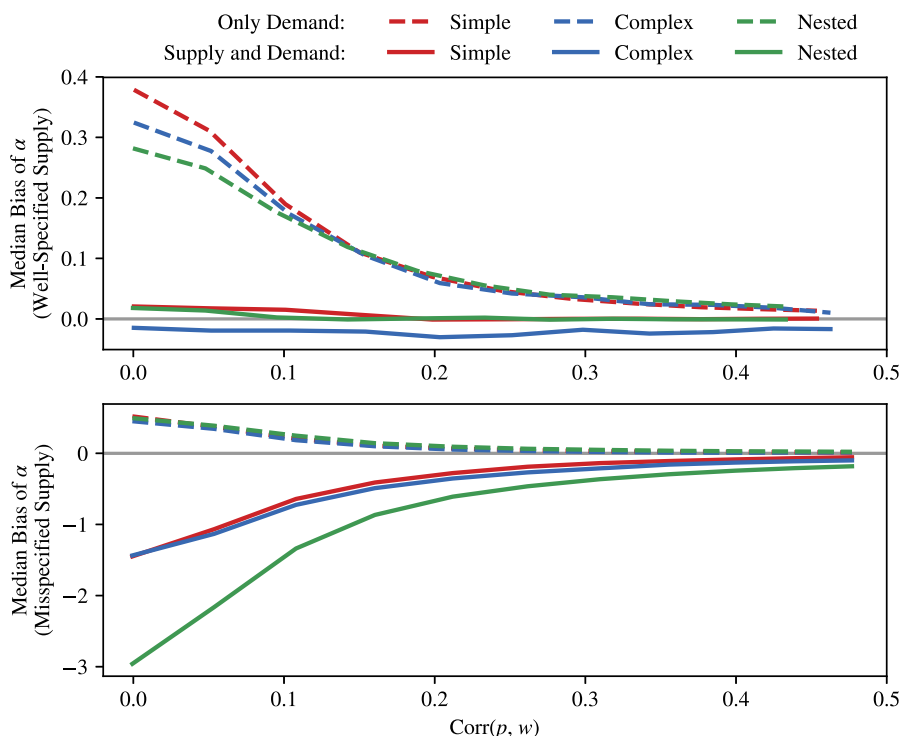
<sup>88</sup> PyBLP makes this type of procedure quick to set up with post-estimation methods that compute mean utilities and market shares under different integration configurations.

<sup>89</sup> A likely important distinction is that their simulations all appear to include “strong cost shifters” in which case all estimators perform quite well, whereas our base specification includes a “weak cost shifter”  $w_{ji}$  with  $\text{Corr}(p_{ji}, w_{ji}) \approx 0.2$ .

<sup>90</sup> See Remark 3. The punchline is that  $E[p_i|Z_i]$  can be approximated with a linear projection onto instruments, or with the nonlinear estimate from solving for  $(p_i^*, s_i^*)$  in equilibrium. The addition of the supply side to the optimal IV problem facilitates the latter.



FIGURE 2

INSTRUMENT STRENGTH AND MISSPECIFICATION [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

Each plot documents how bias of the linear parameter on price,  $\alpha$ , decreases with the strength of the cost shifter  $w_{jt}$ , which is included as a demand-side instrument. To weaken or strengthen the instrument, we vary its supply-side parameter from  $\gamma_w = 0$  to  $\gamma_w = 1$ , and report the correlation this induces between  $w_{jt}$  and prices  $p_{jt}$ . Reported bias values are medians across 1000 different simulations. The top plot reports results for the simulation configurations described in Section 5. In the bottom plot, we simulate data according to perfect competition (i.e., prices are set equal to marginal costs instead of those that satisfy Bertrand-Nash first order conditions), but continue to estimate the model under the assumption of imperfect competition. For all problems, we use the “approximate” version of the feasible optimal instruments and a Gauss-Hermite product rule that exactly integrates polynomials of degree 17 or less.

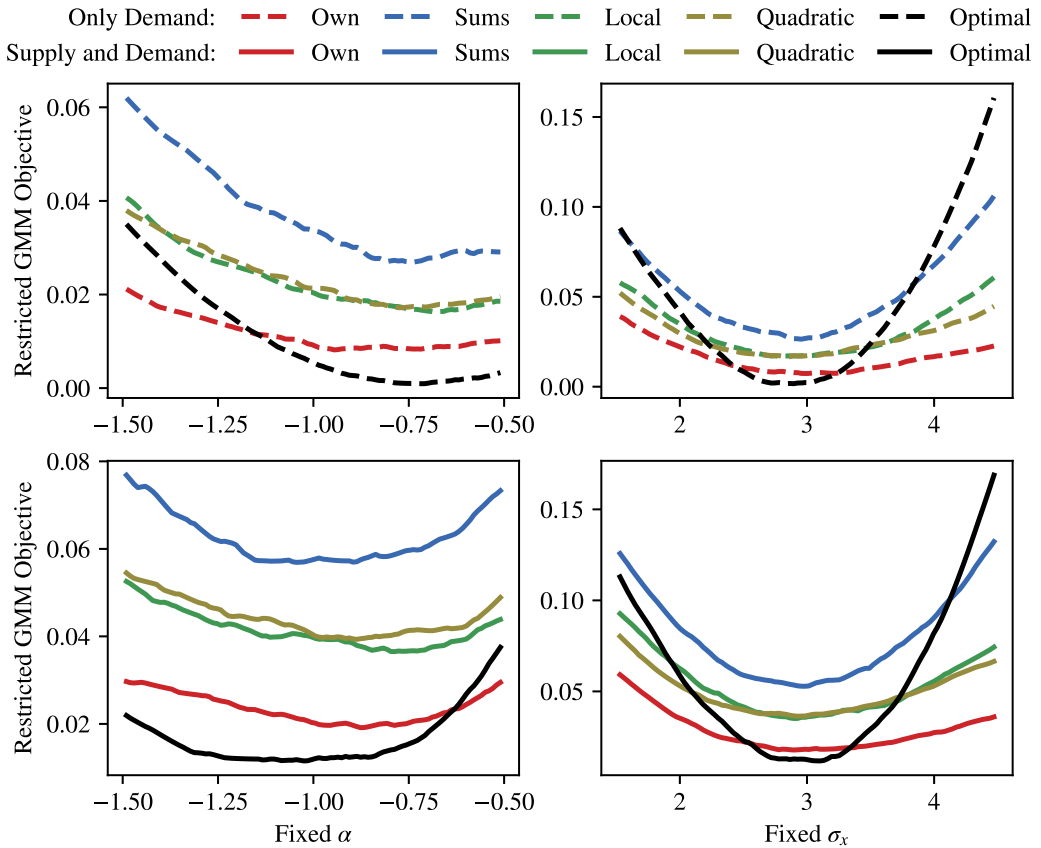
We also show that the incremental value of the supply-side restrictions is largest when the exogenous cost shifter  $w_{jt}$  is weakest. We plot results in Figure 2 where we vary the magnitude of the coefficient  $\gamma_w$ , which governs how responsive marginal costs are to the exogenous cost shifter. Reducing this coefficient reduces the correlation between  $p_{jt}$  and  $w_{jt}$ . When the cost-shifting instrument is very weak with  $\text{Corr}(p_{jt}, w_{jt}) \approx 0.05$ , this increases both the bias and the variance of the price parameter  $\alpha$  consistent with Armstrong (2016). However, if we include the correctly specified supply restrictions (with optimal instruments) we are able to eliminate the bias and substantially reduce the variance of the estimates.<sup>91</sup> This is consistent with the “folklore” around Berry, Levinsohn, and Pakes (1995), where the parameters were difficult to estimate absent supply moments, and it is in contrast to Reynaert and Verboven (2014) who do not find substantial benefits of including supply-side restrictions once optimal instruments are employed.<sup>92</sup> The

<sup>91</sup> We should be cautious because Berry and Haile (2014) suggest that absent any cost-shifters, the model may not be nonparametrically identified. In our Monte Carlo examples, even as the coefficient  $\gamma_w$  becomes small,  $\gamma_x > 0$  so that the observed characteristics do a good job explaining  $p_{jt}$  (the effective “first-stage” for price remains nontrivial).

<sup>92</sup> There are several caveats/differences: simultaneous supply and demand is not the main focus of their work; they construct optimal instruments so that the model is just-identified, whereas our approach constructs them to be over-identified; their excluded cost-shifter is effectively stronger than ours.

FIGURE 3

PROFIED GMM OBJECTIVE WITH ALTERNATIVE INSTRUMENTS [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]



Each plot profiles the GMM objective with respect to a single parameter for the Simple simulation. We fix either  $\alpha$  or  $\sigma_x$ , re-optimize over other parameters, and plot median restricted objective values over 100 different simulations. The left column profiles the objective over the price parameter  $\alpha$ , whereas the right column profiles over the random coefficient  $\sigma_x$ . The top row uses moments from demand alone, whereas the bottom row uses both supply and demand moments. Own instruments are  $[1, x_{jt}, w_{jt}, x_{jt}^2, w_{jt}^2, x_{jt} \cdot w_{jt}]$ . Sums include own and competitor product characteristics. Local and Quadratic instruments follow the definitions in Gandhi and Houde (2019). Optimal instruments are the “approximate” version from Algorithm 2. All instruments are defined in Table 2. For all problems, we use a Gauss-Hermite product rule that exactly integrates polynomials of degree 17 or less.

second panel of Figure 2 suggests that we should be cautious. We generate the data from an assumption of perfect competition, and then impose the Bertrand-Nash multiproduct-oligopoly price setting assumption when construct the supply moments. We illustrate that incorporating moments from an incorrectly specified supply side is worse than not incorporating supply moments because it induces bias in the  $\alpha$  parameter.

We find that that estimates are less sensitive to the precise method used to compute the feasible approximation to the optimal instruments, and we report those results in Table B3. We also find that computing feasible optimal IV with the recipe of Berry, Levinsohn, and Pakes (1999) under the wrong model of firm conduct outperforms computing  $E[p_{jt}|Z_t]$  via a “first stage” linear regression when *only* demand moments are used. We document this in Figure B1.

We further illustrate these advantages in Figure 3 where we plot the profiled GMM objective function using our Simple simulation. In this exercise, we hold fixed either  $\alpha$  or  $\sigma_x$  and

re-optimize the GMM objective over the other parameters. We then plot the profiled objective over a grid of values for  $\alpha$  or  $\sigma_x$ . We repeat this exercise for various sets of instruments, and with or without supply moments. The resulting plots indicate that the approximation to the optimal instruments and the inclusion of the supply moments makes the resulting objective function steeper about the minimum. This suggests that stronger instruments may aid both in numerical optimization and parametric identification of parameters as well as improved efficiency.<sup>93</sup> The minimum under the optimal instruments is generally closer to zero. This allows us to reject misspecified models of supply and fail to reject correctly specified models with the LR test in (32) under the approximation to the optimal IV.

Consistent with Gandhi and Houde (2019), our recommendation is to start with differentiation IV in a first stage including some version of “expected price” and, assuming firm conduct is known, to compute feasible optimal instruments in a second stage. The substantial small sample benefits of including optimal instruments suggest that they should be employed more widely, particularly when there are multiple random coefficients.

□ **Problem size.** We might also be interested in how the BLP problem scales as we vary its size. We display our results in Figure 4. As we increase the number of markets  $T$ , we can substantially improve both the bias and the efficiency of our estimates for both  $\alpha$  and  $\sigma_x$ . For a small number of markets and without additional moments from the supply side, the bias in  $\alpha$  can be substantial. For more than  $T > 100$  markets we get similar (asymptotic) performance for models with and without the additional supply-side restrictions.

In the online appendix we describe how the both the computational time and econometric performance of the estimator is affected by the scale of the problem. We find that computational time is roughly linear in the number of markets  $T$ , and that it grows at a rate closer to  $\sqrt{J_i}$  as we increase the number of products. When we include both supply and demand the computational time appears to grow more quickly than  $J_i$ .

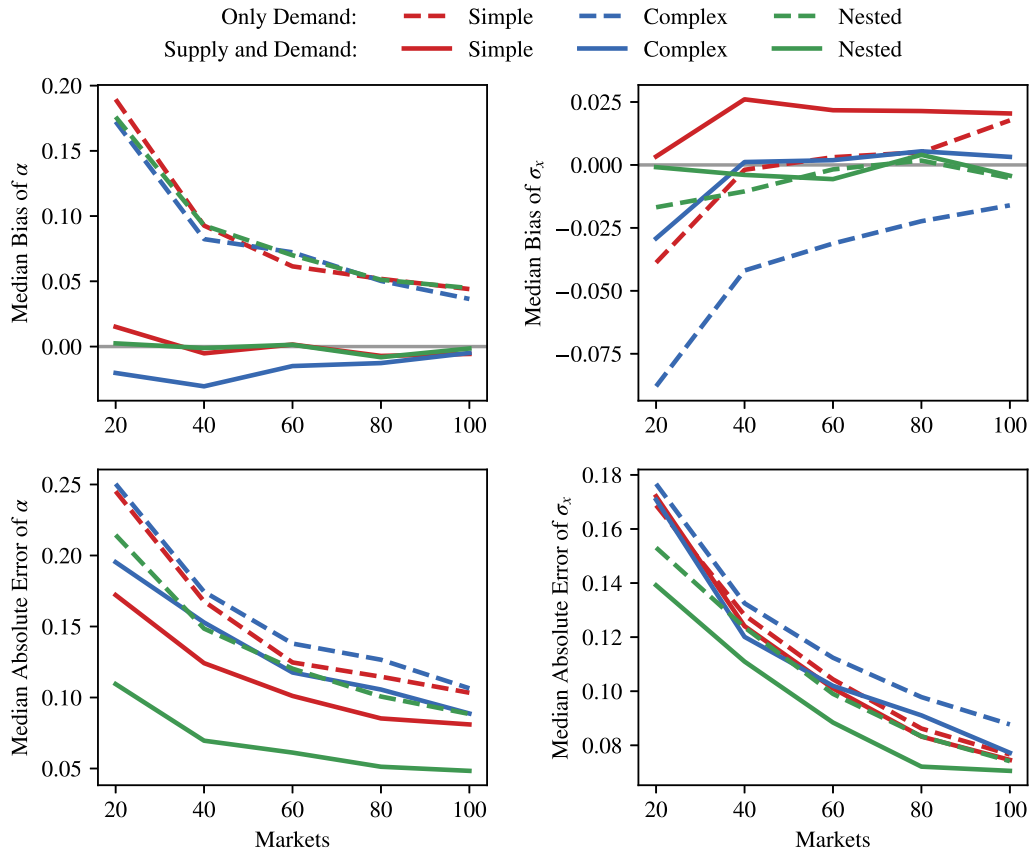
Our Monte Carlo exercises are fairly simple and benefit from variation in the number of products per market, but when using only demand-side restrictions, it appears as if  $T = 40$  markets is roughly enough to obtain “reasonable” parameter estimates in terms of bias and efficiency. These results are somewhat in contrast to those in Armstrong (2016) who finds that when  $T$  is small, as  $J_i$  becomes large the instruments become weak and the estimator performs poorly. We find that when we increase the size of  $J_i$  the estimator performs better rather than worse. We attribute this to three main differences: we allow for variation in the number of products per firm across markets, we use the feasible approximation to the optimal instruments, and in some specifications we include the additional supply moments.

□ **Optimization algorithms.** We consider an array of different optimization algorithms and report our results regarding convergence in Table 6. With the exception of the derivative-free Nelder-Mead algorithm, most of the optimizers reliably find an optimum that satisfies first and second order conditions. Our preferred algorithms are Knitro’s Interior/Direct algorithm and BFGS-based algorithms in SciPy as they provide the best speed and reliability. In the online appendix we consider an even wider array of algorithms and report parameters estimates.

Our findings are somewhat different from those of Knittel and Metaxoglou (2014), who report that different optimization algorithms find a variety of local minima and often fail to converge to a valid minima. We obtain essentially the opposite result. For all derivative-based optimization algorithms in Table 6, we find that more than 99% of all simulation runs converge to a local minimum, which we define as having an  $L^\infty$  norm of the gradient sufficiently close to zero and a positive semi-definite Hessian matrix. We should caution that our simulations are simple enough to be run thousands of times, so it may not be surprising that most optimization

<sup>93</sup> In practice, “weak identification” in nonlinear models arises when the objective function becomes flat with respect to the parameters (Stock and Wright, 2000).

FIGURE 4

PROBLEM SCALING [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

These plots document how bias and variance of parameter estimates decrease with the number of markets  $T$ . The top row plots median bias across 1000 simulations, whereas the bottom row plots median absolute error across the same simulations. The left column reports results for  $\alpha$  and the right column reports results for  $\sigma_x$ . For all problems, we use the “approximate” version of the feasible optimal instruments and a Gauss-Hermite product rule that exactly integrates polynomials of degree 17 or less.

software packages appear to work well in low dimensions. It may also be the case that various numerical fixes and improvements to the fixed point iteration problem may have resolved some of the issues with optimization.<sup>94</sup> In general, strong instruments and minimal numerical error lead to steep and smooth objective functions, which are easier to optimize.

In practice, we recommend placing box constraints on parameters, using gradient-based optimization routines, and setting tight optimization tolerances. We also recommend trying multiple optimizers and starting values to check for agreement.

## 6. Replication exercises

- Here we provide replications using PyBLP for well-known BLP applications.

<sup>94</sup> We use the approximation to the optimal instruments, which tends to make the objective steep about the maximum as demonstrated in Figure 3. We also substantially reduce numerical error by using a Gauss-Hermite product rule that exactly integrates polynomials of degree 17 or less.



TABLE 6 Optimization Algorithms

Simulation	Supply	$ \theta_2 $	Software	Algorithm	Gradient	Termination	Percent of Runs		Median, First GMM Step		
							Converged	PSD Hessian	Seconds	Evaluations	$q = \bar{g}'W\bar{g}$ $  \nabla q  _\infty$
Simple	No	1	Knitro	Interior/Direct	Yes	$  \nabla q  _\infty$	100.0%	100.0%	0.2	4	1.10E-08 8.30E-07
Simple	No	1	SciPy	L-BFGS-B	Yes	$  \nabla q  _\infty$	100.0%	100.0%	0.2	4	8.19E-09 7.28E-07
Simple	No	1	SciPy	BFGS	Yes	$  \nabla q  _\infty$	100.0%	100.0%	0.6	11	1.58E-08 1.03E-06
Simple	No	1	SciPy	TNC	Yes	$  \theta_2^n - \theta_2^{n-1}  _\infty$	99.9%	99.8%	0.6	10	3.89E-24 9.61E-15
Simple	No	1	SciPy	Nelder-Mead	No	$  \theta_2^n - \theta_2^{n-1}  _\infty$	66.5%	100.0%	19.6	115	1.08E-24 4.70E-15
Simple	Yes	2	Knitro	Interior/Direct	Yes	$  \nabla q  _\infty$	100.0%	100.0%	0.6	5	2.17E-06 3.54E-06
Simple	Yes	2	SciPy	L-BFGS-B	Yes	$  \nabla q  _\infty$	100.0%	100.0%	0.4	4	2.18E-06 3.32E-06
Simple	Yes	2	SciPy	BFGS	Yes	$  \nabla q  _\infty$	100.0%	100.0%	2.0	11	2.61E-06 5.26E-06
Simple	Yes	2	SciPy	TNC	Yes	$  \theta_2^n - \theta_2^{n-1}  _\infty$	99.8%	100.0%	2.2	18	2.15E-06 5.12E-11
Simple	Yes	2	SciPy	Nelder-Mead	No	$  \theta_2^n - \theta_2^{n-1}  _\infty$	53.3%	100.0%	31.7	251	2.14E-06 9.69E-13
Complex	No	3	Knitro	Interior/Direct	Yes	$  \nabla q  _\infty$	100.0%	96.9%	0.7	6	1.92E-07 4.11E-06
Complex	No	3	SciPy	L-BFGS-B	Yes	$  \nabla q  _\infty$	100.0%	93.6%	0.4	6	1.83E-07 3.59E-06
Complex	No	3	SciPy	BFGS	Yes	$  \nabla q  _\infty$	100.0%	96.5%	1.9	26	2.25E-07 5.14E-06
Complex	No	3	SciPy	TNC	Yes	$  \theta_2^n - \theta_2^{n-1}  _\infty$	100.0%	86.9%	1.8	20	5.08E-20 1.61E-12
Complex	No	3	SciPy	Nelder-Mead	No	$  \theta_2^n - \theta_2^{n-1}  _\infty$	54.0%	74.6%	30.1	275	1.18E-24 1.39E-14
Complex	Yes	4	Knitro	Interior/Direct	Yes	$  \nabla q  _\infty$	100.0%	93.7%	1.9	9	3.20E-06 6.11E-06
Complex	Yes	4	SciPy	L-BFGS-B	Yes	$  \nabla q  _\infty$	100.0%	94.1%	1.4	9	3.12E-06 5.57E-06
Complex	Yes	4	SciPy	BFGS	Yes	$  \nabla q  _\infty$	100.0%	94.2%	4.6	28	3.19E-06 6.36E-06
Complex	Yes	4	SciPy	TNC	Yes	$  \theta_2^n - \theta_2^{n-1}  _\infty$	99.5%	99.5%	5.7	31	2.87E-06 4.02E-10
Complex	Yes	4	SciPy	Nelder-Mead	No	$  \theta_2^n - \theta_2^{n-1}  _\infty$	45.5%	99.5%	64.6	480	2.80E-06 1.83E-12
RCNL	No	2	Knitro	Interior/Direct	Yes	$  \nabla q  _\infty$	100.0%	100.0%	1.7	10	1.67E-08 5.72E-06
RCNL	No	2	SciPy	L-BFGS-B	Yes	$  \nabla q  _\infty$	100.0%	99.9%	1.9	11	1.52E-09 1.47E-06
RCNL	No	2	SciPy	BFGS	Yes	$  \nabla q  _\infty$	100.0%	100.0%	4.3	25	2.64E-09 1.95E-06
RCNL	No	2	SciPy	TNC	Yes	$  \theta_2^n - \theta_2^{n-1}  _\infty$	100.0%	99.6%	4.1	22	3.56E-19 1.16E-11
RCNL	No	2	SciPy	Nelder-Mead	No	$  \theta_2^n - \theta_2^{n-1}  _\infty$	56.4%	98.7%	60.5	243	1.27E-25 2.53E-14
RCNL	Yes	3	Knitro	Interior/Direct	Yes	$  \nabla q  _\infty$	100.0%	100.0%	3.4	13	2.83E-06 9.95E-06
RCNL	Yes	3	SciPy	L-BFGS-B	Yes	$  \nabla q  _\infty$	100.0%	100.0%	3.2	12	2.66E-06 3.85E-06
RCNL	Yes	3	SciPy	BFGS	Yes	$  \nabla q  _\infty$	100.0%	100.0%	9.3	37	2.73E-06 4.41E-06

(Continued)

TABLE 6 Continued

Simulation	Supply	$ \theta_2 $	Software	Algorithm	Gradient	Termination	Percent of Runs		Median, First GMM Step	
							Converged	PSD Hessian	Seconds	Evaluations
RCNL	Yes	3	SciPy	TNC	Yes	$\ \theta_2^n - \theta_2^{n-1}\ _\infty$	100.0%	100.0%	7.2	24
RCNL	Yes	3	SciPy	Nelder-Mead	No	$\ \theta_2^n - \theta_2^{n-1}\ _\infty$	39.0%	100.0%	115.2	423

Note: This table documents optimization convergence statistics over 1000 simulated datasets for different optimization algorithms. For comparison's sake, we only perform one GMM step and do not set parameter bounds. We report each algorithm's convergence rate (the percent of times the algorithm reported that it successfully found an optimum before reaching 1000 iterations) and the percent of times the final Hessian matrix was positive semidefinite. We also report medians for the number of seconds needed to solve each problem, the number of objective evaluations, the final GMM objective value, and the  $L^\infty$  norm of the gradient. We configure three algorithms to terminate with gradient-based  $L^\infty$  norms of 1E-5: Interior/Direct from Knitro, along with L-BFGS-B (limited-memory BFGS) and BFGS from SciPy. Because SciPy's derivative-free Nelder-Mead algorithm does not support gradient-based termination, we instead configure it to terminate with an absolute parameter-based  $L^\infty$  norm of 1E-5 and for comparison's sake use the same configuration for SciPy's TNC (truncated Newton) algorithm. For all problems, we use the "approximate" version of the feasible optimal instruments and a Gauss-Hermite product rule that exactly integrates polynomials of degree 17 or less.

FIGURE 5

NEVO (2000b) REPLICATION CODE [Color figure can be viewed at wileyonlinelibrary.com]

```

import numpy as np
import pandas as pd

import pyblp

problem = pyblp.Problem(
    product_formulations=(
        pyblp.Formulation('0 + prices', absorb='C(product_ids)'),          # Linear demand
        pyblp.Formulation('1 + prices + sugar + mushy'),                  # Nonlinear demand
    ),
    agent_formulation=pyblp.Formulation('0 + income + income_squared + age + child'), # Demographics
    product_data=pd.read_csv(pyblp.data.NEVO_PRODUCTS_LOCATION),
    agent_data=pd.read_csv(pyblp.data.NEVO_AGENTS_LOCATION)
)

results = problem.solve(
    sigma=np.diag([0.3302, 2.4526, 0.0163, 0.2441]),          # Starting values for unobserved heterogeneity
    pi=[
        [ 5.4819,  0,          0.2037, 0 ],
        [15.8935, -1.2000, 0,      2.6342],
        [-0.2506,  0,          0.0511, 0 ],
        [ 1.2650,  0,          -0.8091, 0 ]
    ],
    method='1s',                                              # One-step GMM
    optimization=pyblp.Optimization('bfgs', {'gtol': 1e-5})  # Gradient-based termination tolerance
)

elasticities = results.compute_elasticities()
markups = results.compute_markups()

```

This Python code demonstrates how to construct and solve the problem from Nevo (2000b) in PyBLP. Names in the formulation objects correspond to variable names in the datasets, which are packaged with PyBLP and in this example are loaded into memory with the Python package `pandas`. Although not an explicit dependency of PyBLP, `pandas` is a convenient package for loading data into data frames, and comes pre-packaged in Anaconda installations. In addition to `pandas` data frames, PyBLP can also handle other data types such as NumPy structured arrays or simple dictionaries. Most estimation outputs are stored as attributes of the problem results class. Post-estimation outputs such as elasticities and markups can be computed with class methods. Estimation results are reported in Table 7.

□ **Nevo (2000b).** In our first replication we estimate the model of Nevo (2000b) on its publicly available “fake data.” This problem is notable because it includes a combination of observable demographics, unobserved heterogeneity, and product fixed effects. We demonstrate how to construct and solve the problem with PyBLP in Figure 5.

We estimate the model three times: once following the original Nevo (2000b) example,<sup>95</sup> a second time with an optimization tolerance of  $1\text{E}-5$  instead of  $1\text{E}-4$ ,<sup>96</sup> and a third time using the demand side only feasible optimal instruments described above. We obtain nearly identical estimates of mean own-price elasticities and markups for each specification despite different parameter estimates for the Price  $\times$  Income terms that are nearly collinear in the data.

□ **Berry, Levinsohn, and Pakes (1995, 1999).** For our second replication, we consider the problem in Berry, Levinsohn, and Pakes (1995), which lacks demographic interactions and product fixed effects but adds a supply side and allows the price coefficient to vary with income. We provide the PyBLP formulation of the problem in Figure 6. This time, we access Python through the R package `reticulate` (Allaire et al., 2017) to demonstrate how PyBLP can be used in R.

<sup>95</sup> Replication is straightforward because the original instruments are provided along with the data and the reported estimates are from one-step GMM with the two-stage least squares (2SLS) weighting matrix  $(Z'Z)^{-1}$ . The data we use are distributed with PyBLP.

<sup>96</sup> It is well known that the code accompanying Nevo (2000b) set tolerances too loose. However, with a scaled GMM objective value of  $Nq(\hat{\theta}) = 4.56$ , our results with a tighter tolerance are identical to those reported by Dubé, Fox, and Su (2012) using the MPEC approach.

FIGURE 6

BERRY, LEVINSOHN, AND PAKES (1995, 1999) REPLICATION CODE [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

```
library(readr)
library(reticulate)

pyblp <- import('pyblp')

problem <- pyblp$Problem(
  product_formulations = tuple(
    pyblp$Formulation('1 + hpwt + air + mpd + space'),      # Linear demand
    pyblp$Formulation('1 + prices + hpwt + air + mpd + space'), # Nonlinear demand
    pyblp$Formulation('1 + log(hpwt) + air + log(mpg) + log(space) + trend') # Supply
  ),
  agent_formulation = pyblp$Formulation('0 + I(1 / income)'), # Price interaction
  costs_type = 'log', # Log-linear costs
  product_data = read_csv(pyblp$data$BLP_PRODUCTS_LOCATION),
  agent_data = read_csv(pyblp$data$BLP_AGENTS_LOCATION)
)

results <- problem$solve(
  sigma = diag(c(3.612, 0, 4.628, 1.818, 1.050, 2.056)), # Starting values for unobserved heterogeneity
  pi = rbind(0, -43.501, 0, 0, 0, 0), # Starting value for the term on price
  initial_update = TRUE, # Update weight matrix at starting values
  costs_bounds = tuple(0.001, NULL), # Constrain marginal costs to be positive
  W_type = 'clustered', # Cluster by automobile model
  se_type = 'clustered'
)

elasticities = results$compute_elasticities()
markups = results$compute_markups()

instrument_results = results$compute_optimal_instruments(method = 'approximate')
updated_problem = instrument_results$to_problem()
```

This R code demonstrates how to construct and solve the problem from Berry, Levinsohn, and Pakes (1995, 1999) in PyBLP. Python functions are called with the R package `reticulate`. Similar packages that allow for Python interoperability are available in many other languages. Names in the formulation objects correspond to variable names in the datasets, which are packaged with PyBLP and in this example are loaded into memory with the R package `readr`. Most estimation outputs are stored as attributes of the problem results class. Post-estimation outputs such as elasticities and markups can be computed with class methods. Here, the “approximate” technique from Algorithm 2 is used to construct feasible optimal instruments. This gives an optimal instruments results object that is converted into an updated problem, which can be solved like any other. Estimation results are reported in Table 8.

The widespread availability of such packages for between-language interoperability means that it is often straightforward to incorporate functionality from other languages.<sup>97</sup>

Our configuration for the Berry, Levinsohn, and Pakes (1995) problem differs more substantially from the original article because parts of the original configuration are not included with the data.<sup>98</sup> We estimate the model once with the original article’s sums of characteristics BLP instruments and importance sampling integration rule, and a second time with feasible optimal instruments and 10,000 scrambled Halton draws in each market. We report our results in Table 8. We obtain broadly similar parameter estimates. Estimates under Halton draws and the feasible approximation to the optimal instruments suggest somewhat less preference heterogeneity, which leads to slightly less elastic demand and larger markups than in the original specification.

□ **Knittel and Metaxoglou (2014).** With the Nevo (2000b) configuration and a demand-only version of the Berry, Levinsohn, and Pakes (1995) problem, we conduct a more extensive

<sup>97</sup> For more on interoperability, see Footnote 3.

<sup>98</sup> We obtain the data from the replication package for Andrews, Gentzkow, and Shapiro (2017) and obtain very similar results as this earlier replication. Following Berry, Levinsohn, and Pakes (1999), we replace the original article’s  $\log(y_i - p_j)$  term with its first-order linear approximation  $p_j/y_i$ . Otherwise, there are individuals for whom  $p_j > y_i$ , creating a host of problems.

TABLE 7 Nevo (2000b) Replication

		Published Estimates	Replication	Tighter Tolerance	Best Practices
Means	Price	-32.433 (7.743)	-32.404 (7.729)	-62.729 (14.803)	-27.489 (4.383)
Standard Deviations	Price	1.848 (1.075)	1.851 (1.070)	3.313 (1.340)	2.910 (0.669)
	Constant	0.377 (0.129)	0.376 (0.129)	0.558 (0.163)	0.196 (0.085)
	Sugar	0.004 (0.012)	0.003 (0.012)	0.006 (0.014)	0.028 (0.008)
	Mushy	0.081 (0.205)	0.080 (0.204)	0.093 (0.185)	0.324 (0.110)
	Interactions	16.598 (172.334)	16.457 (172.237)	588.318 (270.441)	15.957 (98.164)
	Price $\times$ income	-0.659 (8.955)	-0.655 (8.951)	-30.192 (14.101)	-1.282 (5.119)
	Price $\times$ income squared	11.625 (5.207)	11.543 (5.166)	11.054 (4.123)	4.551 (2.405)
	Constant $\times$ income	3.089 (1.213)	3.100 (1.203)	2.292 (1.209)	6.253 (0.541)
	Constant $\times$ age	1.186 (1.016)	1.172 (1.001)	1.284 (0.631)	0.162 (0.207)
	Sugar $\times$ income	-0.193 (0.005)	-0.193 (0.045)	-0.385 (0.121)	-0.289 (0.037)
	Sugar $\times$ age	0.029 (0.036)	0.030 (0.036)	0.052 (0.026)	0.046 (0.014)
	Mushy $\times$ income	1.468 (0.697)	1.462 (0.693)	0.748 (0.802)	0.998 (0.303)
	Mushy $\times$ age	-1.514 (1.103)	-1.502 (1.091)	-1.353 (0.667)	-0.523 (0.188)
Mean own-price elasticity			-3.700	-3.618	-3.685
Mean markup			0.360	0.364	0.363
GMM objective		6.60E-03	6.61E-03	2.02E-03	2.03E-04
GMM objective scaled by $N$		1.49E+01	1.49E+01	4.56E+00	4.59E-01

Note: This table reports replication results for the model of Nevo (2000b) described in Section 6. From left to right, we report estimates from the original article, our replication results, additional results for when we reduce the BFGS optimization algorithm's gradient-based  $L^\infty$  norm from  $1\text{E}-4$  to a slightly tighter  $1\text{E}-5$ , and a final set of results using best estimation practices: a tighter termination tolerance and the "approximate" version of the feasible optimal instruments. Standard errors are in parentheses.

replication exercise meant to mimic Knittel and Metaxoglou (2014). We solve these two problems with 50 random starting values and multiple optimization algorithms.<sup>99</sup>

We first replicate some of the difficulties encountered by Knittel and Metaxoglou (2014) by using very loose optimization tolerances.<sup>100</sup> Early termination of optimization algorithms creates dispersion across recovered parameter estimates, objective values, and implied elasticities across optimizers and starting values. However, tighter optimization tolerances suffice to eliminate all dispersion for the problem in Nevo (2000b).<sup>101</sup> In addition to tight tolerances, using quadrature instead of a small number of pseudo-Monte Carlo draws also eliminates all dispersion in the

<sup>99</sup> For a longer description of our replication exercise, please refer to the online appendix.

<sup>100</sup> Knittel and Metaxoglou (2014) use a tolerance of  $1\text{E}-3$  for changes in the parameter vector and the objective function. Because of the loose objective function tolerance in particular, the optimization routines often terminate too early. We replicate this behavior with loose  $L^\infty$  gradient- and parameter-based tolerances of  $1\text{E}-1$ .

<sup>101</sup> We use  $L^\infty$  gradient- and parameter-based tolerances of  $1\text{E}-4$ . In their online appendix, Knittel and Metaxoglou (2014) also report that when using a tighter tolerance and a gradient-based routine, they manage to eliminate essentially all dispersion for the problem in Nevo (2000b) but not Berry, Levinsohn, and Pakes (1995).

TABLE 8 Berry, Levinsohn, and Pakes (1995, 1999) Replication

		Published Estimates	Replication	Best Practices
Means	Constant	-7.061 (0.941)	-7.284 (2.807)	-6.679 (1.304)
	HP/weight	2.883 (2.019)	3.460 (1.415)	2.774 (0.833)
	Air	1.521 (0.891)	-0.999 (2.101)	0.572 (0.349)
	MP\$	-0.122 (0.320)	0.421 (0.250)	0.340 (0.098)
	Size	3.460 (0.610)	4.178 (0.658)	3.920 (0.322)
		3.612 (1.485)	2.025 (6.065)	2.962 (1.637)
Standard deviations	Constant	3.612 (1.485)	2.025 (6.065)	2.962 (1.637)
	HP/weight	4.628 (1.885)	6.101 (2.200)	1.388 (2.107)
	Air	1.818 (1.695)	3.956 (2.110)	1.424 (0.435)
	MP\$	1.050 (0.272)	0.254 (0.549)	0.072 (1.002)
	Size	2.056 (0.585)	1.908 (1.108)	0.231 (3.837)
		43.501 (6.427)	44.842 (9.216)	45.898 (11.748)
Term on price	$\ln(y - p)$	43.501 (6.427)	44.842 (9.216)	45.898 (11.748)
Supply-side terms	Constant	0.952 (0.194)	2.760 (0.116)	2.785 (0.104)
	$\ln(\text{HP/weight})$	0.477 (0.056)	0.897 (0.072)	0.731 (0.071)
	Air	0.619 (0.038)	0.423 (0.087)	0.528 (0.040)
	$\ln(\text{MPG})$	-0.415 (0.055)	-0.525 (0.073)	-0.651 (0.071)
	$\ln(\text{size})$	-0.046 (0.081)	-0.261 (0.210)	-0.472 (0.125)
	Trend	0.019 (0.002)	0.027 (0.003)	0.018 (0.002)
			-3.928	-3.461
Mean own-price elasticity			0.316	0.346
Mean markup			2.24E-01	1.06E-01
GMM objective			4.97E+02	2.36E+02
GMM objective scaled by $N$				

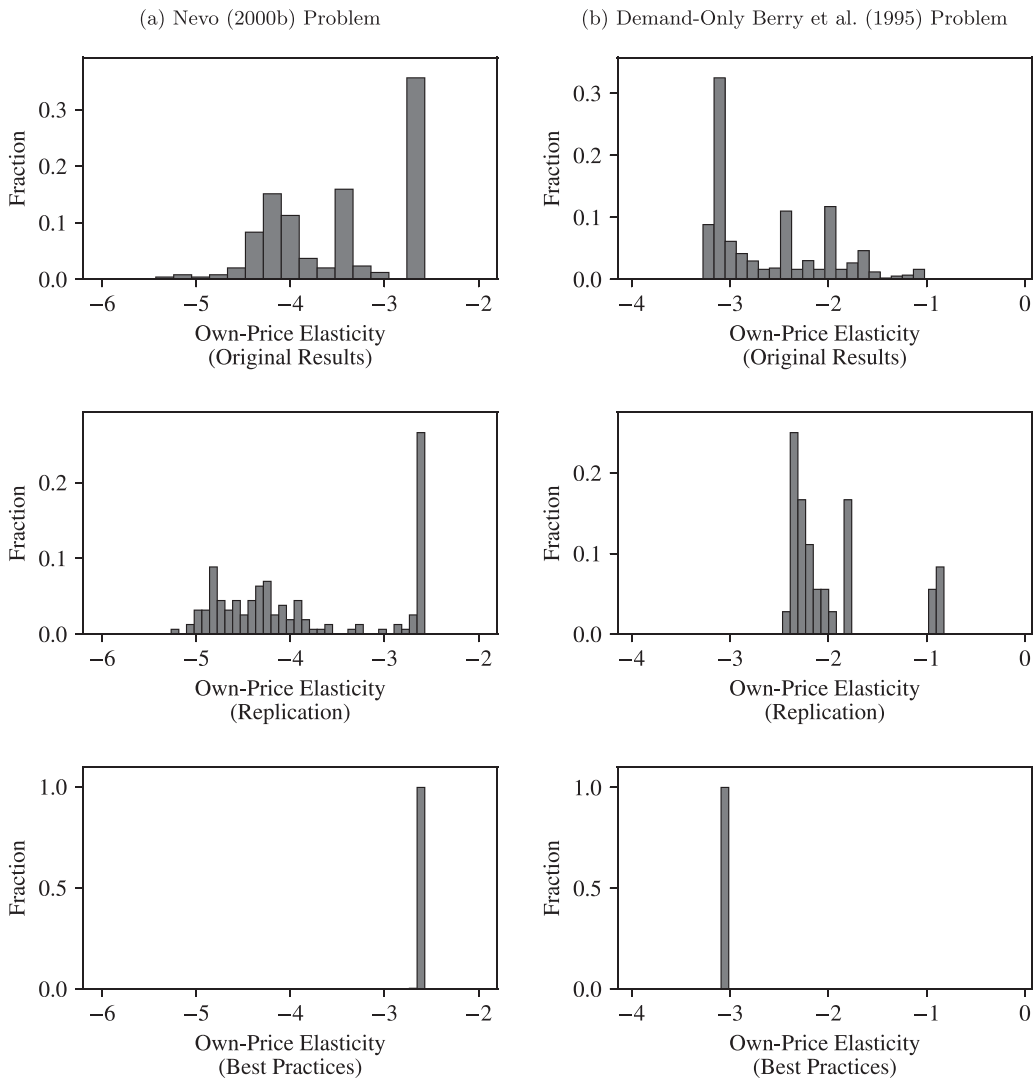
Note: This table reports replication results for the model of Berry, Levinsohn, and Pakes (1995, 1999) described in Section 6. From left to right, we report estimates from the original article, our replication results, and results using best estimation practices: 10,000 scrambled Halton draws in each market and the “approximate” version of the feasible optimal instruments. Standard errors are in parentheses and are clustered by automobile model.

demand-only version of the Berry, Levinsohn, and Pakes (1995) problem. This is consistent with Brunner et al. (2017), who find that simulation error contributes substantial instability to this particular configuration.

In other words, our best practices eliminate much of the difficulties encountered by Knittel and Metaxoglou (2014). In line with our Monte Carlo experiments, when properly configured, the choice of optimization routine (open-source or otherwise) does not seem to be particularly important for the nested fixed point algorithm. To demonstrate this result graphically, in Figure 7 we present histograms of own prices elasticities obtained from our replication exercise below the corresponding figures from Knittel and Metaxoglou (2014). The “Best Practices” histograms exhibit essentially no dispersion across optimization routines and starting values.

FIGURE 7

KNITTEL AND METAXOGLLOU (2014) HISTOGRAMS FOR MEDIAN PRODUCT OWN-PRICE ELASTICITIES [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]



All figures report the own price elasticity for the median product. The three figures on the left are for the problem in Nevo (2000b). Figures on the right are for a demand-only version of the problem in Berry, Levinsohn, and Pakes (1995) described in Knittel and Metaxoglou (2014). The top two figures are reproduced as fair use from Knittel and Metaxoglou (2014). The bottom four are produced by PyBLP where each observation is one trial from 50 starting values and seven optimization algorithm configurations supported by Knitro and SciPy. The middle figures replicate some of the difficulties in Knittel and Metaxoglou (2014) by using loose optimization tolerances. The bottom two figures eliminate these difficulties with best estimation practices. It suffices to use tight optimization tolerances for the problem in Nevo (2000b). Difficulties for the demand-only version of the problem in Berry, Levinsohn, and Pakes (1995) can be eliminated by additionally using a Gauss-Hermite product rule that exactly integrates polynomials of degree 11 or less instead of 50 pseudo-Monte Carlo draws in each market. For additional details, please consult the online appendix.



## 7. Conclusion

■ Our goal has been to review recent methodological developments related to the BLP problem, and to collect and evaluate them not only in a single article, but also in a single software package, PyBLP. We have provided a list of best practices with numerical evidence to support our recommendations, and we have implemented them as defaults in PyBLP. Our hope is that these practices (and estimation of BLP models in general) can now be made available to more researchers, and provide a common platform that should facilitate replication. For researchers who wish to implement and estimate similar models that are not among the wide range of BLP-type models supported by PyBLP, we hope that this article and our well-documented code will serve as a good starting point.

In addition, we present some methodological results that we believe to be novel. We show how with a slight reformulation of the nested fixed point problem, it is possible to include high dimensional fixed effects in models with simultaneous supply and demand. We also provide a somewhat different expression for optimal instruments than the prior literature (Reynaert and Verboven, 2014), which makes clear the over-identifying restrictions implied by the supply side. Also novel, we find that optimal instruments when combined with a correctly specified supply side are extremely valuable. Consistent with prior work, we find the gains to optimal instruments to be substantial such that they should nearly always be employed. Thankfully, we have made this process extremely straightforward in PyBLP.

Somewhat reassuringly, we find that under our best practices, including correctly specified supply restrictions and approximations to the optimal instruments, finite sample performance of the BLP estimator appears to be quite good and perhaps better than previously believed.

## Appendix A: Derivations

□ **Concentrating out linear parameters.** Our objective is to concentrate out  $\hat{\beta}(\theta_2)$  and  $\hat{\gamma}(\theta_2)$ . Define  $Y_{jt}^D$ ,  $Y_{jt}^S$ ,  $X_{jt}^D$ , and  $X_{jt}^S$  as follows:

$$\begin{aligned} Y_{jt}^D &\equiv \hat{\delta}_{jt}(\theta_2) + \alpha p_{jt} = [x_{jt}, v_{jt}]\beta + \xi_{jt} \equiv X_{jt}^D \beta + \xi_{jt}, \\ Y_{jt}^S &\equiv p_{jt} - \hat{\eta}_{jt}(\theta_2) = [x_{jt}, w_{jt}]\gamma + \omega_{jt} \equiv X_{jt}^S \gamma + \omega_{jt}. \end{aligned} \quad (A1)$$

Stacking the system across observations yields:<sup>102</sup>

$$\underbrace{\begin{bmatrix} Y_D \\ Y_S \end{bmatrix}}_{2N \times 1} = \underbrace{\begin{bmatrix} X_D & 0 \\ 0 & X_S \end{bmatrix}}_{2N \times (K_1 + K_2)} \underbrace{\begin{bmatrix} \beta \\ \gamma \end{bmatrix}}_{(K_1 + K_2) \times 1} + \underbrace{\begin{bmatrix} \xi \\ \omega \end{bmatrix}}_{2N \times 1}. \quad (A2)$$

Adding the  $N \times M_D$  instruments for demand,  $Z_D$ , and the  $N \times M_S$  instruments for supply,  $Z_S$ , yields  $M = M_D + M_S$  moment restrictions

$$E \begin{bmatrix} Z_D' (Y_{jt}^D - X_{jt}^D \beta) \\ Z_S' (Y_{jt}^S - X_{jt}^S \gamma) \end{bmatrix} = 0, \quad (A3)$$

which have sample analogues

$$\underbrace{\frac{1}{N} \begin{bmatrix} Z_D' & 0 \\ 0 & Z_S' \end{bmatrix}}_{\tilde{Y}} \underbrace{\begin{bmatrix} Y_D \\ Y_S \end{bmatrix}}_{2N \times 1} - \underbrace{\frac{1}{N} \begin{bmatrix} Z_D' X_D & 0 \\ 0 & Z_S' X_S \end{bmatrix}}_{\tilde{X}} \underbrace{\begin{bmatrix} \beta \\ \gamma \end{bmatrix}}_{(K_1 + K_2) \times 1} \quad (A4)$$

Now we can simply perform a GMM regression of  $\tilde{Y}$  on  $\tilde{X}$  with the same  $M \times M$  weighting matrix  $W$  used in the overall problem:

$$\begin{bmatrix} \hat{\beta}(\theta_2) \\ \hat{\gamma}(\theta_2) \end{bmatrix} = (\tilde{X}' W \tilde{X})^{-1} \tilde{X}' W \tilde{Y}. \quad (A5)$$

□ **Analytic derivative calculations.** *This derivation appears to be novel to the literature for the case of simultaneous estimation of supply and demand.*

<sup>102</sup> Note that cannot perform independent regressions unless we are willing to assume that  $\text{Cov}(\xi_{jt}, \omega_{jt}) = 0$ .

The gradient of the GMM objective function is

$$\nabla q(\theta_2) = 2G(\theta_2)'Wg(\theta_2).$$

The challenging piece here is the Jacobian of the GMM objective:

$$\underbrace{G(\theta_2)}_{M \times K_2} = \frac{1}{N} \underbrace{\begin{bmatrix} Z_D' & 0 \\ 0 & Z_S' \end{bmatrix}}_{M \times 2N} \underbrace{\begin{bmatrix} \frac{\partial \xi}{\partial \theta_2} \\ \frac{\partial \omega}{\partial \theta_2} \end{bmatrix}}_{2N \times K_2}.$$

We can write:<sup>103</sup>

$$\begin{bmatrix} \frac{\partial \xi}{\partial \theta_2} \\ \frac{\partial \omega}{\partial \theta_2} \end{bmatrix} = \begin{bmatrix} \frac{\partial \delta}{\partial \theta_2} \\ -f'_{MC}(\cdot) \frac{\partial \eta}{\partial \theta_2} \end{bmatrix}. \quad (\text{A6})$$

The  $f'_{MC}(\cdot)$  in (A6) comes from (7) where  $f_{MC}(\cdot)$  is typically linear with  $f'_{MC}(c_{jt}) = 1$ , or logarithmic with  $f'_{MC}(c_{jt}) = 1/c_{jt}$ .

For the demand moments, after invoking the implicit function theorem, this has a convenient block structure that can be separated market by market  $t$ :<sup>104</sup>

$$\underbrace{\frac{\partial \delta_t}{\partial \theta_2}}_{J_t \times K_2} = - \left( \underbrace{\begin{bmatrix} \frac{\partial \mathbf{s}_t}{\partial \delta_t}(\theta_2) \end{bmatrix}}_{J_t \times J_t} \right)^{-1} \underbrace{\begin{bmatrix} \frac{\partial \mathbf{s}_t}{\partial \theta_2}(\theta_2) \end{bmatrix}}_{J_t \times K_2}.$$

Differentiating the supply moments is challenging because demand derivatives  $\frac{\partial s_{jt}}{\partial p_{kt}}(\xi_t(\theta_2), \theta_2)$  in the matrix  $\Delta_t(\xi_t(\theta_2), \theta_2)$  depend on both  $\xi_t(\theta_2)$  and  $\theta_2$  directly. To avoid excessive tensor product notation, consider the derivative with respect to an element within  $\theta_2$  labeled  $\theta_\ell$ :<sup>105</sup>

$$\begin{aligned} \underbrace{\frac{\partial \eta_t}{\partial \theta_\ell}}_{J_t \times 1} &= -\Delta_t^{-1} \frac{\partial \Delta_t}{\partial \theta_\ell} \Delta_t^{-1} \mathbf{s}_t - \Delta_t^{-1} \frac{\partial \Delta_t}{\partial \xi_t} \frac{\partial \xi_t}{\partial \theta_\ell} \Delta_t^{-1} \mathbf{s}_t + \Delta_t^{-1} \underbrace{\frac{\partial \mathbf{s}_t}{\partial \theta_\ell}}_0 \\ &= - \underbrace{\Delta_t^{-1}}_{(J_t \times J_t)} \underbrace{\frac{\partial \Delta_t}{\partial \theta_\ell}}_{(J_t \times J_t)} \underbrace{\eta_t}_{(J_t \times 1)} - \underbrace{\Delta_t^{-1}}_{(J_t \times J_t)} \underbrace{\frac{\partial \Delta_t}{\partial \xi_t}}_{(J_t \times J_t \times J_t)} \underbrace{\frac{\partial \xi_t}{\partial \theta_\ell}}_{(J_t \times 1)} \underbrace{\eta_t}_{(J_t \times 1)}. \end{aligned}$$

This expression is complicated because the supply-side structural error term  $\omega_t$  and the markup  $\eta_t$  depend both directly on  $\theta_2$  and indirectly on  $\theta_2$  through  $\xi_t$ .

□ **Quantity dependent marginal costs.** The original Berry, Levinsohn, and Pakes (1995) paper incorporates quantity dependent marginal costs, which allow for increasing or decreasing returns to scale:

$$\log(p_{jt} - \eta_{jt}(\theta_2)) = \log c_{jt} = [x_{jt}, w_{jt}] \gamma + \gamma_q \log q_{jt} + \omega_{jt}.$$

The obvious problem is that  $\log q_{jt}$  is endogenous in that it depends on  $(\xi_t, \omega_t)$ . A second question is whether the firm takes into account the fact that selling an additional unit *changes the marginal cost*, which would imply an additional term in the first-order condition:

$$\begin{aligned} s_{jt}(\mathbf{p}_t) + \sum_{k \in J_{ft}} \frac{\partial s_{kt}}{\partial p_{jt}}(\mathbf{p}_t) \cdot (p_{kt} - c_{kt}) - s_{kt}(\mathbf{p}_t) \cdot \frac{\partial c_{kt}}{\partial q_{kt}} \cdot \frac{\partial s_{kt}}{\partial p_{jt}}(\mathbf{p}_t) \cdot M_t &= 0, \\ s_{jt}(\mathbf{p}_t) + \sum_{k \in J_{ft}} \frac{\partial s_{kt}}{\partial p_{jt}}(\mathbf{p}_t) \cdot \left( p_{kt} - c_{kt} - \underbrace{M_t \cdot s_{kt}(\mathbf{p}_t)}_{q_{kt}(\mathbf{p}_t)} \cdot \frac{\partial c_{kt}}{\partial q_{kt}} \right) &= 0, \\ s_{jt}(\mathbf{p}_t) + \sum_{k \in J_{ft}} \frac{\partial s_{kt}}{\partial p_{jt}}(\mathbf{p}_t) \cdot (p_{kt} - c_{kt} \cdot (1 + \gamma_q)) &= 0. \end{aligned}$$

<sup>103</sup> During optimization, this equality does not hold because we concentrate over the linear parameters. Abbreviating (A6) as  $\frac{\partial L}{\partial \theta_2} = \frac{\partial R}{\partial \theta_2}$  and using notation from Appendix A.1,  $\frac{\partial L}{\partial \theta_2} = (I - X(\tilde{X}'W\tilde{X})^{-1}\tilde{X}'WZ)' \frac{\partial R}{\partial \theta_2} \neq \frac{\partial R}{\partial \theta_2}$ . We thank Luis Armona and Daniel Stackman for pointing this out. Fortunately, one can use orthogonality between  $L$  and the projection matrix to show  $\nabla q \propto \frac{\partial L}{\partial \theta_2} ZWZ'L = \frac{\partial R}{\partial \theta_2} ZWZ'L$ , i.e. it is fine to use (A6) when computing the gradient.

<sup>104</sup> The matrix inverse of  $\frac{\partial s_t}{\partial \theta_t}(\theta_2)$  is guaranteed by the diagonal dominance of system of equations with respect to  $\delta_{jt}$ . As long as the outside good has a positive share, we have that for each  $j$ ,  $|\frac{\partial s_{jt}}{\partial \delta_{jt}}| > \sum_{k \neq j} |\frac{\partial s_{kt}}{\partial \delta_{jt}}|$ . In practice, as shares become small, there may still be numerical issues.

<sup>105</sup> In the markup  $\eta_{jt}$ , the  $s_t$  is data and thus does not depend on parameters.

If  $\log q_{jt}$  enters the (log) marginal cost function, the implied marginal costs are all increased proportionally by  $(1 + \gamma_q)$ . Otherwise the first order condition depends both on the functional form  $f_{MC}(\cdot)$  and how  $q_{jt}$  enters the equation:

$$s_{jt}(p_t) + \sum_{k \in J_{jt}} \frac{\partial s_{kt}}{\partial p_{jt}}(p_t) \cdot \left( p_{kt} - c_{kt} - q_{kt} \frac{\partial c_{kt}}{\partial q_{kt}} \right) = 0.$$

This does not appear to be how the existing literature such as Berry, Levinsohn, and Pakes (1995, 1999) address quantity dependent marginal costs. Instead, the first-order condition treats the marginal cost as if it were constant (i.e., assuming  $\frac{\partial c_{kt}}{\partial q_{kt}} = 0$  above), but then when the parameters  $\gamma$  are recovered allows for an extra term on  $q_{jt}$  or  $\log q_{jt}$ . Thus the firm treats marginal costs as if they were constant when setting prices, but the marginal costs are still quantity dependent.<sup>106</sup>

$$f_{MC}(p_{jt} - \eta_{jt}(\theta_2)) = [x_{jt}, w_{jt}] \gamma + \gamma_q \log q_{jt} + \omega_{jt}.$$

Also note that  $q_{jt}$  cannot be included in the set of instruments for supply,  $Z_t^S$ , and as an endogenous variable, it increases the number of required instruments. The good news is that the conventional quantity shifters (BLP instruments) should be relevant here.

Appendix B: Additional tables and figures

TABLE B1 Fixed Point Tricks

Problem	Median $s_{0t}$	Overflow Safe	Type	Initial $\delta_t$	Mean Milliseconds	Mean Evaluations	Percent Converged
Simple simulation	0.91	Yes	$\delta_t$	$\delta_t^0$	2.56	15.94	100.00%
Simple simulation	0.91	No	$\delta_t$	$\delta_t^0$	1.66	15.94	100.00%
Simple simulation	0.91	Yes	$\exp(\delta_t)$	$\delta_t^0$	2.57	15.94	100.00%
Simple simulation	0.91	No	$\exp(\delta_t)$	$\delta_t^0$	1.66	15.94	100.00%
Simple simulation	0.91	Yes	$\delta_t$	$\delta_t^{n-1}$	1.78	13.93	100.00%
Complex simulation	0.91	Yes	$\delta_t$	$\delta_t^0$	3.32	16.14	100.00%
Complex simulation	0.91	No	$\delta_t$	$\delta_t^0$	2.32	16.15	100.00%
Complex simulation	0.91	Yes	$\exp(\delta_t)$	$\delta_t^0$	3.33	16.16	100.00%
Complex simulation	0.91	No	$\exp(\delta_t)$	$\delta_t^0$	2.34	16.15	100.00%
Complex simulation	0.91	Yes	$\delta_t$	$\delta_t^{n-1}$	2.74	14.74	100.00%
RCNL simulation	0.92	Yes	$\delta_t$	$\delta_t^0$	8.45	33.49	100.00%
RCNL simulation	0.92	No	$\delta_t$	$\delta_t^0$	6.06	33.49	100.00%
RCNL simulation	0.92	Yes	$\exp(\delta_t)$	$\delta_t^0$	8.24	33.49	100.00%
RCNL simulation	0.92	No	$\exp(\delta_t)$	$\delta_t^0$	6.09	33.50	100.00%
RCNL simulation	0.92	Yes	$\delta_t$	$\delta_t^{n-1}$	6.37	27.32	100.00%
Nevo example	0.54	Yes	$\delta_t$	$\delta_t^0$	4.05	24.06	100.00%
Nevo example	0.54	No	$\delta_t$	$\delta_t^0$	2.67	24.06	100.00%
Nevo example	0.54	Yes	$\exp(\delta_t)$	$\delta_t^0$	4.07	24.06	100.00%
Nevo example	0.54	No	$\exp(\delta_t)$	$\delta_t^0$	2.69	24.06	100.00%
Nevo example	0.54	Yes	$\delta_t$	$\delta_t^{n-1}$	3.10	18.27	100.00%
BLP example	0.89	Yes	$\delta_t$	$\delta_t^0$	31.54	40.61	100.00%
BLP example	0.89	No	$\delta_t$	$\delta_t^0$	26.97	40.69	100.00%
BLP example	0.89	Yes	$\exp(\delta_t)$	$\delta_t^0$	31.64	40.61	100.00%
BLP example	0.89	No	$\exp(\delta_t)$	$\delta_t^0$	26.57	40.69	100.00%
BLP example	0.89	Yes	$\delta_t$	$\delta_t^{n-1}$	29.08	37.95	100.00%

Note: This table documents the impact of common tricks used in the literature on solving the nested fixed point. Reported values are medians across 100 different simulations and 10 identical runs of the two example problems. We report the number of milliseconds and contraction evaluations needed to solve the nested fixed point, averaged across all markets and objective evaluations of one GMM step. We also report convergence rates: the percent of times no numerical errors were encountered a limit of 1000 iterations was not reached. Overflow safe results are those that use the log-sum-exp (LSE) function. The two fixed point types are the standard linear contraction over  $\delta_{jt}$  and the exponentiated version over  $\exp(\delta_{jt})$ . An initial  $\delta_t^0$  means that the contraction always starts at the solution to the logit model, whereas  $\delta_t^{n-1}$  is the “hot-start” version where starting values are those that solved the fixed point for the previous guess of  $\theta_2$ . We use the SQUAREM algorithm with an absolute  $L^\infty$  norm tolerance of 1E-14. Simulations are configured as in Section 5. The example problems from Nevo (2000b) and Berry, Levinsohn, and Pakes (1995, 1999) are the replications described in Section 6.

<sup>106</sup> Allowing for extra term in the first-order condition also likely violates the conditions necessary for uniqueness of the pricing equilibrium, particularly if  $\gamma_q$  is negative and there are increasing returns to scale.

TABLE B2 Impact of Alternative Integration Methods on Parameter Estimates

Simulation	Supply	Integration	$I_t$	True Value				Median Bias				Median Absolute Error				
				Seconds	$\alpha$	$\sigma_x$	$\sigma_p$	$\rho$	$\alpha$	$\sigma_x$	$\sigma_p$	$\rho$	$\alpha$	$\sigma_x$	$\sigma_p$	$\rho$
Simple	No	Monte Carlo	100	1.0	-1	3			0.233	-0.691			0.298	0.691		
Simple	No	Monte Carlo	1000	3.1	-1	3			0.198	-0.132			0.251	0.191		
Simple	No	MLHS	1000	3.2	-1	3			0.188	-0.051			0.241	0.167		
Simple	No	Halton	1000	3.2	-1	3			0.186	-0.050			0.241	0.165		
Simple	No	Importance	1000	21.6	-1	3			0.181	0.018			0.242	0.169		
Simple	No	Product rule	9 <sup>1</sup>	0.8	-1	3			0.189	-0.039			0.245	0.169		
Simple	Yes	Monte Carlo	100	2.7	-1	3			0.113	-0.705			0.243	0.705		
Simple	Yes	Monte Carlo	1000	8.8	-1	3			0.021	-0.102			0.180	0.182		
Simple	Yes	MLHS	1000	8.8	-1	3			0.020	-0.015			0.172	0.162		
Simple	Yes	Halton	1000	9.2	-1	3			0.020	-0.015			0.170	0.162		
Simple	Yes	Importance	1000	27.2	-1	3			-0.001	0.065			0.176	0.181		
Simple	Yes	Product rule	9 <sup>1</sup>	2.2	-1	3			0.015	0.003			0.172	0.172		
Complex	No	Monte Carlo	100	1.9	-1	3	0.2		0.304	-0.776	-0.091		0.317	0.778	0.110	
Complex	No	Monte Carlo	1000	5.3	-1	3	0.2		0.193	-0.190	-0.012		0.253	0.223	0.098	
Complex	No	MLHS	1000	5.1	-1	3	0.2		0.191	-0.103	-0.018		0.254	0.182	0.102	
Complex	No	Halton	1000	5.7	-1	3	0.2		0.141	-0.120	0.033		0.241	0.190	0.121	
Complex	No	Importance	1000	28.3	-1	3	0.2		0.180	-0.017	-0.028		0.254	0.173	0.108	
Complex	No	Product rule	9 <sup>2</sup>	1.6	-1	3	0.2		0.172	-0.088	-0.011		0.250	0.177	0.169	
Complex	Yes	Monte Carlo	100	5.2	-1	3	0.2		0.114	-0.702	-0.137		0.250	0.713	0.141	
Complex	Yes	Monte Carlo	1000	15.0	-1	3	0.2		0.029	-0.126	-0.040		0.194	0.204	0.106	
Complex	Yes	MLHS	1000	15.2	-1	3	0.2		0.015	-0.051	-0.050		0.195	0.171	0.146	
Complex	Yes	Halton	1000	17.1	-1	3	0.2		-0.025	-0.053	0.024		0.194	0.162	0.127	
Complex	Yes	Importance	1000	38.6	-1	3	0.2		-0.050	0.071	-0.020		0.208	0.190	0.112	
Complex	Yes	Product rule	9 <sup>2</sup>	4.7	-1	3	0.2		-0.020	-0.029	0.004		0.195	0.171	0.169	

(Continued)

TABLE B2 (Continued)

Simulation	Supply	Integration	$I_t$	True Value				Median Bias				Median Absolute Error			
				Seconds	$\alpha$	$\sigma_x$	$\sigma_p$	$\rho$	$\alpha$	$\sigma_x$	$\sigma_p$	$\rho$	$\alpha$	$\sigma_x$	$\sigma_p$
RCNL	No	Monte Carlo	100	5.7	-1	3		0.5	0.236	-0.645		0.046	0.268	0.645	0.051
RCNL	No	Monte Carlo	1000	18.9	-1	3		0.5	0.182	-0.131		0.001	0.219	0.182	0.022
RCNL	No	MLHS	1000	19.0	-1	3		0.5	0.177	-0.026		-0.007	0.216	0.160	0.021
RCNL	No	Halton	1000	19.9	-1	3		0.5	0.174	-0.024		-0.008	0.214	0.155	0.021
RCNL	No	Importance	1000	40.2	-1	3		0.5	0.086	0.854		-0.110	0.241	0.854	0.110
RCNL	No	Product rule	9 <sup>1</sup>	4.3	-1	3		0.5	0.176	-0.017		-0.007	0.214	0.153	0.021
RCNL	Yes	Monte Carlo	100	12.5	-1	3		0.5	0.072	-0.573		0.046	0.124	0.573	0.048
RCNL	Yes	Monte Carlo	1000	45.8	-1	3		0.5	0.020	-0.096		0.007	0.112	0.156	0.019
RCNL	Yes	MLHS	1000	45.8	-1	3		0.5	0.008	-0.005		0.000	0.109	0.138	0.017
RCNL	Yes	Halton	1000	47.6	-1	3		0.5	0.010	-0.008		0.000	0.109	0.139	0.017
RCNL	Yes	Importance	1000	66.0	-1	3		0.5	-0.071	0.862		-0.099	0.130	0.863	0.099
RCNL	Yes	Product rule	9 <sup>1</sup>	9.3	-1	3		0.5	0.002	-0.001		0.000	0.109	0.139	0.018

Note: This table documents bias and variance of parameter estimates over 1000 simulated datasets for different numerical integration methods and numbers of integration nodes  $I_t$ . Importance sampling is based on Halton draws. For descriptions of the integration rules, please refer to Figure 1. For all problems, we use the “approximate” version of the feasible optimal instruments.

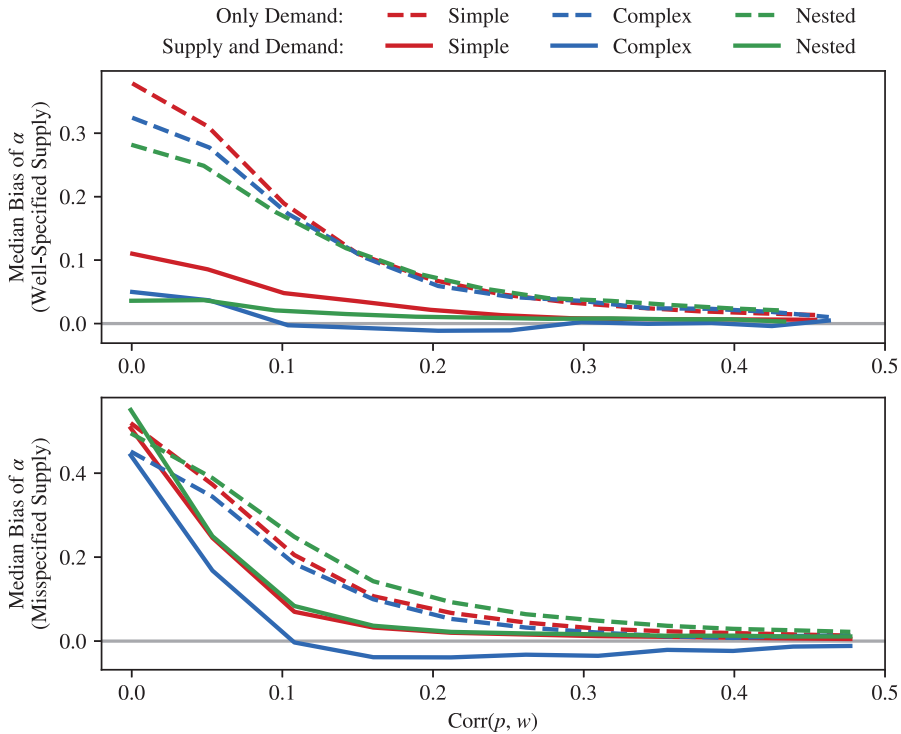
TABLE B3 Form of Feasible Optimal Instruments

Simulation	Supply	Approach	True Value				Median Bias				Median Absolute Error			
			Seconds	$\alpha$	$\sigma_x$	$\sigma_p$	$\rho$	$\alpha$	$\sigma_x$	$\sigma_p$	$\rho$	$\alpha$	$\sigma_x$	$\sigma_p$
Simple	No	Approximate	0.8	-1	3			0.189	-0.039			0.245	0.169	
Simple	No	Asymptotic	4.8	-1	3			0.189	-0.038			0.245	0.169	
Simple	No	Empirical	4.8	-1	3			0.188	-0.035			0.245	0.168	
Simple	Yes	Approximate	2.2	-1	3			0.015	0.003			0.172	0.172	
Simple	Yes	Asymptotic	18.1	-1	3			0.021	0.003			0.192	0.172	
Simple	Yes	Empirical	18.2	-1	3			0.026	0.007			0.181	0.171	
Complex	No	Approximate	1.6	-1	3	0.2		0.172	-0.088	-0.011		0.250	0.177	0.169
Complex	No	Asymptotic	6.5	-1	3	0.2		0.177	-0.084	-0.008		0.251	0.175	0.165
Complex	No	Empirical	6.5	-1	3	0.2		0.171	-0.085	-0.012		0.246	0.175	0.168
Complex	Yes	Approximate	4.7	-1	3	0.2		-0.020	-0.029	0.004		0.195	0.171	0.169
Complex	Yes	Asymptotic	29.7	-1	3	0.2		0.001	-0.085	-0.029		0.210	0.209	0.168
Complex	Yes	Empirical	29.5	-1	3	0.2		-0.008	-0.074	-0.016		0.211	0.199	0.168
RCNL	No	Approximate	4.3	-1	3		0.5	0.176	-0.017		-0.007	0.214	0.153	0.021
RCNL	No	Asymptotic	9.4	-1	3		0.5	0.175	-0.022		-0.007	0.214	0.153	0.021
RCNL	No	Empirical	9.6	-1	3		0.5	0.173	-0.020		-0.007	0.215	0.151	0.021
RCNL	Yes	Approximate	9.3	-1	3		0.5	0.002	-0.001		0.000	0.109	0.139	0.018
RCNL	Yes	Asymptotic	46.4	-1	3		0.5	0.010	0.005		-0.000	0.113	0.142	0.018
RCNL	Yes	Empirical	46.2	-1	3		0.5	0.011	-0.004		-0.000	0.113	0.145	0.017

Note: This table documents bias and variance of parameter estimates over 1000 simulated datasets for different forms of optimal instruments. The “approximate” approach replaces the structural errors with their unconditional expectations of zero. For the “asymptotic” approach we take 100 draws from the estimated normal distribution for  $(\xi_{jt}, \omega_{jt}) \sim N(0, \hat{\Omega})$ . For the “empirical” approach we sample with replacement from the joint distribution of  $(\xi_{jt}, \hat{\omega}_{jt})$ . For all problems, we use a Gauss-Hermite product rule that exactly integrates polynomials of degree 17 or less.

FIGURE B1

IMPACT OF SUPPLY-SIDE MOMENTS ON OPTIMAL IV FOR DEMAND-ONLY PROBLEMS [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]



Each plot documents how bias of the linear parameter on price,  $\alpha$ , decreases with the strength of the cost shifter  $w_j$ . Reported bias values are medians across 1000 different simulations. Unlike in Figure 2, once feasible optimal instruments have been constructed, only demand moments are used during estimation. Here, dashed lines mean that instruments were constructed only with demand moments, and solid lines mean that supply moments were used as well when constructing instruments. In the bottom plot we simulate data according to perfect competition (i.e., prices are set equal to marginal costs instead of those that satisfy Bertrand-Nash first order conditions), but continue to construct feasible optimal instruments under the assumption of imperfect competition. For all problems, we use the “approximate” version of the feasible optimal instruments and a Gauss-Hermite product rule that exactly integrates polynomials of degree 17 or less.

## References

- AI, C. and CHEN, X. “Efficient Estimation of Models with Conditional Moment Restrictions Containing Unknown Functions.” *Econometrica*, Vol. 71 (2003), pp. 1795–1843. <https://onlinelibrary.wiley.com/doi/abs/10.1111/1468-0262.00470>.
- ALLAIRE, J., USHEY, K., TANG, Y. and EDELBUEITTEL, D. Reticulate: R Interface to Python. <https://github.com/rstudio/reticulate>, Accessed March 01, 2020.
- AMEMIYA, T. “The Maximum Likelihood and the Nonlinear Three-Stage Least Squares Estimator in the General Nonlinear Simultaneous Equation Model.” *Econometrica*, Vol. 45 (1977), pp. 955–968. <http://www.jstor.org/stable/1912684>.
- ANDREWS, L., GENTZKOW, M. and SHAPIRO, J.M. “Measuring the Sensitivity of Parameter Estimates to Estimation Moments.” *Quarterly Journal of Economics*, Vol. 132 (2017), pp. 1553–1592.
- ARMSTRONG, T. “Large Market Asymptotics for Differentiated Product Demand Estimators with Economic Models of Supply.” *Econometrica*, Vol. 84 (2016), pp. 1961–1980.
- BACKUS, M., CONLON, C. and SINKINSON, M. “Common Ownership and Competition in the Ready-To-Eat Cereal Industry.” Working Paper, 2020.
- BAYER, P., FERREIRA, F. and McMILLAN, R. “A Unified Framework for Measuring Preferences for Schools and Neighborhoods.” *Journal of Political Economy*, Vol. 115 (2007), pp. 588–638. <https://doi.org/10.1086/522381>.



- BERRY, S. "Estimating Discrete-Choice Models of Product Differentiation." *RAND Journal of Economics*, Vol. 25 (1994), pp. 242–261.
- BERRY, S., LEVINSOHN, J. and PAKES, A. "Automobile Prices in Market Equilibrium." *Econometrica*, Vol. 63 (1995), pp. 841–890.
- . "Voluntary Export Restraints on Automobiles: Evaluating a Trade Policy." *American Economic Review*, Vol. 89 (1999), pp. 400–430. <http://www.aeaweb.org/articles?id=10.1257/aer.89.3.400>.
- . "Differentiated Products Demand Systems from a Combination of Micro and Macro Data: The New Car Market." *Journal of Political Economy*, Vol. 112 (2004a), pp. 68–105.
- BERRY, S., LINTON, O.B. and PAKES, A. "Limit Theorems for Estimating the Parameters of Differentiated Product Demand Systems." *Review of Economic Studies*, Vol. 71 (2004b), pp. 613–654. <http://ideas.repec.org/a/bla/restud/v71y2004ip613-654.html>.
- BERRY, S. and PAKES, A. "The Pure Characteristics Demand Model." *International Economic Review*, Vol. 48 (2007), pp. 1193–1225.
- BERRY, S.T. and HAILE, P.A. "Identification in Differentiated Products Markets Using Market Level Data." *Econometrica*, Vol. 82 (2014), pp. 1749–1797. <https://onlinelibrary.wiley.com/doi/abs/10.3982/ECTA9027>.
- BHAT, C.R. "Quasi-Random Maximum Simulated Likelihood Estimation of the Mixed Multinomial Logit Model." *Transportation Research Part B: Methodological*, Vol. 35 (2001), pp. 677–693. <http://www.sciencedirect.com/science/article/pii/S01912615000014X>.
- BONNET, C. and DUBOIS, P. "Inference on Vertical Contracts Between Manufacturers and Retailers Allowing for Non-linear Pricing and Resale Price Maintenance." *The RAND Journal of Economics*, Vol. 41 (2010), pp. 139–164. <http://www.jstor.org/stable/40649274>.
- BOYD, J. and MELLMAN, R.E. "The Effect of Fuel Economy Standards on the U.S. Automotive Market: An Hedonic Demand Analysis." *Transportation Research Part A: General*, Vol. 14 (1980), pp. 367–378. <http://www.sciencedirect.com/science/article/pii/0191260780900552>.
- BRENKERS, R. and VERBOVEN, F. "Liberalizing a Distribution System: The European Car Market." *Journal of the European Economic Association*, Vol. 4 (2006), pp. 216–251. <https://onlinelibrary.wiley.com/doi/abs/10.1162/jeea.2006.4.1.216>.
- BRESNAHAN, T.F. "The Oligopoly Solution Concept is Identified." *Economics Letters*, Vol. 10 (1982), pp. 87–92. <http://www.sciencedirect.com/science/article/pii/0165176582901215>.
- BRUNNER, D. Numerical Integration in Random Coefficient Models of Demand. Master's thesis, Heinrich-Heine-Universität Düsseldorf, 2017.
- BRUNNER, D., HEISS, F., ROMAHN, A. and WEISER, C. "Reliable Estimation of Random Coefficient Logit Demand Models." Technical report, University of Düsseldorf, 2017.
- CAPLIN, A. and NALEBUFF, B. "Aggregation and Imperfect Competition: On the Existence of Equilibrium." *Econometrica*, Vol. 59 (1991), pp. 25–59. <http://www.jstor.org/stable/2938239>.
- CARDELL, N. and DUNBAR, F.C. "Measuring the Societal Impacts of Automobile Downsizing." *Transportation Research Part A: General*, Vol. 14 (1980), pp. 423–434. <http://www.sciencedirect.com/science/article/pii/0191260780900606>.
- CHAMBERLAIN, G. "Asymptotic Efficiency in Estimation with Conditional Moment Restrictions." *Journal of Econometrics*, Vol. 34 (1987), pp. 305–334.
- CONLON, C. "The MPEC Approach to Empirical Likelihood Estimation of Demand." Working Paper, 2017.
- CONLON, C. and RAO, N. "The Price of Liquor is Too Damn High: The Effects of Post and Hold Pricing." Working Paper, 2017.
- CORREIA, S. "Linear Models with High-Dimensional Fixed Effects: An Efficient and Feasible Estimator." Working Paper, 2016.
- DONALD, S.G., IMBENS, G.W. and NEWEY, W.K. "Choosing Instrumental Variables in Conditional Moment Restriction Models." *Journal of Econometrics*, Vol. 152 (2009), pp. 28–36. <http://www.sciencedirect.com/science/article/pii/S0304407609000566>.
- DUBÉ, J.P.H., FOX, J.T. and SU, C.L. "Improving the Numerical Performance of BLP Static and Dynamic Discrete Choice Random Coefficients Demand Estimation." *Econometrica*, Vol. 80 (2012), pp. 2231–2267.
- FAN, Y. "Ownership Consolidation and Product Characteristics: A Study of the US Daily Newspaper Market." *American Economic Review*, Vol. 103 (2013), pp. 1598–1628. <http://www.aeaweb.org/articles?id=10.1257/aer.103.5.1598>.
- FONG, D.C.L. and SAUNDERS, M. "LSMR: An Iterative Algorithm for Sparse Least-Squares Problems." *SIAM Journal on Scientific Computing*, Vol. 33 (2011), pp. 2950–2971.
- FOX, J.T., IL KIM, K., RYAN, S.P. and BAJARI, P. "The Random Coefficients Logit Model is Identified." *Journal of Econometrics*, Vol. 166 (2012), pp. 204–212. <http://www.sciencedirect.com/science/article/pii/S0304407611001655>.
- FREYBERGER, J. "Asymptotic Theory for Differentiated Products Demand Models with Many Markets." *Journal of Econometrics*, Vol. 185 (2015), pp. 162–181. <http://www.sciencedirect.com/science/article/pii/S0304407614002474>.
- GALLEGO, G., HUH, W.T., KANG, W. and PHILLIPS, R. "Price Competition with the Attraction Demand Model: Existence of Unique Equilibrium and Its Stability." *Manufacturing & Service Operations Management*, Vol. 8 (2006), pp. 359–375. <https://pubsonline.informs.org/doi/abs/10.1287/msom.1060.0115>.

- GANDHI, A. and HOUDE, J. "Measuring Substitution Patterns in Differentiated Products Industries." Working Paper, 2019.
- GORTMAKER, J. and TARASCINA, A. *PyHDFE*. <https://github.com/jeffgortmaker/pyhdf>. Accessed on March 1 2020.
- GRIGOLON, L. and VERBOVEN, F. "Nested Logit or Random Coefficients Logit? A Comparison of Alternative Discrete Choice Models of Product Differentiation." *The Review of Economics and Statistics*, Vol. 96 (2014), pp. 916–935. [https://doi.org/10.1162/REST\\_a\\_00420](https://doi.org/10.1162/REST_a_00420).
- GUIMARÃES, P. and PORTUGAL, P. "A Simple Feasible Procedure to Fit Models with High-Dimensional Fixed Effects." *Stata Journal*, Vol. 10 (2010), pp. 628–649. <https://ideas.repec.org/a/tsj/stataj/v10y2010i4p628-649.html>.
- HALTON, J.H. "On the Efficiency of Certain Quasi-Random Sequences of Points in Evaluating Multi-Dimensional Integrals." *Numerische Mathematik*, Vol. 2 (1960), pp. 84–90. <https://doi.org/10.1007/BF01386213>.
- HANSEN, L.P. "Large Sample Properties of Generalized Method of Moments Estimators." *Econometrica*, Vol. 50 (1982), pp. 1029–1054.
- HEISS, F. "The Panel Probit Model: Adaptive Integration on Sparse Grids." *Advances in Econometrics*, Vol. 26 (2010), pp. 41–64.
- HEISS, F. and WINSCHER, V. "Likelihood Approximation by Numerical Integration on Sparse Grids." *Journal of Econometrics*, Vol. 144 (2008), pp. 62–80. <http://www.sciencedirect.com/science/article/B6VC0-4RJYVBD-1/2/12335dfcbca363c96edf0931cc1b02f>.
- HESS, S., TRAIN, K.E. and POLAK, J.W. "On the Use of a Modified Latin Hypercube Sampling (MLHS) Method in the Estimation of a Mixed Logit Model for Vehicle Choice." *Transportation Research Part B: Methodological*, Vol. 40 (2006), pp. 147–163. <https://ideas.repec.org/a/eee/transb/v40y2006i2p147-163.html>.
- HO, K. and PAKES, A. "Hospital Choices, Hospital Prices, and Financial Incentives to Physicians." *American Economic Review*, Vol. 104 (2014), pp. 3841–84. <http://www.aeaweb.org/articles?id=10.1257/aer.104.12.3841>.
- HONG, H., LI, H. and LI, J. "BLP Estimation Using Laplace Transformation and Overlapping Simulation Draws." *Journal of Econometrics*, (Forthcoming).
- HOUDE, J.F. "Spatial Differentiation and Vertical Mergers in Retail Markets for Gasoline." *American Economic Review*, Vol. 102 (2012), pp. 2147–82. <http://www.aeaweb.org/articles?id=10.1257/aer.102.5.2147>.
- IARIA, A. and WANG, A. "Identification and Estimation of Demand for Bundles." Working Paper, 2019.
- JOHNSON, S.G. *PyCall* (2019). <https://github.com/JuliaPy/PyCall.jl>.
- JUDD, K. *Numerical Methods in Economics*, Vol. 1, 1 ed. Cambridge, Mass.: The MIT Press, 1998. <https://EconPapers.repec.org/RePEc:mtp:titles:0262100711>.
- JUDD, K.L. and SKRAINKA, B. "High Performance Quadrature Rules: How Numerical Integration Affects a Popular Model of Product Differentiation." CeMMAP Working Papers CWP03/11, Centre for Microdata Methods and Practice, Institute for Fiscal Studies, 2011.
- KNITTEL, C.R. and METAXOGLU, K. "Estimation of Random-Coefficient Demand Models: Two Empiricists' Perspective." *Review of Economics and Statistics*, Vol. 96 (2014), pp. 34–59.
- KONVALOV, A. and SANDOR, Z. "On Price Equilibrium with Multi-Product Firms." *Economic Theory*, Vol. 44 (2010), pp. 271–292. <http://www.jstor.org/stable/40864781>.
- LEE, J. and SEO, K. "A Computationally Fast Estimator for Random Coefficients Logit Demand Models Using Aggregate Data." *The RAND Journal of Economics*, Vol. 46 (2015), pp. 86–102. <http://www.jstor.org/stable/43895583>.
- . "Revisiting the nested fixed-point algorithm in BLP random coefficients demand estimation." *Economics Letters*, Vol. 149 (2016), pp. 67–70.
- LEE, R.S. "Vertical Integration and Exclusivity in Platform and Two-Sided Markets." *American Economic Review*, Vol. 103 (2013), pp. 2960–3000. <http://www.aeaweb.org/articles?id=10.1257/aer.103.7.2960>.
- MCFADDEN, D. and TRAIN, K. "Mixed MNL Models for Discrete Response." *Journal of Applied Econometrics*, Vol. 15 (2000), pp. 447–470.
- MILLER, N. and WEINBERG, M. "Understanding the Price Effects of the MillerCoors Joint Venture." *Econometrica*, Vol. 85 (2017), pp. 1763–1791.
- MIRAVETE, E.J., SEIM, K. and THURK, J. "Market Power and the Laffer Curve." *Econometrica*, Vol. 86 (2018), pp. 1651–1687. <https://onlinelibrary.wiley.com/doi/abs/10.3982/ECTA12307>.
- MORE, J.J., GARBOW, B.S. and HILLSTROM, K.E. "User Guide for MINPACK-1." Technical report, Argonne National Laboratory, 1980.
- MORROW, W.R. and SKERLOS, S.J. "On the Existence of Bertrand-Nash Equilibrium Prices Under Logit Demand." *CoRR*, vol. abs/1012.5832 (2010).
- . "Fixed-Point Approaches to Computing Bertrand-Nash Equilibrium Prices Under Mixed-Logit Demand." *Operations Research*, Vol. 59 (2011), pp. 328–345. <http://orjournal.informs.org/content/59/2/328.abstract>.
- NEVO, A. "Mergers with Differentiated Products: The Case of the Ready-to-Eat Cereal Industry." *RAND Journal of Economics*, Vol. 31 (2000a), pp. 395–421.
- . "A Practitioner's Guide to Estimation of Random Coefficients Logit Models of Demand (Including Appendix)." *Journal of Economics and Management Strategy*, Vol. 9 (2000b), pp. 513–548.
- . "Measuring Market Power in the Ready-to-Eat Cereal Industry." *Econometrica*, Vol. 69 (2001), pp. 307–342.
- NEWKEY, W. "Generalized Method of Moments Specification Testing." *Journal of Econometrics*, Vol. 29 (1985), pp. 229–256. <https://EconPapers.repec.org/RePEc:eee:econom:v:29:y:1985:i:3:p:229-256>.

- NEWKEY, W.K. "Semiparametric Efficiency Bounds." *Journal of Applied Econometrics*, Vol. 5 (1990), pp. 99–135. <https://onlinelibrary.wiley.com/doi/abs/10.1002/jae.3950050202>.
- NEWKEY, W.K. and SMITH, R.J. "Higher Order Properties of GMM and Generalized Empirical Likelihood Estimators." *Econometrica*, Vol. 72 (2004), pp. 219–255.
- NIELSON, C. "Targeted Vouchers, Competition Among Schools, and the Academic Achievement of Poor Students." Working Paper, 2017.
- OWEN, A. "Scrambled Net Variance for Integrals of Smooth Functions." *Annals of Statistics*, Vol. 25 (1997), pp. 1541–1562.
- OWEN, A.B. "Multidimensional Variation for Quasi-Monte Carlo." *World Scientific* (2005), pp. 49–74. [https://www.worldscientific.com/doi/abs/10.1142/9789812567765\\_0004](https://www.worldscientific.com/doi/abs/10.1142/9789812567765_0004).
- . "A Randomized Halton Algorithm in R." 2017.
- PETRIN, A. "Quantifying the Benefits of New Products: The Case of the Minivan." *Journal of Political Economy*, Vol. 110 (2002), pp. 705–729.
- REYNAERT, M. and VERBOVEN, F. "Improving the Performance of Random Coefficients Demand Models: The Role of Optimal Instruments." *Journal of Econometrics*, Vol. 179 (2014), pp. 83–98. <https://EconPapers.repec.org/RePEc:eee:econom:v:179:y:2014:i:1:p:83-98>.
- REYNAERTS, J., VARADHAN, R. and NASH, J.C. "Enhancing the Convergence Properties of the BLP (1995) Contraction Mapping." Vives Discussion Paper Series 35, Katholieke Universiteit Leuven, Faculteit Economie en Bedrijfswetenschappen, Vives, 2012. <http://ideas.repec.org/p/ete/vivwps/35.html>.
- SALANIE, B. and WOLAK, F. "Fast, Robust, and Approximately Correct: Estimating Mixed Demand Systems." Working Paper, 2019.
- SKRAINKA, B. *Three Essays on Product Differentiation*. PhD dissertation, University College London, 2012a.
- SKRAINKA, B.S. "A Large Scale Study of the Small Sample Performance of Random Coefficient Models of Demand." Working Paper, 2012b.
- SOMAINI, P. and WOLAK, F. "An Algorithm to Estimate the Two-Way Fixed Effects Model." *Journal of Econometric Methods*, Vol. 5 (2016), pp. 143–152.
- STOCK, J.H. and WRIGHT, J.H. "GMM with Weak Identification." *Econometrica*, Vol. 68 (2000), pp. 1055–1096.
- SU, C.L. and JUDD, K.L. "Constrained Optimization Approaches to Estimation of Structural Models." *Econometrica*, Vol. 80 (2012), pp. 2213–2230.
- TRAIN, K. "Halton Sequences for Mixed Logit." Technical report, University of California, Berkeley, 2000.
- . *Discrete Choice Methods with Simulation*. Cambridge, UK: Cambridge University Press, 2009. <https://EconPapers.repec.org/RePEc:cup:cbooks:9780521766555>.
- VARADHAN, R. and ROLAND, C. "Simple and Globally Convergent Methods for Accelerating the Convergence of Any EM Algorithm." *Scandinavian Journal of Statistics*, Vol. 35 (2008), pp. 335–353. <http://www.jstor.org/stable/41548597>.
- VILLAS-BOAS, S.B. "Vertical Relationships between Manufacturers and Retailers: Inference with Limited Data." *The Review of Economic Studies*, Vol. 74 (2007), pp. 625–652. <http://www.jstor.org/stable/4626153>.

## Supporting information

Additional supporting information may be found online in the Supporting Information section at the end of the article.

**Figure OA1:** Histograms for Median Product Own-Price Elasticities

**Figure OA2:** Histograms for Top Product Own-Price Elasticities

**Figure OA3:** Histograms for Median Product Markups

**Figure OA4:** Histograms for Top Product Markups

**Figure OA5:** Histograms for GMM Objective Values

**Figure OA6:** Integration Error: Importance Sampling

**Figure OA7:** Integration Error: Importance Sampling and Relative RMSE

**Figure OA8:** Instrument Strength and Misspecification: Variance

**Figure OA9:** Impact of Supply-Side Moments on Optimal IV for Demand-Only Problem: Variance

**Figure OA10:** Instrument Strength and Misspecification: Non-Optimal Instruments, Bias

**Figure OA11:** Instrument Strength and Misspecification: Non-Optimal Instruments, Variance

**Figure OA12:** Profiled GMM Objective with Alternative Instruments: Complex Simulation

**Figure OA13:** Profiled GMM Objective with Alternative Instruments: RCNL Simulation

**Figure OA14:** Problem Scaling:  $\sigma_p$  and  $\phi$

**Figure OA15:** Scaling the Number of Products per Firm:  $\alpha$  and  $\sigma_x$

**Figure OA16:** Scaling the Number of Products per Firm:  $\sigma_p$  and  $\rho$

**Figure OA17:** Scaling the Number of Firms per Market:  $\alpha$  and  $\sigma_x$

**Figure OA18:** Scaling the Number of Firms per Market:  $\sigma_p$  and  $\rho$

**Table OA1:** Optimization Algorithms: Knittel and Metaxoglou (2014) Replication

**Table OA2:** Instrument Strength: Well-Specified Supply

**Table OA3:** Instrument Strength: Misspecified Supply

**Table OA4:** Instrument Strength: Well-Specified Supply, Non-Optimal Instruments

**Table OA5:** Instrument Strength: Misspecified Supply, Non-Optimal Instruments

**Table OA6:** Problem Scaling: Summary

**Table OA7:** Standard Errors: Alternative Instruments

**Table OA8:** Standard Errors: Alternative Integration Methods

**Table OA9:** Standard Errors: Problem Scaling

**Table OA10:** Post-Estimation Outputs: Alternative Instruments

**Table OA11:** Post-Estimation Outputs: Alternative Integration Methods

**Table OA12:** Post-Estimation Outputs: Problem Scaling

**Table OA13:** Merger Simulation: Alternative Instruments

**Table OA14:** Merger Simulation: Alternative Integration Methods

**Table OA15:** Merger Simulation: Problem Scaling

**Table OA16:** Optimization Algorithms: Additional Routines

**Table OA17:** Optimization Algorithms: Sums of Characteristics BLP Instruments

**Table OA18:** Optimization Algorithms: Parameter Estimates

Data S1